

ML Hands-on Workshop @ SPCE

Instructor: Santosh Chapaneri

Jan 2021

NumPy

- Python lists are great. They can store strings, integers, or mixtures.
- NumPy arrays though are **multi-dimensional** and most **engineering** python libraries use them instead.
- They store the **same type of data** in each element and **cannot change size**.

In [1]: `import numpy as np`

```
x = np.zeros(5)
print(x)
```

```
[0. 0. 0. 0. 0.]
```

In [2]: `x = np.zeros((5,2))`
`print(x)`

```
[[0. 0.]
 [0. 0.]
 [0. 0.]
 [0. 0.]
 [0. 0.]]
```

In [3]: `print(np.arange(3, 10))` *# Does not include end point*
`print(np.linspace(0, 1, 25))` *# Includes end point*

```
[3 4 5 6 7 8 9]
[0.          0.04166667 0.08333333 0.125          0.16666667 0.20833333
 0.25          0.29166667 0.33333333 0.375          0.41666667 0.45833333
 0.5           0.54166667 0.58333333 0.625          0.66666667 0.70833333
 0.75          0.79166667 0.83333333 0.875          0.91666667 0.95833333
 1.           ]
```

In [4]: `print(np.logspace(0, 1, 25))` *# Log spaced numbers*

```
[ 1.          1.10069417  1.21152766  1.33352143  1.46779927  1.6155981
 1.77827941  1.95734178  2.15443469  2.37137371  2.61015722  2.87298483
 3.16227766  3.48070059  3.83118685  4.21696503  4.64158883  5.10896977
 5.62341325  6.18965819  6.81292069  7.49894209  8.25404185  9.08517576
10.          ]
```

Numpy Arithmetic Operations

```
In [5]: # Trivial math
x = np.array( [[50.0, 60.0], [70.0, 80.0]] )
y = np.array( [[10.0, 20.0], [30.0, 40.0]] )
```

```
print( np.add(x, y) )
print( np.subtract(x, y) )
print( np.multiply(x, y) )
print( np.divide(x, y) )
```

```
[[ 60.  80.]
 [100. 120.]]
[[40. 40.]
 [40. 40.]]
[[ 500. 1200.]
 [2100. 3200.]]
[[5.      3.      ]
 [2.33333333 2.    ]]
```

```
In [6]: # Element-wise sqrt of matrix
print( np.sqrt(x) )
```

```
[[7.07106781 7.74596669]
 [8.36660027 8.94427191]]
```

```
In [7]: # Dot product
v = np.array([9.0, 10.0])
w = np.array([11.0, 12.0])

print( np.dot(v, w) ) # 9 x 11 + 10 x 12
```

```
219.0
```

```
In [8]: # Dot product of matrices
print(x)
print(y)
print( np.dot(x, y) ) # 50 x 10 + 60 x 30 = 2300, etc.
```

```
[[50. 60.]
 [70. 80.]]
[[10. 20.]
 [30. 40.]]
[[2300. 3400.]
 [3100. 4600.]]
```

```
In [9]: # Sum along row or column
print(x)
print( 'Sum along columns:', np.sum(x, axis = 0) )
print( 'Sum along rows:', np.sum(x, axis = 1) )
```

```
[[50. 60.]
 [70. 80.]]
Sum along columns: [120. 140.]
Sum along rows: [110. 150.]
```

```
In [10]: # Matrix reshaping
a = np.arange(40).reshape(5, 8)
print(a)
```

```
[[ 0  1  2  3  4  5  6  7]
 [ 8  9 10 11 12 13 14 15]
 [16 17 18 19 20 21 22 23]
 [24 25 26 27 28 29 30 31]
 [32 33 34 35 36 37 38 39]]
```

```
In [11]: # ALL rows and 4th column
print( a[:, 3] )
```

```
[ 3 11 19 27 35]
```

```
In [12]: # What's happening here?
print( a[1::2, ::3] )
```

```
[[ 8 11 14]
 [24 27 30]]
```

```
In [13]: # And here?
print( a[-3:, -3:] )
```

```
[[21 22 23]
 [29 30 31]
 [37 38 39]]
```

```
In [14]: a = np.arange(36).reshape(3, 4, 3)
# read from left to right: arr is a 3 by 4 by 3 array
a
```

```
Out[14]: array([[[ 0,  1,  2],
 [ 3,  4,  5],
 [ 6,  7,  8],
 [ 9, 10, 11]],

 [[12, 13, 14],
 [15, 16, 17],
 [18, 19, 20],
 [21, 22, 23]],

 [[24, 25, 26],
 [27, 28, 29],
 [30, 31, 32],
 [33, 34, 35]])
```

```
In [15]: # 3 x 2 array, all elements having same value
c = np.full( (3, 2), 8)
print(c)
```

```
[[8 8]
 [8 8]
 [8 8]]
```

```
In [16]: # Identity matrix
d = np.eye(3)
print(d)
```

```
[[1. 0. 0.]
 [0. 1. 0.]
 [0. 0. 1.]]
```

```
In [17]: # Matrix of random numbers
e = np.random.random( (4, 3) )
print(e)
```

```
[[0.00814017 0.20417684 0.67825128]
 [0.96779944 0.29122055 0.65857505]
 [0.72885564 0.22464709 0.40440525]
 [0.17064696 0.10418323 0.36574012]]
```

Code Vectorization

```
In [18]: # Apply cos on each element of the list
from math import pi
x = np.linspace(-pi, pi, 4)
print( np.cos(x) )
```

```
[-1.    0.5   0.5  -1. ]
```

```
In [19]: # Apply sqrt on each element of the List
z = [i**2 for i in range(1,11)]
print(z)
print( np.sqrt(z) )

[1, 4, 9, 16, 25, 36, 49, 64, 81, 100]
[ 1.  2.  3.  4.  5.  6.  7.  8.  9. 10.]
```

NumPy Exercises

Q1) Create a 3x3 matrix with values ranging from 0 to 8.

```
In [20]: Z = np.arange(9).reshape(3, 3)
print(Z)

[[0 1 2]
 [3 4 5]
 [6 7 8]]
```

Q2) Find indices of non-zero elements from [1,2,0,0,4,0]. Hint: Use np.nonzero

```
In [21]: nz = np.nonzero([1,2,0,0,4,0])
print(nz)

(array([0, 1, 4]),)
```

Q3) Create a 2d array with 1 on the border and 0 inside.

```
In [22]: Z = np.ones((10,10))
Z[1:-1, 1:-1] = 0
print(Z)

[[1.  1.  1.  1.  1.  1.  1.  1.  1.  1.]
 [1.  0.  0.  0.  0.  0.  0.  0.  0.  1.]
 [1.  0.  0.  0.  0.  0.  0.  0.  0.  1.]
 [1.  0.  0.  0.  0.  0.  0.  0.  0.  1.]
 [1.  0.  0.  0.  0.  0.  0.  0.  0.  1.]
 [1.  0.  0.  0.  0.  0.  0.  0.  0.  1.]
 [1.  0.  0.  0.  0.  0.  0.  0.  0.  1.]
 [1.  0.  0.  0.  0.  0.  0.  0.  0.  1.]
 [1.  1.  1.  1.  1.  1.  1.  1.  1.  1.]]
```

Q4) Consider a random 10 x 2 matrix representing cartesian coordinates, convert them to polar coordinates.

```
In [23]: Z = np.random.random((10,2))
X,Y = Z[:,0], Z[:,1]
R = np.sqrt(X**2 + Y**2)
T = np.arctan2(Y, X)
print(R)
print(T)

[0.58380699 0.57655513 0.15749278 0.66881743 0.74085048 0.5542201
 0.69795708 0.41117945 0.6933996 0.90116015]
[0.9375869 0.23236367 0.19147987 0.39716354 0.02352703 0.13156396
 1.24203099 1.00732616 0.33272647 0.58279335]
```

Pandas

Dataframes

- According to Pandas documentation: *Two-dimensional size-mutable, potentially heterogeneous tabular data structure with labeled axes (rows and columns). Arithmetic operations align on both row and column labels.*
- In human terms, this means that a **dataframe has rows and columns**, can **change size**, and possibly **has mixed data types**.

Peek at the DataFrame contents

- `df.info()` # index & data types
- `df.head(i)` # get first i rows
- `df.tail(i)` # get last i rows
- `df.describe()` # summary stats cols

A very powerful feature in Pandas is **groupby**.

- This function allows us to **group together rows that have the same value in a particular column**.
- Then, we can aggregate this group-by object to compute statistics in each group.

MovieLens 100k movie rating data:

- main page: <http://grouplens.org/datasets/movielens/> (<http://grouplens.org/datasets/movielens/>)
- 100,000 ratings from 1000 users on 1700 movies

In [24]: `import pandas as pd`

```
users = pd.read_csv('u.user', sep='|', index_col='user_id')
```

In [25]: `# print the first 5 rows`
`users.head()`

Out[25]:

	age	gender	occupation	zip_code
user_id				
1	24	M	technician	85711
2	53	F	other	94043
3	23	M	writer	32067
4	24	M	technician	43537
5	33	F	other	15213

```
In [26]: # print the first 10 rows
users.head(10)
```

```
Out[26]:
```

	age	gender	occupation	zip_code
user_id				
1	24	M	technician	85711
2	53	F	other	94043
3	23	M	writer	32067
4	24	M	technician	43537
5	33	F	other	15213
6	42	M	executive	98101
7	57	M	administrator	91344
8	36	M	administrator	05201
9	29	M	student	01002
10	53	M	lawyer	90703

```
In [27]: # print the last 5 rows
users.tail()
```

```
Out[27]:
```

	age	gender	occupation	zip_code
user_id				
939	26	F	student	33319
940	32	M	administrator	02215
941	20	M	student	97229
942	48	F	librarian	78209
943	22	M	student	77841

```
In [28]: # column names
users.columns
```

```
Out[28]: Index(['age', 'gender', 'occupation', 'zip_code'], dtype='object')
```

```
In [29]: # data types of each column
users.dtypes
```

```
Out[29]: age          int64
gender          object
occupation      object
zip_code        object
dtype: object
```

```
In [30]: # number of rows and columns
users.shape
```

```
Out[30]: (943, 4)
```

```
In [31]: # select one column using the DataFrame attribute
users.gender
```

```
Out[31]: user_id
1      M
2      F
3      M
4      M
5      F
..
939    F
940    M
941    M
942    F
943    M
Name: gender, Length: 943, dtype: object
```

```
In [32]: # summarize (describe) the DataFrame
users.describe() # describe all numeric columns
```

```
Out[32]:
```

	age
count	943.000000
mean	34.051962
std	12.192740
min	7.000000
25%	25.000000
50%	31.000000
75%	43.000000
max	73.000000

```
In [33]: users.describe(include=['object']) # describe all object columns
```

```
Out[33]:
```

	gender	occupation	zip_code
count	943	943	943
unique	2	21	795
top	M	student	55414
freq	670	196	9

```
In [34]: users.describe(include='all') # describe all columns
```

```
Out[34]:
```

	age	gender	occupation	zip_code
count	943.000000	943	943	943
unique	NaN	2	21	795
top	NaN	M	student	55414
freq	NaN	670	196	9
mean	34.051962	NaN	NaN	NaN
std	12.192740	NaN	NaN	NaN
min	7.000000	NaN	NaN	NaN
25%	25.000000	NaN	NaN	NaN
50%	31.000000	NaN	NaN	NaN
75%	43.000000	NaN	NaN	NaN
max	73.000000	NaN	NaN	NaN

```
In [35]: # count the number of occurrences of each value
users.gender.value_counts() # most useful for categorical variables
```

```
Out[35]: M    670
         F    273
         Name: gender, dtype: int64
```

```
In [36]: users.age.value_counts() # can also be used with numeric variables
```

```
Out[36]: 30    39
         25    38
         22    37
         28    36
         27    35
         ..
         11     1
         10     1
         73     1
         66     1
         7      1
         Name: age, Length: 61, dtype: int64
```

```
In [37]: # Boolean filtering: only show users with age < 20
young_bool = users.age < 20      # create a Series of booleans...
users[young_bool]                # ...and use that Series to filter rows
```

Out[37]:

	age	gender	occupation	zip_code
user_id				
30	7	M	student	55436
36	19	F	student	93117
52	18	F	student	55105
57	16	M	none	84010
67	17	M	student	60402
...
872	19	F	student	74078
880	13	M	student	83702
887	14	F	student	27249
904	17	F	student	61073
925	18	F	salesman	49036

77 rows × 4 columns

```
In [38]: users[users.age < 20]      # or, combine into a single step
```

Out[38]:

	age	gender	occupation	zip_code
user_id				
30	7	M	student	55436
36	19	F	student	93117
52	18	F	student	55105
57	16	M	none	84010
67	17	M	student	60402
...
872	19	F	student	74078
880	13	M	student	83702
887	14	F	student	27249
904	17	F	student	61073
925	18	F	salesman	49036

77 rows × 4 columns


```
In [39]: # for each occupation in 'users', count the number of occurrences  
users.occupation.value_counts()
```

```
Out[39]: student      196  
other      105  
educator    95  
administrator  79  
engineer    67  
programmer  66  
librarian   51  
writer      45  
executive   32  
scientist   31  
artist      28  
technician  27  
marketing   26  
entertainment 18  
healthcare  16  
retired     14  
salesman    12  
lawyer      12  
none        9  
homemaker   7  
doctor      7  
Name: occupation, dtype: int64
```

```
In [40]: # for each occupation, calculate the mean age  
users.groupby('occupation').age.mean()
```

```
Out[40]: occupation  
administrator  38.746835  
artist         31.392857  
doctor         43.571429  
educator       42.010526  
engineer       36.388060  
entertainment  29.222222  
executive      38.718750  
healthcare     41.562500  
homemaker      32.571429  
lawyer         36.750000  
librarian      40.000000  
marketing      37.615385  
none          26.555556  
other          34.523810  
programmer     33.121212  
retired        63.071429  
salesman       35.666667  
scientist      35.548387  
student        22.081633  
technician     33.148148  
writer         36.311111  
Name: age, dtype: float64
```

```
In [41]: # for each occupation, calculate the minimum and maximum ages
users.groupby('occupation').age.agg(['min', 'max'])
```

Out[41]:

	min	max
occupation		
administrator	21	70
artist	19	48
doctor	28	64
educator	23	63
engineer	22	70
entertainment	15	50
executive	22	69
healthcare	22	62
homemaker	20	50
lawyer	21	53
librarian	23	69
marketing	24	55
none	11	55
other	13	64
programmer	20	63
retired	51	73
salesman	18	66
scientist	23	55
student	7	42
technician	21	55
writer	18	60

```
In [42]: # for each combination of occupation and gender, calculate the mean age
users.groupby(['occupation', 'gender']).age.mean()
```

```
Out[42]: occupation  gender
administrator  F      40.638889
               M      37.162791
artist         F      30.307692
               M      32.333333
doctor         M      43.571429
educator       F      39.115385
               M      43.101449
engineer       F      29.500000
               M      36.600000
entertainment  F      31.000000
               M      29.000000
executive      F      44.000000
               M      38.172414
healthcare     F      39.818182
               M      45.400000
homemaker      F      34.166667
               M      23.000000
lawyer         F      39.500000
               M      36.200000
librarian      F      40.000000
               M      40.000000
marketing      F      37.200000
               M      37.875000
none           F      36.500000
               M      18.600000
other          F      35.472222
               M      34.028986
programmer     F      32.166667
               M      33.216667
retired        F      70.000000
               M      62.538462
salesman       F      27.000000
               M      38.555556
scientist      F      28.333333
               M      36.321429
student        F      20.750000
               M      22.669118
technician     F      38.000000
               M      32.961538
writer         F      37.631579
               M      35.346154
Name: age, dtype: float64
```

Pandas Exercise: IMDB Data

1. Read in 'imdb_1000.csv' and store it in a DataFrame named movies
2. Check the number of rows and columns
3. Check the data type of each column
4. Calculate the average movie duration
5. Sort by duration to find the shortest and longest movies (*Hint: use sort_values*)
6. Count how many movies have each of the content ratings
7. Calculate the average star rating for movies 2 hours or longer, and compare that with the average star rating for movies shorter than 2 hours
8. Calculate the average duration for each genre

```
In [43]: # 1. Read in 'imdb_1000.csv' and store it in a DataFrame named movies
movies = pd.read_csv('imdb_1000.csv')
```

```
In [44]: # 2. Check the number of rows and columns
movies.shape
```

```
Out[44]: (979, 6)
```

```
In [45]: # 3. Check the data type of each column
movies.dtypes
```

```
Out[45]: star_rating    float64
title                object
content_rating       object
genre                object
duration             int64
actors_list          object
dtype: object
```

```
In [46]: # 4. Calculate the average movie duration
movies.duration.mean()
```

```
Out[46]: 120.97957099080695
```

```
In [47]: # 5. Sort the DataFrame by duration to find the shortest and longest movies
movies.sort_values('duration').head(1)
```

```
Out[47]:
```

	star_rating	title	content_rating	genre	duration	actors_list
389	8.0	Freaks	UNRATED	Drama	64	[u'Wallace Ford', u'Leila Hyams', u'Olga Bacla...]

```
In [48]: movies.sort_values('duration').tail(1)
```

```
Out[48]:
```

	star_rating	title	content_rating	genre	duration	actors_list
476	7.8	Hamlet	PG-13	Drama	242	[u'Kenneth Branagh', u'Julie Christie', u'Dere...]

```
In [49]: # 6. Count how many movies have each of the content ratings
movies.content_rating.value_counts()
```

```
Out[49]: R          460
PG-13       189
PG          123
NOT RATED   65
APPROVED    47
UNRATED     38
G           32
PASSED      7
NC-17       7
X           4
GP          3
TV-MA       1
Name: content_rating, dtype: int64
```

```
In [50]: # 7. calculate the average star rating for movies 2 hours or longer,
# and compare that with the average star rating for movies shorter than 2 hours
movies[movies.duration >= 120].star_rating.mean()
```

```
Out[50]: 7.948898678414082
```

```
In [51]: movies[movies.duration < 120].star_rating.mean()
```

```
Out[51]: 7.838666666666657
```

```
In [52]: # 8. Calculate the average duration for each genre
movies.groupby('genre').duration.mean()
```

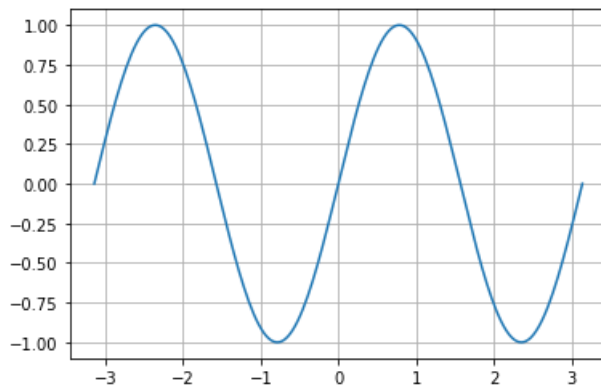
```
Out[52]: genre
Action      126.485294
Adventure   134.840000
Animation    96.596774
Biography    131.844156
Comedy       107.602564
Crime        122.298387
Drama        126.539568
Family       107.500000
Fantasy      112.000000
Film-Noir    97.333333
History       66.000000
Horror        102.517241
Mystery      115.625000
Sci-Fi       109.000000
Thriller     114.200000
Western      136.666667
Name: duration, dtype: float64
```

Matplotlib - Plotting

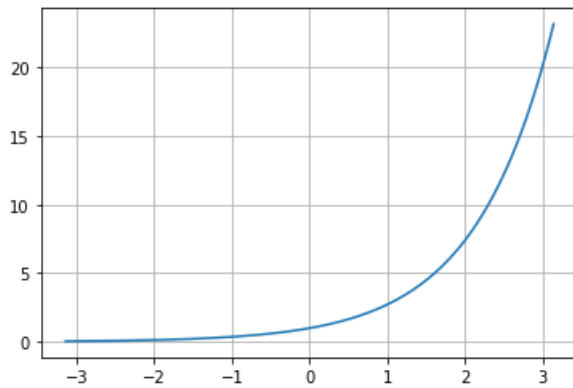
Basic plotting

```
In [53]: import matplotlib.pyplot as plt
import numpy as np
from math import pi

x = np.linspace(-pi, pi, 200)
y = np.sin(2*x)
plt.plot(x, y)
plt.grid()
plt.show()
```

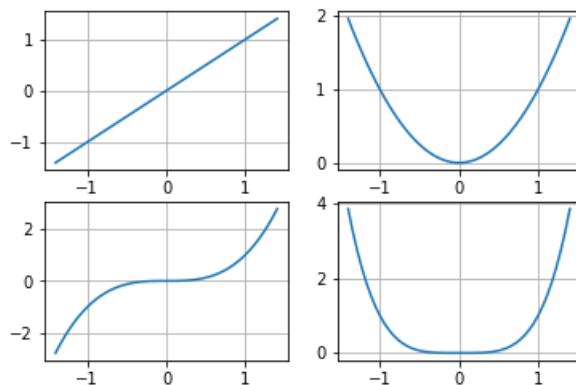


```
In [54]: y = np.exp(x) # also, try with exp(-x)
plt.plot(x, y)
plt.grid()
plt.show()
```



Subplots

```
In [55]: x = np.linspace(-1.4, 1.4, 50)
plt.subplot(2, 2, 1) # 2 x 2, top left
plt.plot(x, x)
plt.grid()
plt.subplot(2, 2, 2) # 2 x 2, top right
plt.plot(x, x**2)
plt.grid()
plt.subplot(2, 2, 3) # 2 x 2, bottom left
plt.plot(x, x**3)
plt.grid()
plt.subplot(2, 2, 4) # 2 x 2, bottom right
plt.plot(x, x**4)
plt.grid()
plt.show()
```



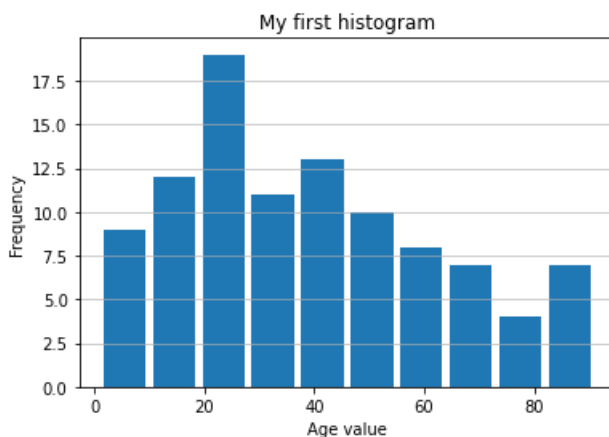
Histograms

```
In [56]: # Histogram of Age
x = [1,1,2,3,3,5,7,8,9,10,10,11,11,13,13,15,16,17,18,18,18,19,20,21,21,23,24,
24,25,25,25,25,26,26,27,27,27,27,29,30,30,31,33,34,34,34,35,36,
36,37,37,38,38,39,40,41,41,42,43,44,45,45,46,47,48,48,49,50,51,52,53,54,
55,55,56,57,58,60,61,63,64,65,66,68,70,71,72,74,75,77,81,83,84,87,89,
90,90,91]

# rwidth: bar's relative width, default 1
plt.hist(x, bins=10, rwidth = 0.85)

# Only the y grid
plt.grid(axis='y', alpha=0.75)

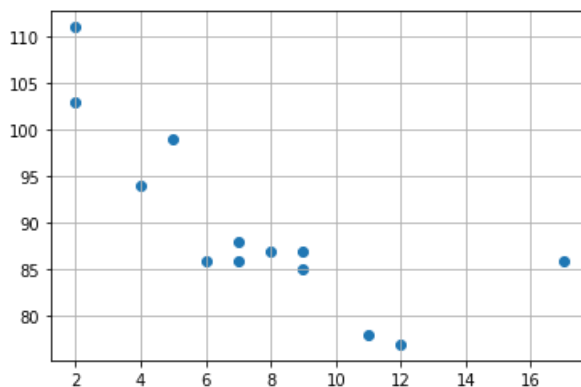
plt.xlabel('Age value')
plt.ylabel('Frequency')
plt.title('My first histogram')
plt.show()
```



Scatter Plot

```
In [57]: x = [5,7,8,7,2,17,2,9,4,11,12,9,6]
y = [99,86,87,88,111,86,103,87,94,78,77,85,86]
# Note: Both x and y must have same Length

plt.scatter(x, y)
plt.grid()
plt.show()
```

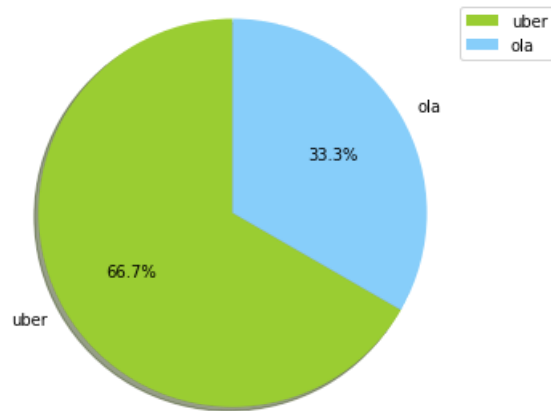


Pie Chart

- Pie chart shows the size of items in a data series, proportional to the sum of the items.
- The data points in a pie chart are shown as a percentage of the whole pie.

```
In [58]: sizes = [10, 5]
mylabels = ['uber', 'ola']

plt.pie(
    sizes,
    labels = mylabels,
    shadow = True,
    colors = ['yellowgreen', 'lightskyblue'],
    startangle = 90,    # rotate conter-clockwise by 90 degrees
    autopct = '%1.1f%', # display fraction as percentage
)
plt.legend(fancybox=True)
plt.axis('equal')      # plot pyplot as circle
plt.tight_layout()
plt.show()
```



```
In [59]: sizes = [10, 5]
mylabels = ['uber', 'ola']
myexplode = [0, 0.2] # space between slices

plt.pie(
    sizes,
    labels = mylabels,
    shadow = True,
    explode = myexplode,
    colors = ['yellowgreen', 'lightskyblue'],
    startangle = 90,    # rotate conter-clockwise by 90 degrees
    autopct = '%1.1f%', # display fraction as percentage
)
plt.legend(fancybox=True)
plt.axis('equal')      # plot pyplot as circle
plt.tight_layout()
plt.show()
```

