

```
1 import numpy as np
2
3 # basic concept
4 #One of the key features of NumPy is its N-dimensional array object, or ndarray,
5 # which is a fast, flexible container for large data sets in Python.
6
7 arr = np.array([1, 2, 3, 4, 5])
8
9 print(arr.dtype)
10 print(arr)
11
12 float_arr = arr.astype(np.float64)
13
14 print(float_arr.dtype)
15 print(float_arr)
16
17 # In addition to np.array, there are a number of other functions for creating new
arrays.
18 # As examples, zeros and ones create arrays of 0's or 1's, respectively, with a
given length or shape.
19 # empty creates an array without initializing its values to any particular value.
20
21 arr = np.ones(10)
22 print(arr)
23
24 arr = np.zeros(10)
25 print(arr)
26
27 arr = np.zeros((3, 6))
28 print(arr)
29
30 arr = np.empty((2, 4, 2))
31 print(arr)
32
33
34 int_array = np.arange(10)
35 print(int_array)
36
37 # array vectorization
38 # Arrays are important because they enable you to express batch operations on data
without writing any for loops.
39 # This is usually called vectorization.
40 # Any arithmetic operations between equal-size arrays applies the operation
elementwise:
41
42 arr = np.array([[1, 2, 3], [4, 5, 6]])
43
44 print(arr)
45 print(arr * arr)
```

```

46 print(arr - arr)
47
48 print(1 / arr)
49 print("arr*0.5=", arr * 0.5)
50
51 # it is different with broadcasting which applied to the operations between different
    arrays.
52 # Arithmetic operations with scalars are as you would expect, propagating the value
    to each element
53
54 print(1 / arr)
55 print("arr*0.5=", arr * 0.5)
56
57
58 # Numpy Array operation
59
60 # NumPy array indexing is a rich topic, as there are many ways you may want to select
    a subset of your data
61 # or individual elements.
62 # One-dimensional arrays are simple; on the surface they act similarly to Python
    lists:
63
64 arr = np.arange(10)
65 print(arr[5])
66 print(arr[5:8])
67
68 # As you can see, if you assign a scalar value to a slice, as in arr[5:8] = 12,
69 # the value is propagated (or broadcasted henceforth) to the entire selection.
70 # any modifications will be reflected in the source array:
71
72 alist = list(range(10))
73 alist[5] = 10
74 # alist[5:8] = 10 #error
75 print(alist)
76
77 arr[5:8] = 10
78 print(arr)
79
80 new_arr = arr[5:8].copy()
81 print(new_arr)
82
83
84
85
86
87 # slice
88 arr2d = np.array([[1, 2, 3],
89 ..... [4, 5, 6],
90 ..... [7, 8, 9]])

```

```

91
92 print(arr2d)
93 print(arr2d[:2])
94 print(arr2d[:2, 1:])
95 print(arr2d[1, :2])
96 print(arr2d[2, 0])
97 print(arr2d[2, :1])
98
99
100
101 arr2d[:2, 1:] = 0
102 print(arr2d)
103
104
105 # Note that a colon by itself means to take the entire axis,
106 # so you can slice only higher dimensional axes by doing:
107
108 print(arr2d[:, :1])
109
110
111 # boolean Index
112
113 # Let's consider an example where we have some data in an array and an array of
names with duplicates.
114 # I'm going to use here the randn function in numpy.random to generate some random
normally distributed data:
115
116 names = np.array(['Bob', 'Joe', 'Will', 'Bob', 'Will', 'Joe', 'Joe'])
117 data = np.random.randn(7, 4)
118
119 print(names)
120 print(data)
121
122 #names == "Bob"
123 # If we wanted to select all the rows with corresponding name 'Bob'.
124 # Like arithmetic operations, comparisons (such as ==) with arrays are also
vectorized.
125 print("Bob data =", data[names == 'Bob'])
126
127 print("!Bob data =", data[names != 'Bob'])
128
129 print("Partial Bob data =", data[names == 'Bob', 2:])
130
131 print("data < 0", data[data < 0])
132
133
134 # data[data < 0]=0
135 # print("data < 0 is 0", data)
136

```

```

137 data [names != 'Bob'] = 0
138 print(data)
139
140 # fancy indexing
141 # Fancy indexing is a term adopted by NumPy to describe indexing using integer
    arrays.
142
143 arr = np.empty((8,4))
144 for i in range(8):
145     arr[i] = i
146
147 print(arr)
148
149 # To select out a subset of the rows in a particular order,
150 # you can simply pass a list or ndarray of integers specifying the desired order:
151
152 print(arr[[4, 3, 0, 6]])
153
154 print(arr[[2, 3]])
155
156 # Using negative indices select rows from the end:
157 print(arr[[-2, -1]])
158
159
160
161 arr = np.arange(32)
162 print(arr)
163
164 arr = arr.reshape(8, 4)
165 print(arr)
166
167 #arr = np.arange(32).reshape(8, 4)
168
169 # Passing multiple index arrays does something slightly different;
170 # it selects a 1D array of elements corresponding to each tuple of indices:
171 print("intersection elements= ", arr[[1, 5, 7, 2], [0, 3, 1, 2]])
172
173 # Take a moment to understand what just happened: the elements (1, 0), (5, 3), (7, 1
    ), and (2, 2) were selected.
174 # The behavior of fancy indexing in this case is a bit different from what some
    users might have expected
175 # (myself included), which is the rectangular region formed by selecting a subset of
    the matrix's rows and columns
176
177 print(arr[np.ix_([1, 5, 7, 2], [0, 3, 1, 2])]). # generate new array, different with slice
178 print("arr=", arr)
179
180 #transpose
181

```

```

182 # Transposing is a special form of reshaping which similarly returns a view on the
    un- derlying data
183 # without copying anything.
184
185 arr = np. arange(15). reshape((3, 5))
186
187 print("arr = ", arr)
188
189 print(arr.T)
190
191 # When doing matrix computations, you will do this very often,
192 # like for example computing the inner matrix product XTX using np. dot:
193
194 arrDot = np. dot(arr.T, arr)
195 print("arrDot", arrDot)
196
197 # universal function
198 # A universal function, or ufunc, is a function that performs elementwise operations
    on data in ndarrays.
199 # You can think of them as fast vectorized wrappers for simple functions that take
    one or more scalar values
200 # and produce one or more scalar results.
201
202 arr = np. arange(8)+1
203 print(arr)
204
205 print(np. square(arr))
206 print(np. exp(arr))
207 print(np. log(arr))
208
209
210 # case study 1 : mesh computing
211
212 # import matplotlib.pyplot as plt
213 # import numpy as np
214 #
215 # points = np. arange(-5, 5, 0.1) # 1000 points with the same distance
216 # print(points)
217 # xs, ys = np. meshgrid(points, points)
218 #
219 # print("ys=", ys)
220 # print("xs=", xs)
221 #
222 # z = np. sqrt(xs**2+ys**2)
223 # print(z)
224 #
225 # plt. imshow(z)
226 # plt. colorbar()
227 # plt. title('Image plot of  $\sqrt{X^2+Y^2}$  for a grid of values')

```

```
228 #
229
230
231 # logic expression
232 #x if condition else y
233
234
235 arr = np.random.randn(4, 4)
236 print(arr)
237
238 # The numpy.where function is a vectorized version of the ternary expression x if
condi tion else y.
239 # print(np.where(arr > 0.5, 2, -2))
240 print(np.where(arr > 1, 0, arr))
241
242
243 # result = [(2 if c > 0.5 else -2) for c in arr]
244 # print(result)
245
246 alist = list(range(10))
247 print(alist)
248
249 result = [(0 if c > 5 else 1) for c in alist]
250 print(result)
251
252 # A set of mathematical functions which compute statistics about an entire array or
about the data
253 # along an axis are accessible as array methods.
254 # Aggregations (often called reductions) like sum, mean, and standard deviation std
can either be used by
255 # calling the array instance method or using the top level NumPy function
256
257 # sum, mean, std, var, min, max , argmin, argmax, cumsum, cumprod
258 arr = np.arange(9).reshape(3, 3)
259
260 print(arr)
261
262 print(arr.sum())
263 print(np.sum(arr))
264
265 print(arr.mean())
266 print(arr.std())
267 print(arr.var())
268 print(arr.min())
269 print(arr.max())
270 print(arr.argmin())
271 print(arr.argmax())
272
273 # Cumulative sum of elements starting from 0
```

```

274 a = np.array([[1, 2, 3], [4, 5, 6]])
275 print(np.cumsum(a))
276
277 col = np.cumsum(a, axis=0) . . . . . # sum over rows for each of the 3 columns
278 print(col)
279
280 row = np.cumsum(a, axis=1) . . # sum over columns for each of the 2 rows
281 print(row)
282
283 arr = np.random.randn(8)
284 # result = arr.sort()
285 # print(result)
286 print(np.sort(arr))
287
288 #unique
289 # NumPy has some basic set operations for one-dimensional ndarrays.
290 # Probably the most commonly used one is np.unique, which returns the sorted unique
    values in an array:
291
292 names = np.array(['Bob', 'Joe', 'Will', 'Bob', 'Will', 'Joe', 'Joe'])
293 print(np.unique(names))
294 print(sorted(set(names)))
295
296 #####
297 #
298 #
299 #
300 #####
301
302 from pandas import Series, DataFrame
303 import pandas as pd
304
305 # Thus, whenever you see pd. in code, it's referring to pandas.
306 # Series and DataFrame are used so much that I find it easier to import them into
    the local namespace.
307
308
309 # A Series is a one-dimensional array-like object containing an array of data (of
    any NumPy data type)
310 # and an associated array of data labels, called its index.
311 # The simplest Series is formed from only an array of data:
312
313 obj = Series([4, 7, -5, 3])
314 print(obj)
315 print(obj.values)
316 print(obj.index)
317
318 obj2 = Series([4, 7, 5, 3], index=['d', 'b', 'a', 'c'])
319 print(obj2)

```

```

320 print(obj2['a'])
321 print(obj2[['a', 'c']]) # double [[]]
322 print(obj2[obj2 > 0])
323 print(obj2*2)
324 print(np.exp(obj2))
325
326 # Should you have data contained in a Python dict, you can create a Series from it
by passing the dict
327 sdata = {'Ohio': 35000, 'Texas': 71000, 'Oregon': 16000, 'Utah': 5000}
328 obj3 = Series(sdata)
329 print(obj3)
330
331 # extract partial data
332 states = ['California', 'Ohio', 'Oregon', 'Texas']
333 obj4 = Series(sdata, index=states)
334 print(obj4)
335
336 # check non value
337
338 print(pd.isnull(obj4))
339 print(pd.notnull(obj4))
340 print(obj4.isnull())
341
342 # Both the Series object itself and its index have a name attribute
343 obj4.name = 'population'
344 obj4.index.name = 'state'
345 print(obj4)
346
347 # Series' s index can be altered in place by assignment:
348 print(obj)
349 obj.index = ['Bob', 'Steve', 'Jeff', 'Ryan']
350 print(obj)
351
352 #####
353 #
354 # Pandas - DataFrame
355 #
356 #####
357
358 data = {'state': ['Ohio', 'Ohio', 'Ohio', 'Nevada', 'Nevada'], 'year': [2000,
    2001, 2002, 2001, 2002],
359 'pop': [1.5, 1.7, 3.6, 2.4, 2.9]}
360
361 frame = DataFrame(data)
362
363 print(frame)
364
365
366 print(DataFrame(data, columns=['year', 'state', 'pop']))

```



```

367 print(DataFrame(data, columns=['year', 'state', 'pop', 'debt']))
368
369 # A column in a DataFrame can be retrieved as a Series either by dict-like notation
    or by attribute:
370 print(frame.columns)
371 print(frame['state'])
372 print(frame.state)
373
374 print(frame)
375
376 #Rows can also be retrieved by position or name by a couple of methods, such as the
    ix indexing field
377 print(frame.ix[3])
378
379 frame['debt'] = 16.5
380 print(frame)
381
382 # For example, the empty 'debt' column could be assigned a scalar value or an array
    of values
383 frame['debt'] = np.arange(5.)
384 print(frame)
385
386 # When assigning lists or arrays to a column, the value's length must match the
    length of the DataFrame.
387 # If you assign a Series, it will be instead conformed exactly to the DataFrame's
    index, inserting missing values in any holes:
388
389 val = Series([-1.2, -1.5, -1.7], index=[2, 4, 5])
390 frame['debt'] = val
391 print(frame)
392
393 #Assigning a column that doesn't exist will create a new column.
394
395 frame['eastern'] = 1
396 print(frame)
397
398
399 frame['marks'] = frame.state == 'Ohio'
400 del frame['eastern']
401 print(frame)
402
403 # Index Objects
404 obj = Series(range(3), index=['a', 'b', 'c'])
405 print(obj)
406
407 # Index objects are immutable index[1] = 'd'
408
409 # Reindexing
410 # Calling reindex on this Series rearranges the data according to the new index,

```

```

411 # introducing missing values if any index values were not already present:
412
413 obj2 = obj.reindex(['a', 'b', 'c', 'd', 'e'])
414 print(obj2)
415
416 obj2 = obj.reindex(['a', 'b', 'c', 'd', 'e'], fill_value=0)
417 print(obj2)
418
419 # For ordered data like time series, it may be desirable to do some interpolation or
    filling of values when reindexing.
420 # The method option allows us to do this, using a method such as ffill which forward
    fills the values:
421
422 obj3 = Series(['blue', 'purple', 'yellow'], index=[0, 2, 4])
423 print(obj3)
424 obj3 = obj3.reindex(range(6), method='ffill')
425 print(obj3)
426
427 # ffill or pad : Fill (or carry) values forward, bfill or backfill : Fill (or carry
    ) values backward
428
429 # With DataFrame, reindex can alter either the (row) index, columns, or both.
430
431 frame = DataFrame(np.arange(9).reshape((3, 3)), index=['a', 'c', 'd'], columns=['
    Ohio', 'Texas', 'California'])
432 print(frame)
433
434
435 # When passed just a sequence, the rows are reindexed in the result:
436 frame2 = frame.reindex(['a', 'b', 'c', 'd'])
437 print(frame2)
438
439 # The columns can be reindexed using the columns keyword:
440 states = ['Texas', 'Utah', 'California']
441 frame = frame.reindex(columns=states)
442
443 print(frame)
444
445 # Both can be reindexed in one shot, though interpolation will only apply row-wise(
    axis 0)
446 frame = frame.reindex(index=['a', 'b', 'c', 'd'], method='ffill', columns=states)
447 print(frame)
448
449
450 # Dropping entries from an axis
451
452 obj = Series(np.arange(5.), index=['a', 'b', 'c', 'd', 'e'])
453 new_obj = obj.drop('c')
454 print(new_obj)

```

```

455
456 # With DataFrame, index values can be deleted from either axis:
457
458 data = DataFrame(np.arange(16).reshape((4, 4)), index=['Ohio', 'Colorado', 'Utah',
    'New York'], columns=['one', 'two', 'three', 'four'])
459
460 data.drop(['Colorado', 'Ohio'])
461
462 print(data)
463 data.drop('two', axis=1)
464 print(data)
465 # Summarizing and Computing Descriptive Statistics
466 print(data.describe())
467
468 print(data.sum())
469 print(data.sum(axis=1))
470
471 data.ix["ohio"] = None
472 print(data)
473 data1 = data.mean(axis=0, skipna=True)
474 print(data1)
475
476 #like idxmin and idxmax, return indirect statistics like the index value where the
    minimum or maximum values are attained:
477 print(data.idxmax())
478
479 #import pandas.io.data as web
480
481 # from pandas import Series, DataFrame
482 # import pandas as pd
483
484 import pandas_datareader.data as web
485 import datetime
486
487 # start = datetime.datetime(2013, 1, 1)
488 # end = datetime.datetime(2017, 3, 7)
489 # df = web.DataReader("GOOGL", 'yahoo', start, end)
490 #
491 # print("df=", df)
492 #
493 # price = df['Adj Close']
494 #
495 # # price = DataFrame({tic: df['Adj Close'] for tic, data in df.iteritems()})
496 #
497 # print("price =", price)
498 #
499 # returns = price.pct_change()
500 # print(returns.tail())
501

```

```

502 start = datetime.datetime(2017, 1, 1)
503 end = datetime.datetime(2017, 3, 7)
504 all_data = {}
505
506 for ticker in ['AAPL', 'IBM', 'GOOG']:
507     all_data[ticker] = web.DataReader(ticker, 'yahoo', start, end)
508
509
510
511 # draw graph
512
513 import matplotlib as mpl
514 import matplotlib.pyplot as plt
515 #matplotlib inline
516
517 np.random.seed(1000)
518 y = np.random.standard_normal(20)
519 x = range(len(y))
520
521 plt.plot(x, y)
522
523 price = DataFrame({tic: data['Adj Close'] for tic, data in all_data.items()})
524 volume = DataFrame({tic: data['Volume'] for tic, data in all_data.items()})
525
526 print("price = \n", price.tail(5))
527 print("price AAPL= \n", price['AAPL'])
528 print("volume = \n", volume.tail(5))
529
530 #plt.plot(price['AAPL'])
531
532 returns = price.pct_change()
533 print(returns.tail(10))
534
535 # The corr method of Series computes the correlation of the overlapping, non-NA,
536 # aligned-by-index values in two Series. Relatedly, cov computes the covariance:
537 cov = returns.AAPL.corr(returns.IBM)
538 print(cov)
539
540 # DataFrame's corr and cov methods, on the other hand,
541 # return a full correlation or covariance matrix as a DataFrame, respectively:
542
543 print(returns.corr())
544 print(returns.cov())
545
546 # Using DataFrame's corrwith method, you can compute pairwise correlations between
a DataFrame's columns or rows with another Series or DataFrame.
547 # Passing a Series returns a Series with the correlation value computed for each
column:
548

```

```
549 print("IBM\n", returns.corrwith(returns.IBM))
550
551 #Passing a DataFrame computes the correlations of matching column names.
552 # Here I compute correlations of percent changes with volume:
553 print("returns with volume\n", returns.corrwith(volume))
554
555 #dates = []
556
557 # for x in range(len(df)):
558 #     #print(str(df.index[x]))
559 #     newdate = str(df.index[x])
560 #     newdate = newdate[0:10]
561 #     #print("newdate =", newdate)
562 #     dates.append(newdate)
563 #
564 # df['dates'] = dates
565
566 #print (df.head(5))
567 #print (df.tail(10))
568
569 # obj = Series([1, 1. 01, 1. 01, 1. 01])
570 # print(obj.pct_change())
571
572
573
```