# Reading and Writing Files

Han-fen Hu

UNLV

---

# Outline

☐ Files and File Paths

☐ os.path Module
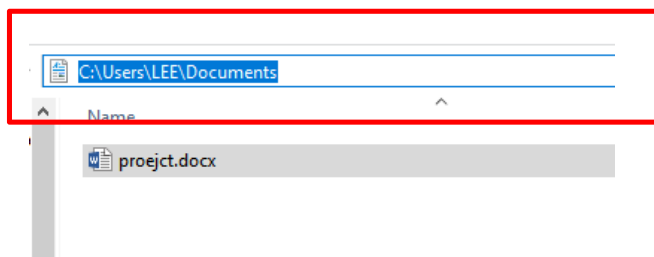
☐ File Read/Write Process

☐ Saving variables

UNLV

# Files

☐ A way to persistently save the data

☐ A way to have input from a static storage

☐ An essential step for automate things

☐ A file's contents can be considered as a single string value

UNLV

---

# Files and File Paths (1)

☐ A file has two key properties: a filename (usually written as one word) and a path

– Filename includes the name and file extension

– Path specifies the location of a file on the computer

C:\Users\LEE\Documents          Path

Name

proejct.docx

UNLV

# Files and File Paths (2)

◻ On Windows, paths are written using backslashes (\) as the separator

◻ OS X and Linux, however, use the forward slash (/) as their path separator.

◻ If you want your programs to work on all operating systems, you will have to write your Python programs to handle both cases

◻ `os.sep` variable

  – set to the correct folder-separating slash for the computer running the program

UNLV

# Files and File Paths (3)

◻ Current Working Directory

  – Every program that runs on your computer has a current working directory (cwd)

  – Any filenames or paths that do not begin with the root folder are assumed to be under the current working directory

UNLV

# cwd Example (1)

❑ Get the current working directory as a string value with the `os.getcwd()`

❑ Change the working directory with os.chdir()

```
1  import os
2  print(os.getcwd())
3
4  # change the working directory
5  os.chdir('C:\\Users\\Default\\Documents')
6
7  print(os.getcwd())
```

```
D:\backup\Fall2020_MIS740
C:\Users\Default\Documents
```

7

---

# cwd Example (2)

❑ Python will display an error if you try to change to a directory that does not exist

```
1  import os
2  print(os.getcwd())
3
4  # change the working directory
5  os.chdir('C:\\Users\\HH\\Documents')
6
7  print(os.getcwd())
```

```
C:\Users\Default\Documents

---------------------------------------------------------------------
FileNotFoundError                         Traceback (most recent call last)
<ipython-input-3-d724a238135f> in <module>
      3
      4 # change the working directory
----> 5 os.chdir('C:\\Users\\HH\\Documents')
      6
      7 print(os.getcwd())

FileNotFoundError: [WinError 3] The system cannot find the path specified: 'C:\\Users\\HH\\Documents'
```
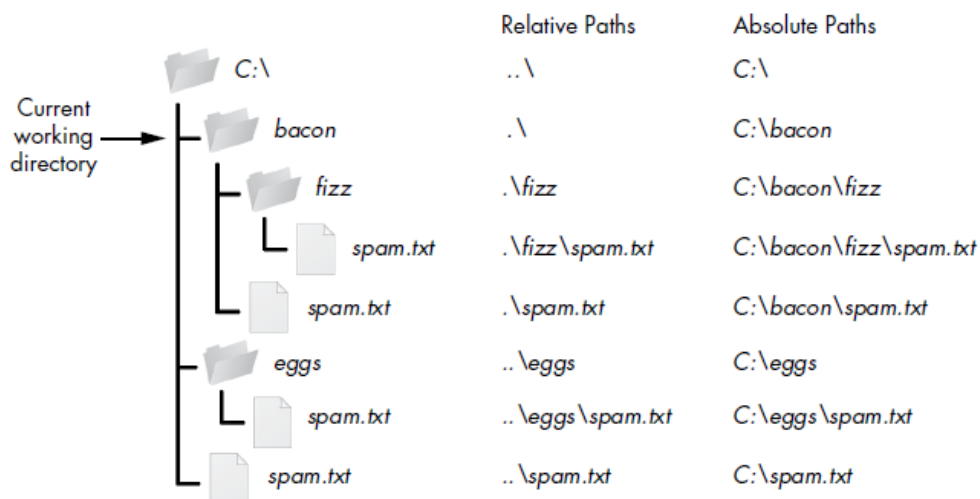
8

# Absolute vs. Relative Paths

☐ Absolute path always begins with the root folder

☐ Relative path is relative to the program's current working directory

– dot (.) folder: shorthand for "this directory."

– dot-dot (..) folders: Means "the parent folder."

– Not real folders but special names that can be used in a path.

UNLV

---

# Absolute vs. Relative Paths: Example

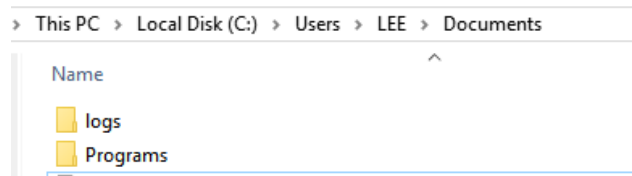| | Relative Paths | Absolute Paths |
|---|---|---|
| C:\ | ..\ | C:\ |
| bacon | .\ | C:\bacon |
| fizz | .\fizz | C:\bacon\fizz |
| spam.txt | .\fizz\spam.txt | C:\bacon\fizz\spam.txt |
| spam.txt | .\spam.txt | C:\bacon\spam.txt |
| eggs | ..\eggs | C:\eggs |
| spam.txt | ..\eggs\spam.txt | C:\eggs\spam.txt |
| spam.txt | ..\spam.txt | C:\spam.txt |

Current working directory → bacon

UNLV

# Creating New Folders

☐ Your programs can create new folders (directories) with the `os.makedirs()`

```
1  import os
2  # change the working directory
3  os.chdir('C:\\Users\\LEE\\Documents')
4
5  # create a folder under current wokring directory
6  os.makedirs('logs')
7
8  # create a folder with absolue path
9  os.makedirs('C:\\Users\\LEE\\Documents\\Programs')
```

> This PC > Local Disk (C:) > Users > LEE > Documents

Name
⌃
📁 logs
📁 Programs

UNLV

---

# os.path Module (1)

☐ Contains many helpful functions related to filenames and file paths

– merging, normalizing and retrieving path names in python

☐ Full documentation: http://docs.python.org/3/library/os.path.html

UNLV

# os.path Module (2)

❑ Checking Path Validity

   – Many Python functions will crash with an error if you supply them with a path that does not exist.

❑ os.path.exists(*path*)

   – Return True if the file or folder referred to in the argument exists and will return False if it does not exist

```python
1   import os
2   # change the working directory
3   os.chdir('C:\\Users\\LEE\\Documents')
4
5   if not os.path.exists('logs'):
6       # create a folder under current wokring directory
7       os.makedirs('logs')
8       print('Foleder "logs" created')
9   else:
10      print('Foleder "logs" already exists')
```

13

# os.path Module (3)

❑ os.path.join()

   – Build paths in a way that will work on any operating system

      • \ for windows; / for OS X and Linux

```python
1   import os # for file and path operation
2   import datetime # for getting current date time
3
4   # get the current month
5   currentMonth = str(datetime.datetime.today().month)
6
7   # generate the path for storing the log files
8   path = os.path.join('app','logs', currentMonth)
9
10  # print it out to verify its correctness
11  print(path)
```

app\logs\10

14

# os.path Module (4)

- ❑ **os.path.abspath(path)**
  - – Return a string of the absolute path of the argument
  - – Converts a relative path into an absolute one

- ❑ **os.path.isabs(path)**
  - – Return True if the argument is an absolute path and False if it is a relative path

- ❑ **os.path.relpath(path, *start*)**
  - – Return a string of a relative path from the *start* path to path.
  - – If *start* is not provided, the current working directory is used as the start path

UNLV

# os.path Module (5)

```python
import os
# show the current working directory
print("Current working directory: "+os.getcwd())

# convert  .. to absotue path, save the string to the variable
absolutePath = os.path.abspath('..')
# print the variable
print("The absolute path of .. : "+absolutePath)

# convert .\\logs to absotue path and print it
print("The absolute path of .\\logs : "+os.path.abspath('.\\logs'))

# show whether . is an absolute path, should be false
print(os.path.isabs('.'))
# convert . to absolute path,
# and check whether the conversion result is an absolute path
print(os.path.isabs(os.path.abspath('.')))
```

```
Current working directory: C:\Users\LEE\Documents
The absolute path of .. : C:\Users\LEE
The absolute path of .\logs : C:\Users\LEE\Documents\logs
False
True
```

UNLV

# os.path Module (6)

```
1  import os
2  # show the current working directory (CWD)
3  print("Current working directory: "+os.getcwd())
4
5  # get the path of CWD, relative to C:\\, assign the result ot a variable
6  relativePath = os.path.relpath(os.getcwd(), 'C:\\')
7  # print the variable
8  print(relativePath)
9
10 # show the path of CWD, relative to C:\\Windowss
11 print(os.path.relpath(os.getcwd(), 'C:\\Windows'))
12
13
14 # show the path of C:\\Windowss, relative to CWD
15 print(os.path.relpath('C:\\Windows', os.getcwd()))
```

```
Current working directory: C:\Users\LEE\Documents
Users\LEE\Documents
..\Users\LEE\Documents
..\..\..\Windows
```

UNLV

---

# os.path Module (7)

☐ os.path.dirname(*path*)

— Return a string of everything that comes before the last slash in the path argument.

☐ os.path.basename(*path*)

— Return a string of everything that comes after the last slash in the path argument

```
C:\Windows\System32\calc.exe
```
Dir name        Base name

UNLV

# os.path Module (8)

```
1  dataFilePath = 'C:\\Users\\LEE\\Documents\\sales2019.xlsx'
2  pathTuple = os.path.split(dataFilePath)
3  print(pathTuple)
```

```
('C:\\Users\\LEE\\Documents', 'sales2019.xlsx')
```

## Is equivalent to

```
1  dataFilePath = 'C:\\Users\\LEE\\Documents\\sales2019.xlsx'
2  pathTuple = (os.path.dirname(dataFilePath), os.path.basename(dataFilePath))
3  print(pathTuple)
```

UNLV

19

---

# os.path Module (9)

☐ os.sep variable

– set to the correct folder-separating slash for the computer running the program

• For Windows

```
1  dataFilePath = 'C:\\Users\\LEE\\Documents\\sales2019.xlsx'
2  print(dataFilePath.split(os.sep))
```

```
['C:', 'Users', 'LEE', 'Documents', 'sales2019.xlsx']
```

• For OS X and Linux

```
1  dataFilePath = '/usr/bin/sales2019.xlsx'
2  print(dataFilePath.split(os.sep))
```

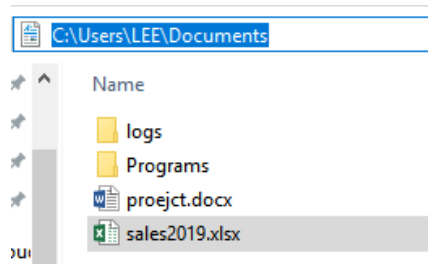```
['', 'usr', 'bin', 'sales2019.xlsx']
```

UNLV

20

# os.path Module (10)

❑ os.listdir(*path*)

– Return a list of filename strings for each file in the *path* argument.

```
1  fileList = os.listdir(os.getcwd())
2  print(fileList)
```

['desktop.ini', 'logs', 'My Music', 'My Pictures', 'My Videos', 'proejct.docx', 'Programs', 'sales2019.xlsx']

C:\Users\LEE\Documents

Name

logs
Programs
proejct.docx
sales2019.xlsx

UNLV

# Exercise / Question

❑ What does the following program do?

```
1  for filename in os.listdir(os.getcwd()):
2      print(os.path.join(os.getcwd(), filename))
```

UNLV

# Types of Files

☐ Plain text files

- – Contain only basic text characters and do not include font, size, or color information
- – Can be opened with Windows's Notepad or OS X's TextEdit application
- – With .txt or .csv file extension

☐ Binary files are all other file types, such as word processing documents, PDFs, images, spreadsheets, and executable programs

- – Every different type of binary file must be handled in its own way

UNLV

---

# File Reading/Writing Process

1. Call the `open()` function to return a File object.

2. Call the `read()` or `write()` method on the `File` object.

3. Close the file by calling the `close()` method on the `File` object.

UNLV

# Open Files

☐ open()

- – Pass it a string path indicating the file you want to open

- – The path can be either an absolute or relative path

- – It returns a `File` object

  - • A File object is simply another type of value in Python, much like the lists and dictionaries

- – The file will be opened in "reading plaintext" mode

  - • Can't write or modify it in any way

UNLV

# Reading the Contents of Files

☐ read()

- – Read the entire contents of a file as a string value

☐ readlines()

- – get a list of string values from the file, one string for each line of text

UNLV

# Close the File

☐ close()

– Close the File.

UNLV

---

☐ Lab

– ReadFileContent
This program reads the content of a file and
show it on the screen

UNLV

# Writing to Files (1)

❑ The file needs to be open in "write plaintext" mode or "append plaintext" mode

- Write mode
  - Overwrite the existing file and start from scratch
  - Pass 'w' as the second argument to open() to open the file in write mode.
- Append mode
  - Append text to the end of the existing file
  - Pass 'a' as the second argument to open() to open the file in append mode.

UNLV

---

# Writing to Files (2)

❑ write(*string*)

- Write *string* to the file
- It does **not** automatically add a newline character to the end of the string

UNLV

# Writing to Files: Example

```python
1   # open the file in the write mode
2   baconFile = open('bacon.txt', 'w')
3   # write to the file; overwrite everything
4   baconFile.write('Hello world!\n')
5   # close the file
6   baconFile.close()
7
8   # open the file in the append mode
9   baconFile = open('bacon.txt', 'a')
10  # write to the file, append the new content to the end
11  baconFile.write('Bacon is not a vegetable.')
12  # close the file
13  baconFile.close()
14
15  # open the file in the read mode
16  baconFile = open('bacon.txt')
17  # read the content
18  content = baconFile.read()
19  # close the file
20  baconFile.close()
21
22  # print the content
23  print(content)
24
```

```
Hello world!
Bacon is not a vegetable.
```

□Lab

– Degree Courses
In this program, the user will enter names of the courses he/she is taking this semester. Write the input value to a file

❑Exercise

– Movie Record
Please write a program that reads the movie
record from a file (MovieBoxOffice.txt) and
shows the content to the user. The user can
add a record by entering a move name and
its box office.

UNLV

# NumPy (1)

Han-fen Hu

UNLV

---

# Outline

☐ Introduction to NumPy

☐ NumPy Arrays

☐ Computation on NumPay Arrays

UNLV

# Introduction to NumPy (1)

- ❑ **Why NumPy?**
  - – One of the most powerful Python libraries
  - – Improve how data is stored and manipulated
  - – Contains a multi-dimensional array and matrix data structures
  - – Pandas relies heavily on NumPy
- ❑ **Purpose**
  - – Store in-memory data in a more efficient way
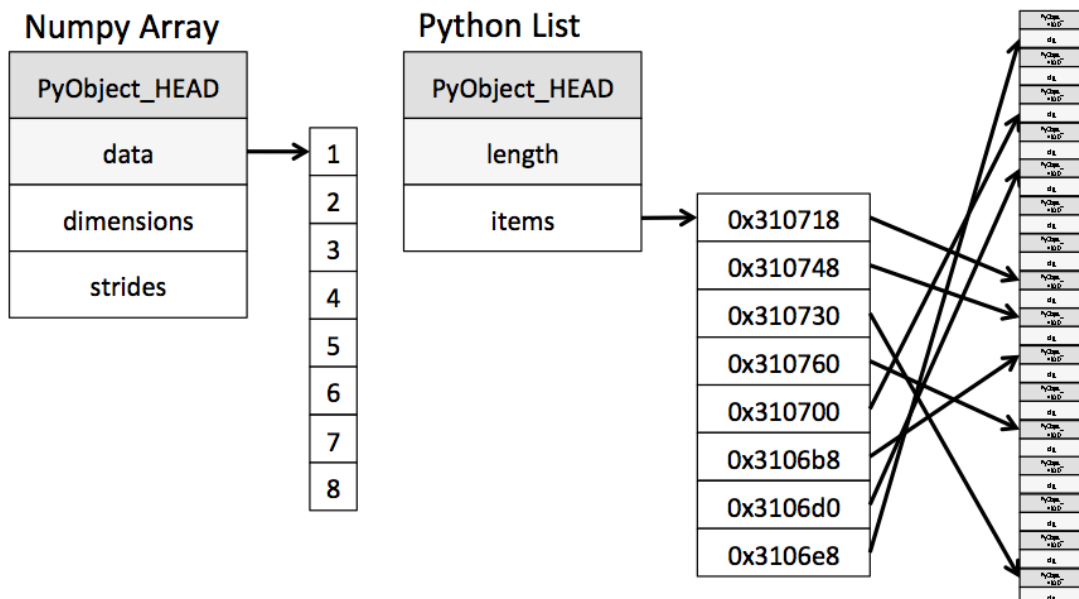  - – Includes a large number of mathematical, algebraic and transformation functions
- ❑ **Installed with Anaconda**

**UNLV**

---

# Introduction to NumPy (2)

- ❑ **A more efficient way to store data**

# NumPy Arrays (1)

☐ A collection of relevant data

☐ Fixed-Type: All items in the array are of the same data type

– Items of a Python list can be of different data types

UNLV

# Popular NumPy Data Types

| Data type | Description |
|---|---|
| bool_ | Boolean (True or False) stored as a byte |
| int_ | Default integer type (normally either int64 or int32) |
| int8 | Byte (-128 to 127) |
| int16 | Integer (-32768 to 32767) |
| int32 | Integer (-2147483648 to 2147483647) |
| int64 | Integer (-9223372036854775808 to 9223372036854775807) |
| float_ | Shorthand for float64. |
| float16 | Half precision float: sign bit, 5 bits exponent, 10 bits mantissa |
| float32 | Single precision float: sign bit, 8 bits exponent, 23 bits mantissa |
| float64 | Double precision float: sign bit, 11 bits exponent, 52 bits mantissa |

UNLV

# Creating NumPy Arrays (1)

☐ From Python lists

```python
1   # import numpy, set the alias as np
2   import numpy as np
3
4   # declare a list
5   ageList = [58,69,32,53,81,60,18,25]
6
7   # convert the list to a numpy array
8   ageArray = np.array(ageList)
9   print(ageArray)
10
11  ageArray = np.array(ageList, dtype='float32')
12  print(ageArray)
```

```
[58 69 32 53 81 60 18 25]
[58. 69. 32. 53. 81. 60. 18. 25.]
```

**UNLV**

# Creating NumPy Arrays (2)

☐ Creating Arrays from Scratch

– Specify the size

– Specify the data type

– Specify the values

```python
1   # import numpy, set the alias as np
2   import numpy as np
3
4   # create an array with all 0s, single-dimension
5   zeroArray = np.zeros(10, dtype='int')
6   # create an array with all 1s, a 3X5 array
7   oneArray = np.ones((3,5), dtype='float')
8   # Create a 3x3 identity matrix
9   eyeArray = np.eye(3)
10  # create an array filled with 3.14, a 2X6 array
11  piArray = np.full((2,6), 3.14)
12
13  print(zeroArray)
14  print()
15  print(oneArray)
16  print()
17  print(eyeArray)
18  print()
19  print(piArray)
```

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

[[1. 1. 1. 1. 1.]
 [1. 1. 1. 1. 1.]
 [1. 1. 1. 1. 1.]]

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

[[3.14 3.14 3.14 3.14 3.14 3.14]
 [3.14 3.14 3.14 3.14 3.14 3.14]]
```

# Creating NumPy Arrays (3)

- ☐ arange(*start, stop, step*)
  - Create an array filled with a linear sequence, with *start*, *stop*, and *step* values

```
1  # Create an array filled with a linear sequence
2  # Starting at 20, ending at 65, stepping by 5
3  checkInAges = np.arange(20, 65, 5)
4  print(checkInAges)
```

```
[20 25 30 35 40 45 50 55 60]
```

- ☐ linspace(*LowerBound, upperBound, numberOfValues*)
  - Create an array with values between two numbers

```
1  # Create an array of 6 values evenly spaced between 5 and 20
2  spacedNumbers = np.linspace(5, 20, 6)
3  print(spacedNumbers)
```

```
[ 5.  8. 11. 14. 17. 20.]
```

UNLV

---

# Creating NumPy Arrays (4)

- ☐ random.random(*numberOfRows, numberOfColumns*)
  - Create an array of uniformly distributed random values between 0 and 1

- ☐ random.normal(*mean, standardDeviation, size*)
  - Create an array of normally distributed random values with *mean* and *standard deviation*

- ☐ random.randint(*LowerBound, upperBound, size*)
  - *Create an array of random integers in the interval*

UNLV

# Creating NumPy Arrays (4)

```
1  # Create a 2x6 array of uniformly distributed
2  # random values between 0 and 1
3  randNumbers = np.random.random((2, 6))
4  print(randNumbers)
```

```
[[0.62425223 0.58639327 0.90045935 0.65651807 0.61259153 0.95764916]
 [0.57857881 0.51963987 0.52926542 0.12380843 0.69996986 0.24306753]]
```

```
1  # Create a 3x5 array of normally distributed random values
2  # with mean 0 and standard deviation 1
3  normalRandNumbers = np.random.normal(0, 1, (3, 5))
4  print(normalRandNumbers)
```

```
[[ 1.90726418  0.85365738  1.22608166 -0.3359527  -0.28892326]
 [-1.58078791  1.1592405   0.25509415  0.75614964  0.97189673]
 [ 0.07902027 -2.62662344 -0.67627699  1.98106937 -0.49556422]]
```

```
1  # Create a 3x2 array of random integers in the interval 1 and 53
2  intRandNumbers = np.random.randint(1, 53, (3, 2))
3  print(intRandNumbers)
```

```
[[19 32]
 [20 41]
 [43 35]]
```

11

---

# NumPy Array Attributes

- ndim: number of dimensions

- shape: the size of each dimension

- size: the total size of the array

- dtype: the data type of the array items

```
1  # Create a 3X4X5 array of random integers in the interval 1 and 100
2  threeDArray = np.random.randint(1, 100, (3, 4, 5))
3  print("arry dimension: ", threeDArray.ndim)
4  print("array shape:", threeDArray.shape)
5  print("array size: ", threeDArray.size)
6  print("array data type: ", threeDArray.dtype)
```

```
arry dimension:  3
array shape: (3, 4, 5)
array size:  60
array data type:  int32
```

12

# Accessing Items with Index (1)

❑ In a one-dimensional array, the *i*th value (counting from zero) can be accessed by specifying the desired index in square brackets, just as with Python lists

❑ NumPy slicing syntax follows that of the standard Python list

UNLV

---

# Accessing Items with Index (2)

```
1   # some array
2   randArray = np.random.randint(1, 100, (10))
3   print(randArray)
4
5   print('The third: '+ str(randArray[2]))
6   print('The second from the right: '+ str(randArray[-2]))
7   print('First three item: ', randArray[:3])
8   print('Item 5 to 7: ', randArray[4:7])
9   print('update the value of the 3 item as 200')
10  randArray[2] =200 # assign the value to a specific item
11  print(randArray)
```

```
[32 23  3 87 40 29 15 31 90 51]
The third: 3
The second from the right: 90
First three item:  [32 23  3]
Item 5 to 7:  [40 29 15]
update the value of the 3 item as 200
[ 32  23 200  87  40  29  15  31  90  51]
```

UNLV

# Question: What's the result?

```
1  import numpy as np
2  twoDArray = np.random.randint(1, 100, (2,5))
3
4  print (twoDArray)
5
6  print(twoDArray[0,1])
7  print(twoDArray[1,-1])
```

```
[[77 14 92  3 10]
 [84 85 31 53 34]]
```

# Accessing Items with Index (2)

☐ We can also slice multi-dimensional arrays

```
1  twoDArray = np.random.randint(1, 100, (6,7))
2  print(twoDArray)
3
4  print(twoDArray[:2, :3]) # first two rows, first three columns
5  print(twoDArray[:,-1:]) # all rows, last column
```

```
[[34 20 59 82 31  2 69]
 [ 2 85 73 82 32 37 91]
 [85 82 53 67 91 91 92]
 [85 69 76 96 25 37 17]
 [ 5 77 75 50 62 60 27]
 [97 87 90 92 55  8  3]]
[[34 20 59]
 [ 2 85 73]]
[[69]
 [91]
 [92]
 [17]
 [27]
 [ 3]]
```

# Subarrays as No-Copy Views (1)

☐ One important–and extremely important– feature about array slices is that they return *views* rather than *copies* of the array data.

– Key aspect of NumPy array slicing that differs from Python list slicing

– When the value in the sliced subarray is updated, the value is reflected in the original array as well

☐ It is useful when we work with large datasets. We can access and process pieces of these datasets without the need to copy the entire dataset.

UNLV

---

# Subarrays as No-Copy Views (2)

```
1   twoDArray = np.random.randint(1, 100, (6,7))
2   print(twoDArray)
3
4   subArray = twoDArray[:2, :2] # 2X2 subarray
5   print(subArray)
6
7   print('Update subarray[1,1] as 300')
8   subArray[1,1]=300
9
10  print(twoDArray)
```

```
[[72 33 98 37 11 74  7]
 [33 66 49 71 40  9 30]
 [ 2 39 32 66 11 50 75]
 [43 13 71 89 19 60 88]
 [ 4 21 47 32 51 18 61]
 [54 49 67 41 33  1 14]]
[[72 33]
 [33 66]]
Update subarray[1,1] as 300
[[ 72  33  98  37  11  74   7]
 [ 33 300  49  71  40   9  30]
 [  2  39  32  66  11  50  75]
 [ 43  13  71  89  19  60  88]
 [  4  21  47  32  51  18  61]
 [ 54  49  67  41  33   1  14]]
```

UNLV

# Subarrays as No-Copy Views (3)

☐ It is sometimes useful to instead explicitly copy the data within an array or a subarray

```
1  twoDArray = np.random.randint(1, 100, (6,7))
2  print(twoDArray)
3
4  subArrayCopy = twoDArray[:2, :2].copy() # 2X2 subarray, but as a copy
5  print(subArrayCopy)
6
7  print('Update subarrayCopy[1,0] as 400')
8  subArrayCopy[1,0]=400
9
10 print('Updated copied subarray')
11 print(subArrayCopy)
12 print('Original Array:')
13 print(twoDArray)
```

```
[[98 43 35 88 63 80 97]
 [53 25 23 59 90 21 61]
 [57 98 33 29 88 68 43]
 [19 20  9 44 69 93 78]
 [79 21 88  1 73 77 78]
 [87 48 37  7 97 13 99]]
[[98 43]
 [53 25]]
Update subarrayCopy[1,0] as 400
Updated copied subarray
[[ 98  43]
 [400  25]]
Original Array:
[[98 43 35 88 63 80 97]
 [53 25 23 59 90 21 61]
 [57 98 33 29 88 68 43]
 [19 20  9 44 69 93 78]
 [79 21 88  1 73 77 78]
 [87 48 37  7 97 13 99]]
```

**UNLV**

---

# Reshape Array

☐ reshape(*newRowSize, newColumnSize*)

– the size of the initial array must match the size of the reshaped array

```
1  twoDArray = np.random.randint(1, 100, (2,5))
2
3  print(twoDArray)
4  # row vector via reshape
5  row=twoDArray.reshape((1, 10))
6  print(row)
```

```
[[34 41  4 89  4]
 [60 76  9 29 28]]
[[34 41  4 89  4 60 76  9 29 28]]
```

**UNLV**

# Array Concatenation (1)

□ np.concatenate(*array1, array2, ….*)

– For uni-dimensional array, add the items to the same array

– For multi-dimensional array, concatenate along a specific axis

```
1  grid = np.array([[1, 2, 3],
2                   [4, 5, 6]])
3  grid2 = np.array([[11, 12, 13],
4                    [14, 15, 16]])
5
6  # concatenate along the first axis (i.e., adding rows)
7  print(np.concatenate([grid, grid], axis=0))
8
9  # concatenate along the first axis (i.e., adding columns)
10 print(np.concatenate([grid, grid2], axis=1))
```

```
1  x = np.array([1, 2, 3])
2  y = np.array([3, 2, 1])
3  z = [99, 99, 99]
4  print(np.concatenate([x, y, z]))
```

```
[ 1  2  3  3  2  1 99 99 99]
```

```
[[1 2 3]
 [4 5 6]
 [1 2 3]
 [4 5 6]]
[[ 1  2  3 11 12 13]
 [ 4  5  6 14 15 16]]
```

21

# Array Concatenation (2)

□ np.vstack (array1, array2, ….)

□ np.hstack (array1, array2, ….)

– Concatenarate arrays of mixed dimensions

```
1  x = np.array([1, 2, 3])
2  grid = np.array([[9, 8, 7],
3                   [6, 5, 4]])
4
5  # vertically stack the arrays
6  print(np.vstack([x, grid]))
7
8  # horizontally stack the arrays
9  y = np.array([[99],
10               [99]])
11 print(np.hstack([grid, y]))
```

```
[[1 2 3]
 [9 8 7]
 [6 5 4]]
[[ 9  8  7 99]
 [ 6  5  4 99]]
```

22

UNLV

# Computation on NumPay Arrays

□ **Universal Functions**

- – Vectorized operations to improve the performance of calculation
- – Instead of using for loops to process each item in an array, the universal function can be used to make repeated calculations on array elements much more efficient

□ **Array arithmetic**

□ **Absolute value**

□ **Exponents and logarithms**

□ **Other functions: Trigonometric functions**

UNLV

23

---

# Array Arithmetic (1)

□ Standard addition, subtraction, multiplication, and division

```
1  x = np.arange(5)
2  print("x      =", x)
3  print("x + 5 =", x + 5)
4  print("x - 5 =", x - 5)
5  print("x * 2 =", x * 2)
6  print("x / 2 =", x / 2)
7  print("x // 2 =", x // 2)   # floor division
8  print("x ** 2 =", x ** 2)   # exponentiation
9  print("x % 2  =", x % 2)    # modulus
```

```
x      = [0 1 2 3 4]
x + 5 = [5 6 7 8 9]
x - 5 = [-5 -4 -3 -2 -1]
x * 2 = [0 2 4 6 8]
x / 2 = [0.  0.5 1.  1.5 2. ]
x // 2 = [0 0 1 1 2]
x ** 2 = [ 0  1  4  9 16]
x % 2  = [0 1 0 1 0]
```

UNLV

24

# Array Arithmetic (2)

◻ The operators can be used together in an expression.

– Standard order of operations is respected

| Operator | Equivalent ufunc | Description |
|---|---|---|
| + | np.add | Addition (e.g., 1 + 1 = 2) |
| - | np.subtract | Subtraction (e.g., 3 - 2 = 1) |
| - | np.negative | Unary negation (e.g., -2) |
| * | np.multiply | Multiplication (e.g., 2 * 3 = 6) |
| / | np.divide | Division (e.g., 3 / 2 = 1.5) |
| // | np.floor_divide | Floor division (e.g., 3 // 2 = 1) |
| ** | np.power | Exponentiation (e.g., 2 ** 3 = 8) |
| % | np.mod | Modulus/remainder (e.g., 9 % 4 = 1) |
| | np.sqrt | Square root (e.g. np.sqrt(9) = 3 |

◻ Lab

– Grade Curving
Write a program that reads a file with scores and applied an equation to curve the grade. The result contains the original and curved grade is written to a new file.

UNLV

☐Exercise

- Height Converter
Please write a program that read a file with some heights in centimeters. Please convert the heights into feet and inches. The result contains the original and converted height should be written to a new file and separated with a comma.

UNLV

# NumPy (2)

Han-fen Hu

UNLV

# Outline

☐ Computation on NumPay Arrays

☐ Aggregations

☐ Broadcasting

☐ Comparisons

UNLV

# Computation on NumPay Arrays

☐ **Universal Functions**

   – Vectorized operations to improve the performance of calculation

   – Instead of using for loops to process each item in an array, the universal function can be used to make repeated calculations on array elements much more efficient

☐ **Array arithmetic**

☐ **Absolute value**

☐ **Exponents and logarithms**

☐ **Other functions: Trigonometric functions**

UNLV

---

# Absolute Value

☐ **np.absolute()**
**np.abs()**

   – Returns the absolute values of the items

```
1  x = np.array([-2, -1, 0, 1, 2])
2  print(np.absolute(x))

[2 1 0 1 2]
```

UNLV

# Exponents and Logarithms

❑ufunc also provides an efficient way to do exponentials and logarithms

```
1  x = [1, 2, 4, 10]
2  print("x       =", x)
3  print("e^x     =", np.exp(x))
4  print("3^x     =", np.power(3, x))
5  print("ln(x)   =", np.log(x)) #natural logarithm
6  print("log2(x)  =", np.log2(x)) #base-2 logarithm
7  print("log10(x) =", np.log10(x)) #base-10 logarithm
```

```
x       = [1, 2, 4, 10]
e^x     = [2.71828183e+00 7.38905610e+00 5.45981500e+01 2.20264658e+04]
3^x     = [    3      9     81 59049]
ln(x)   = [0.         0.69314718 1.38629436 2.30258509]
log2(x)  = [0.         1.         2.         3.32192809]
log10(x) = [0.         0.30103    0.60205999 1.         ]
```

5

UNLV

---

# Aggregations

❑Summing the Values in an Array

   – Multi dimensional aggregates

   – Other aggregation functions

❑Minimum and Maximum

6

UNLV

# Lab

- Average Height of US Presidents
  This program reads the height data from a
  csv file and show the statistics.

- NOTE: This program uses pandas and
  matplotlib that we will cover later in this class

UNLV

# Exercise

- Grade Curving

- In the curvedGrade.csv file, the first column
  is the original scores, and the second
  column shows the curved scores.

  Use the curvedGrade.csv and show how the
  curving changes the distribution of the
  scores, including the min, max, mean,
  standard deviation, and the median.

UNLV

# Binary Operation on Arrays

☐ Binary operation on NumPy arrays are performed on an element-by-element basis

```
1  import numpy as np
2
3  a = np.array([0, 1, 2])
4  b = np.array([5, 5, 5])
5  sum = a+b
6  print(sum)
7
```

[5 6 7]

```
1  a = np.array([0, 1, 2])
2  b = np.array([5, 5, 5])
3
4  difference = b-a
5  print(difference)
```

[5 4 3]

UNLV

# Broadcasting

☐ A set of rules for applying binary ufuncs (e.g., addition, subtraction, multiplication, etc.)

☐ Apply to arrays of different sizes

– Here the one-dimensional array  a  is stretched, or broadcast across the second dimension in order to match the shape of c

```
1  a = np.array([0,1,2])
2  c = np.array([[23,59,61],
3               [55,46,69],
4               [18,43,36]])
5
6  unevenSum = a+c
7  print(unevenSum)
```

[[23 60 63]
 [55 47 71]
 [18 44 38]]

UNLV

10

# Rules of Broadcasting (1)

☐ Rule 1: If the two arrays differ in their number of dimensions, the shape of the one with fewer dimensions is <span style="color:blue">padded with ones on its leading (left) side</span>.

☐ Rule 2: If the shape of the two arrays does not match in any dimension, the array with shape equal to 1 in that dimension is <span style="color:blue">stretched to match</span> the other shape.

☐ Rule 3: If in any dimension the sizes disagree and neither is equal to 1, an error is raised.

UNLV

# Question/ Exercise (1)

☐ What is the result of m+a?

```
1  m = np.ones((2, 3))
2  a = np.array([0,1,2])
3
4  print(m)
5  print(a)
6
7  print(m+a)
```

```
[[1. 1. 1.]
 [1. 1. 1.]]
[0 1 2]
```

UNLV

# Question/ Exercise (2)

❏ What is the result of a+b?

```
1  a = ([[0],[1],[2]])
2  b = np.arange(3)
3  print(a)
4  print(b)
5  print(a+b)
```

```
[[0], [1], [2]]
[0 1 2]
```

```
[[0]          [[0] [0] [0]
 [1]           [1] [1] [1]
 [2]]          [2] [2] [2]]

[0 1 2]       [0 1 2
               0 1 2
               0 1 2]
```

UNLV

---

# Question/ Exercise (3)

❏ What is the result of m+a?

```
1  m = np.ones((3, 2))
2  a = np.arange(3)
3
4  print(m)
5  print(a)
6
7  print(m+a)
```

```
[[1. 1.]
 [1. 1.]
 [1. 1.]]
[0 1 2]
```

UNLV

# Rules of Broadcasting (2)

☐ Broadcasting rules apply to <span style="color:red">any binary ufunc</span>

```
1  ages = np.array([[80,86,88,30],
2                   [50,62,75,23],
3                   [66,58,40,36],
4                   [98,68,93,40]])
5  # using the mean aggregate across the first dimension
6  ageMean = ages.mean(0)
7  print(ageMean)
8
9  # Centering the array
10 agesCentered = ages - ageMean
11 print(agesCentered)
```

```
[73.5  68.5  74.    32.25]
[[  6.5   17.5   14.    -2.25]
 [-23.5   -6.5    1.    -9.25]
 [ -7.5  -10.5  -34.     3.75]
 [ 24.5   -0.5   19.     7.75]]
```

15

☐ Lab

– Curving a Series of Scores

The program reads the original scores from a file and then ask the user to enter the percentage of curving he/she wants to apply. The program then prints the updated score.

NOTE: This program uses pandas that we will cover later in this class

16

# Comparison Operators as ufuncs (1)

☐ NumPy also implements comparison operators as element-wise ufuncs.

 – All six of the standard comparison operations are available

☐ The result of these comparison operators is always an array with a Boolean data type.

# Comparison Operators as ufuncs (2)

```
1  x = np.array([1, 2, 3, 4, 5])
2  print(x < 3)
```

[ True   True False False False]

```
1  x = np.array([1, 2, 3, 4, 5])
2  print(x > 3)
```

[False False False  True   True]

```
1  x = np.array([1, 2, 3, 4, 5])
2  print(x >= 3)
```

[False False  True   True   True]

```
1  x = np.array([1, 2, 3, 4, 5])
2  print(x <= 3)
```

[ True   True   True False False]

```
1  x = np.array([1, 2, 3, 4, 5])
2  print(x != 3)
```

[ True   True False  True   True]

```
1  x = np.array([1, 2, 3, 4, 5])
2  print(x == 3)
```

[False False  True False False]

# Comparison Operators as ufuncs (3)

☐ Element-wise comparison of two arrays

```
1  x = np.array([1, 2, 3, 4, 5])
2  y = np.array([5, 4, 3, 2, 1])
3
4  print (x == y)
```

```
[False False  True False False]
```

☐ Use the comparison with other functions

```
1  ages = np.array([80,86,88,30])
2
3  print(ages >= ages.mean(0))
4
```

```
[ True  True  True False]
```

# Comparison Operators as ufuncs (3)

☐ Comparison Operator will work on arrays of any size and shape

```
1  ages = np.array([[80,86,88,30],
2                   [50,62,75,23],
3                   [66,58,40,36],
4                   [98,68,93,40]])
5
6  # mean of the mean
7  print(ages.mean(0).mean(0))
8  print(ages >= ages.mean(0).mean(0))
```

```
62.0625
[[ True  True  True False]
 [False False  True False]
 [ True False False False]
 [ True  True  True False]]
```

# Working with Boolean Arrays (1)

□ Counting entries with `np.sum()`

   – False is interpreted as 0, and True is
     interpreted as 1

```
1  ages = np.array([[80,86,88,30],
2                   [50,62,75,23],
3                   [66,58,40,36],
4                   [98,68,93,40]])
5
6  # mean of the mean
7  print(ages.mean(0).mean(0))
8  aboveAve = np.sum(ages >= ages.mean(0).mean(0))
9  print(aboveAve)
```

```
62.0625
8
```

---

# Working with Boolean Arrays (2)

□ Check whether any or all the values are
  true using `np.any()` or `np.all()`

```
1  ages = np.array([[80,86,88,30],
2                   [50,62,75,23],
3                   [66,58,40,36],
4                   [98,68,93,40]])
5  print(np.all(ages > 40))
6  print(np.all(ages > 20))
```

```
False
True
```

```
1  ages = np.array([[80,86,88,30],
2                   [50,62,75,23],
3                   [66,58,40,36],
4                   [98,68,93,40]])
5  print(np.any(ages > 90))
6  print(np.any(ages < 20))
```

```
True
False
```

❑Lab

– Counting Rainy Days
The program read a file containing a series of data that represents the amount of precipitation each day for a year in a given city.

UNLV

# NumPy (3)

Han-fen Hu

**UNLV**

---

# Outline

☐ Boolean Logic

☐ Masks

☐ Fancy Indexing

☐ Fast Sorting Arrays

☐ Partial Sorts

☐ Structured Arrays

**UNLV**

# Boolean Operators (1)

☐ Bitwise logic operators can be used together with the comparison operators

- & : and

- | : or

- ~ : not

- ^ : xor (exclusive or)
  - the result evaluates to True if only exactly *one* of the value is True.

| a | b | a ^ b |
|---|---|---|
| False | False | False |
| True | False | True |
| False | True | True |
| True | True | False |

UNLV

3

# Boolean Operators (2)

☐ Using the Keywords and/or Versus the Operators &/|

- and and or gauge the truth or falsehood of entire object,

- & and | refer to elements within each object

```
1  a = np.array([1, 0, 1, 0, 1, 0], dtype=bool)
2  b = np.array([1, 1, 1, 0, 1, 1], dtype=bool)
3  print(a | b)
```

[ True   True   True False   True   True]

```
1  a = np.array([1, 0, 1, 0, 1, 0], dtype=bool)
2  b = np.array([1, 1, 1, 0, 1, 1], dtype=bool)
3  print(a or b)
```

```
---------------------------------------------------------------------
ValueError                              Traceback (most recent call last)
<ipython-input-55-5f3e72e4d095> in <module>
      1 a = np.array([1, 0, 1, 0, 1, 0], dtype=bool)
      2 b = np.array([1, 1, 1, 0, 1, 1], dtype=bool)
----> 3 print(a or b)

ValueError: The truth value of an array with more than one element is ambiguous. Use a.any() or a.all()
```

4

# Comparisons and Masks (1)

❑ Boolean Arrays as Masks

– Use Boolean arrays as masks, to select particular subsets of the data themselves

❑ Masking comes up when you want to extract, modify, count, or otherwise manipulate values in an array based on some criterion

– Can be used to remove all outliers that are above some threshold

UNLV

# Comparisons and Masks (2)

❑ Masking operation

– Index on a Boolean array to filter the data

– Returns a one-dimensional array filled with all the values that meet this condition; i.e., all the values in positions at which the mask array is True.

UNLV

# Comparisons and Masks: Example (1)

```python
1  # generate integers (<10) of random number to fill a 3X4 array
2  x = np.random.randint(10, size=(3, 4))
3
4  print(x)
5  # a Boolean array showing whether the number is less than 5
6  print(x< 5)
7  # use the Boolean array to index the array x
8  print(x[x < 5])
```

```
[[3 0 7 7]
 [7 2 4 8]
 [5 9 3 0]]
[[ True  True False False]
 [False  True  True False]
 [False False  True  True]]
[3 0 2 4 3 0]
```

UNLV

# Comparisons and Masks: Example (2)

```python
1  ages = np.array([[80,86,88,30],
2                   [50,62,75,23],
3                   [66,58,40,36],
4                   [98,68,93,40]])
5  print(ages <= 30)
6
7  youngAges = ages[ages <= 30]
8  print(youngAges)
```

```
[[False False False  True]
 [False False False  True]
 [False False False False]
 [False False False False]]
[30 23]
```

UNLV

□Lab

– Counting Summer Rainy Days
The program read a file containing a series
of data that represents the amount of
precipitation each day for a year in a given
city.

Use a mask to show rainy days in the
summer.

UNLV

# Fancy Indexing

□Ways to access and modify portions of
arrays

– Index

– Slicing

– Boolean masks

– Fancy indexing

• Pass arrays of indices in place of single scalars

UNLV

# Fancy Indexing: Example (1)

```
1  import numpy as np
2
3  x = np.random.randint(10, size=10)
4  print(x)
5  ind = [3, 7, 4]
6  print(x[ind])
```

```
[6 5 8 4 3 8 3 0 6 4]
[4 0 3]
```

```
1  x = np.random.randint(10, size=10)
2  print(x)
3  ind = np.array([[3, 7],
4                  [4, 5]])
5  print(x[ind])
```

```
[2 7 9 3 3 3 5 6 7 0]
[[3 6]
 [3 3]]
```

When using fancy indexing, the shape of the result reflects the shape of the *index arrays*
rather than the shape of the *array being indexed*

# Fancy Indexing: Example (2)

```
1  x = np.array([[ 0,  1,  2,  3],
2                [ 4,  5,  6,  7],
3                [ 8,  9, 10, 11]])
4  row = np.array([0, 1, 2])
5  col = np.array([2, 1, 3])
6  print(x[row, col])
```

```
[ 2  5 11]
```

- Like with standard indexing, the first index refers to the row, and the second to the column
  - The first value in the result is x[0, 2],
  - The second is x[1, 1],
  - and the third is x[2, 3]

# Combined Indexing (1)

☐ Combine fancy and simple indices

  – What is the result?

```
1  x = np.array([[ 0,  1,  2,  3],
2                [ 4,  5,  6,  7],
3                [ 8,  9, 10, 11]])
4  print(x[2, [2, 0, 1]])
```

UNLV

---

# Combined Indexing (2)

☐ Combine fancy indexing with slicing

  – What is the result?

```
1  x = np.array([[ 0,  1,  2,  3],
2                [ 4,  5,  6,  7],
3                [ 8,  9, 10, 11]])
4  print(x[1:, [3, 1, 2]])
```

UNLV

# Combined Indexing (3)

☐ Combine fancy indexing with masking

– What is the result?

```
1  x = np.array([[ 0,  1,  2,  3],
2                [ 4,  5,  6,  7],
3                [ 8,  9, 10, 11]])
4  row = np.array([[0],
5                  [1]])
6  mask = np.array([1, 0, 1, 0], dtype=bool)
7
8  print(x[row, mask])
```

# Updating Values with Fancy Indexing (1)

☐ Just as fancy indexing can be used to access parts of an array, it can also be used to modify parts of an array

```
1  x = np.arange(10)
2  i = np.array([2, 1, 8, 4])
3  x[i] = 99
4  print(x)
```

```
[ 0 99 99  3 99  5  6  7 99  9]
```

# Updating Values with Fancy Indexing (2)

❑ Just as fancy indexing can be used to access parts of an array, it can also be used to modify parts of an array

```
1  x = np.arange(10)
2  i = np.array([2, 1, 8, 4])
3  # assign the value 99 to index 2, 1, 8, 4
4  x[i] = 99
5  print(x)
6  # subtract 10 from index 2, 1 8, 4
7  x[i] -= 10
8  print(x)
9
```

```
[ 0 99 99  3 99  5  6  7 99  9]
[ 0 89 89  3 89  5  6  7 89  9]
```

UNLV

---

# Fast Sorting in NumPy

❑ np.sort(): return a sorted array

❑ np.argsort(): return the *indices* of the sorted elements

```
1  x = np.random.randint(10, size=10)
2  print(x)
3  # print the sorted array
4  print(np.sort(x))
5  # print the indices of the osrted elements
6  print(np.argsort(x))
```

```
[1 1 3 2 2 5 6 1 0 4]
[0 1 1 1 2 2 3 4 5 6]
[8 0 1 7 3 4 2 9 5 6]
```

UNLV

# Sorting along Rows or Columns

☐ Adding the axis argument to the `sort()` function

```
1  s = np.random.randint(0, 10, (4, 6))
2  print(s)
3  print('sorted by row (within column)')
4  print(np.sort(s, axis=0))
5  print('sorted by column (within row)')
6  print(np.sort(s, axis=1))
```

```
[[6 0 4 2 1 5]
 [1 6 0 5 4 7]
 [7 2 9 2 7 5]
 [5 0 4 5 3 5]]
sorted by row (within column)
[[1 0 0 2 1 5]
 [5 0 4 2 3 5]
 [6 2 4 5 4 5]
 [7 6 9 5 7 7]]
sorted by column (within row)
[[0 1 2 4 5 6]
 [0 1 4 5 6 7]
 [2 2 5 7 7 9]
 [0 3 4 5 5 5]]
```

19

UNLV

---

# Question

☐ How do you sort an array by row and column?

```
1  s = np.random.randint(0, 10, (4, 6))
2  print(s)
3  print('sorted by row and column')
4  print(np.sort(np.sort(s, axis=0), axis=1))
5
```

20

UNLV

# Sort an Array in Descending Order

☐ Apply the `sort()` function and negative (-) ufunc

```
 1   # a random number array with 15 values between 0 and 99
 2
 3   s = np.random.randint(0, 100, (15))
 4   print(s)
 5
 6   print('transform the values to the opposite sign and then sort it')
 7   print(np.sort(-s))
 8
 9   print('transform the sorted values to the opposite sign')
10   print(-np.sort(-s))
11
12   print('Completed reverse sorting')
```

```
[95 53 23 71 84 43 98 25 52 95 61 95 15 83 80]
transform the values to the opposite sign and then sort it
[-98 -95 -95 -95 -84 -83 -80 -71 -61 -53 -52 -43 -25 -23 -15]
transform the sorted values to the opposite sign
[98 95 95 95 84 83 80 71 61 53 52 43 25 23 15]
Completed reverse sorting
```

---

# Partial Sorts: Partitioning (1)

☐ Sometimes we're not interested in sorting the entire array, but simply want to find the *k* smallest values in the array

☐ `np.partition()`

– Takes an array and a number *K*

– The result is a new array with the smallest *K* values to the left of the partition, and the remaining values to the right in arbitrary order

– Can partition along different axis of a multidimensional array

# Partial Sorts: Partitioning (2)

```
1  import numpy as np
2  # a random number array with 10 values between 0 and 99
3  s = np.random.randint(0, 100, (10))
4  print("Original array: ")
5  print(s)
6  # partition to get the smallest 3 number
7  p = np.partition(s, 3)
8  print("particall sorted array: ")
9  print(p)
```

```
Original array:
[90 37 54 65 50 25 78 39  1 83]
particall sorted array:
[37  1 25 39 50 54 78 65 83 90]
```

The first three values in the resulting array are the three smallest in the array.
The remaining array positions contain the remaining values

UNLV

# Partial Sorts: Partitioning (3)

```
1  # a 4x6 random number array with values between 0 and 99
2  twoDArray = np.random.randint(0, 100, (4,6))
3  print("Original array: ")
4  print(twoDArray)
5
6  # partition to get the smallest 3 number
7  # when axis is not provided, it is default to axis = 0
8  # sorted by row value
9  p = np.partition(twoDArray, 3, axis = 0)
10 print("particall sorted array (the 3 smallest value of each column on the top): ")
11 print(p)
12
13 # partition to get the smallest 3 number
14 # when axis  = 1, sorted by column value
15 p2 = np.partition(twoDArray, 3, axis = 1)
16 print("particall sorted array (the 3 smallest value of each row on the left): ")
17 print(p2)
```

```
Original array:
[[67 31 46  5 88 75]
 [28 50  6 22 82 22]
 [15 69 66 97 28 79]
 [54  9 17  9 20 64]]
particall sorted array (the 3 smallest value of each column on the top):
[[15  9  6  5 28 22]
 [28 31 17  9 20 64]
 [54 50 46 22 82 75]
 [67 69 66 97 88 79]]
particall sorted array (the 3 smallest value of each row on the left):
[[ 5 46 31 67 88 75]
 [22  6 22 28 50 82]
 [28 15 66 69 79 97]
 [ 9  9 17 20 54 64]]
```

axis = 1
the first three slots in each row contain the
smallest values from that row

UNLV

# Partial Sorts: Partitioning (4)

☐ `np.argpartition()`

– Similar to how `np.argsort()` works

– return the *indices* of the sorted elements

UNLV

# Structured Arrays (1)

☐ To store heterogeneous data

– As compared to `np.array()` that stores only data of the same type

☐ Arrays with compound data types

☐ Store associated data in the same array

UNLV

# Example: Regular Array for Associated Date

- Three arrays
- Use the index to reference values for the same individual

```
1  name = ['Alice', 'Bob', 'Cathy', 'Doug']
2  age = [25, 45, 37, 19]
3  salary = [55000.0, 85500.0, 68000.0, 61500.0]
4
5  print('First individual : ')
6  print(name[0]+", "+str(age[0])+", "+str(salary[0]))
7  print('\nAll: ')
8  for i in range(len(name)):
9      print(name[i]+", "+str(age[i])+", "+str(salary[i]))
```

```
First individual :
Alice, 25, 55000.0


All:
Alice, 25, 55000.0
Bob, 45, 85500.0
Cathy, 37, 68000.0
Doug, 19, 61500.0
```

UNLV

27

---

# Example: Structured Array for Associated Date

- Manage only one array
- Efficient: Arranged together in one convenient block of memory

```
1   name = ['Alice', 'Bob', 'Cathy', 'Doug']
2   age = [25, 45, 37, 19]
3   salary = [55000.0, 85500.0, 68000.0, 61500.0]
4
5   data = np.zeros(4, dtype={'names':('name', 'age', 'salary'),
6                             'formats':('U10', 'i4', 'f8')})
7
8   data['name'] = name
9   data['age'] = age
10  data['salary'] = salary
11  print('First individual : ')
12  print(data[0])
13  print('\nAll: ')
14  for person in data:
15      print(person)
```

```
First individual :
('Alice', 25, 55000.)

All:
('Alice', 25, 55000.)
('Bob', 45, 85500.)
('Cathy', 37, 68000.)
('Doug', 19, 61500.)
```

UNLV

28

# Structured Arrays (2)

☐ Data type used in defining the compound data type

    – U: Unicode string

    – i: integer

    – f: float

☐ Example

    – 'U10': Unicode string of maximum length 10

    – 'i4':  4-byte (i.e., 32 bit) integer

    – 'f8': 8-byte (i.e., 64 bit) float

**UNLV**

---

# Structured Arrays (2)

☐ Individual attributes are still available

☐ We can refer to values either by index or by name

```
1   name = ['Alice', 'Bob', 'Cathy', 'Doug']
2   age = [25, 45, 37, 19]
3   salary = [55000.0, 85500.0, 68000.0, 61500.0]
4
5   data = np.zeros(4, dtype={'names':('name', 'age', 'salary'),
6                             'formats':('U10', 'i4', 'f8')})
7   data['name'] = name
8   data['age'] = age
9   data['salary'] = salary
10  print('All names ')
11  print(data['name'])
12  print('\nName of last person: ')
13  print(data[-1]['name'])
```

```
All names
['Alice' 'Bob' 'Cathy' 'Doug']

Name of last person:
Doug
```

**UNLV**

# Structured Arrays (3)

◻ Boolean masking can be applied to filter data

```python
1  name = ['Alice', 'Bob', 'Cathy', 'Doug']
2  age = [25, 45, 37, 19]
3  salary = [55000.0, 85500.0, 68000.0, 61500.0]
4
5  data = np.zeros(4, dtype={'names':('name', 'age', 'salary'),
6                            'formats':('U10', 'i4', 'f8')})
7  data['name'] = name
8  data['age'] = age
9  data['salary'] = salary
10 # Get names where salary is less then 65000
11 print(data[data['salary'] < 65000]['name'])
```

['Alice' 'Doug']

# Define Structured Array Data Types

◻ `np.dtype()`

◻ Method 1:

– Provide the *names* and the *formats*

```python
1  employeeDType = np.dtype({'names':('name', 'age', 'salary'),
2                            'formats':('U10', 'i4', 'f8')})
3  print(employeeDType)
```

[('name', '<U10'), ('age', '<i4'), ('salary', '<f8')]

◻ Method 2

– Each attribute is specified as a tuple

```python
1  employeeDType = np.dtype([('name', 'U10'), ('age', 'i4'), ('salary', 'f8')])
2
3  print(employeeDType)
```

[('name', '<U10'), ('age', '<i4'), ('salary', '<f8')]

□Lab

– Loyal Customer Lookup
This program read the customer data from a file. The user can enter a criterion to filter customer data by the loyalty points earned.

UNLV

# Pandas (1)

Han-fen Hu

UNLV

---

# Outline

☐ Introduction to Pandas

☐ Series Objects

☐ DataFrame Object

☐ Indexing and Selection

☐ Operating on Data in Pandas

UNLV

# Introduction to Pandas

❑ Pandas is an open-source library providing high-performance, easy-to-use data structures and data analysis tools for the Python

❑ Designed to make working with "relational" or "labeled" data easy and intuitive

❑ https://pandas.pydata.org/

3

**UNLV**

# Pandas Objects

❑ Pandas objects are enhanced versions of NumPy structured arrays

– Rows and columns are identified with labels

❑ Three fundamental Pandas data structures

– Series: one-dimensional array of indexed data

– DataFrame: a generalization of a NumPy array, or as a specialization of a Python dictionary

– Index: an *immutable array* or as an *ordered set*

4

**UNLV**

# Pandas Series Object (1)

◻ One-dimensional array of indexed data

– `Series` wraps both a sequence of values and a sequence of indices, which we can access with the values and index attributes

```
1  import pandas as pd
2
3  data = pd.Series([0.25, 0.5, 0.75, 1.0])
4  print("The series object: ")
5  print(data)
6  print("\nThe values only: ")
7  print(data.values)
8  print("\nThe indices: ")
9  print(data.index)
```

```
The series object:
0    0.25
1    0.50
2    0.75
3    1.00
dtype: float64

The values only:
[0.25 0.5  0.75 1.  ]

The indices:
RangeIndex(start=0, stop=4, step=1)
```

UNLV

---

# Pandas Series Object (2)

◻ Data can be accessed by the associated index via the familiar Python square-bracket notation

– While the Numpy Array has an implicitly defined integer index used to access the values, the Pandas Series has an explicitly defined index associated with the values

```
1  import pandas as pd
2
3  data = pd.Series([0.25, 0.5, 0.75, 1.0])
4  print("The second value: ")
5  print(data[1])
6  print("The second to third value: ")
7  print(data[1:3])
```

```
The second value:
0.5
The second to third value:
1    0.50
2    0.75
dtype: float64
```

UNLV

# Pandas Series Object (3)

☐ The index need not be an integer, but can consist of values of any desired type

```
1  data = pd.Series([0.25, 0.5, 0.75, 1.0],
2                   index=['a', 'b', 'c', 'd'])
3  print(data)
4
5  print('\nThe value for index b: ')
6  print(data['b'])
```

```
a    0.25
b    0.50
c    0.75
d    1.00
dtype: float64

The value for index b:
0.5
```

UNLV

---

# Pandas Series Object (4)

☐ Series can be considered as specialized dictionary

- A dictionary is a structure that maps arbitrary keys to a set of arbitrary values

- A Series is a structure which maps typed keys to a set of typed values

  - Makes it much more efficient than Python dictionaries for certain operations

- Unlike a dictionary, though, the Series also supports array-style operations such as slicing

UNLV

# Slicing of Series Object: Example

```
1  population_dict = {'California': 38332521,
2                     'Texas': 26448193,
3                     'New York': 19651127,
4                     'Florida': 19552860,
5                     'Illinois': 12882135}
6  population = pd.Series(population_dict)
7
8  print(population['California'])
9  print(population['California':'New York'])
```

```
38332521
California    38332521
Texas        26448193
New York     19651127
dtype: int64
```

9

---

# Creating Series Objects

□ pd.Series(data, index=index)

```
1  # index defaults to an integer sequence
2  # data can be a list or array
3  data1= pd.Series([2, 4, 6])
4  print(data1)
```

```
0    2
1    4
2    6
dtype: int64
```

```
1  # data can be a scalar, which is repeated to fill the specified index
2  data2 = pd.Series(5, index=[100, 200, 300])
3  print(data2)
```

```
100    5
200    5
300    5
dtype: int64
```

```
1  # data can be a dictionary
2  # in which index defaults to the sorted dictionary keys
3  data3 = pd.Series({2:'a', 1:'b', 3:'c'})
4  print(data3)
```

```
2    a
1    b
3    c
dtype: object
```

10

# Pandas DataFrame Object (1)

❑ A generalization of a NumPy array

- A two-dimensional array with both flexible row indices and flexible column names

- A sequence of Series objects that share the same index.

```
1  population_dict = {'California': 38332521,'Texas': 26448193,'New York': 19651127,
2                     'Florida': 19552860,'Illinois': 12882135}
3  population = pd.Series(population_dict)
4  area_dict = {'California': 423967, 'Texas': 695662, 'New York': 141297,
5               'Florida': 170312, 'Illinois': 149995}
6  area = pd.Series(area_dict)
7  states = pd.DataFrame({'population': population,
8                         'area': area})
9  print(states)
```

```
            population    area
California    38332521  423967
Texas         26448193  695662
New York      19651127  141297
Florida       19552860  170312
Illinois      12882135  149995
```

# Pandas DataFrame Object (2)

❑ DataFrame has an index attribute that gives access to the index labels

❑ DataFrame has a columns attribute, which is an Index object holding the column labels

```
1  population_dict = {'California': 38332521,'Texas': 26448193,'New York': 19651127,
2                     'Florida': 19552860,'Illinois': 12882135}
3  population = pd.Series(population_dict)
4  area_dict = {'California': 423967, 'Texas': 695662, 'New York': 141297,
5               'Florida': 170312, 'Illinois': 149995}
6  area = pd.Series(area_dict)
7  states = pd.DataFrame({'population': population,
8                         'area': area})
9  print(states.index)
10 print(states.columns)
```

```
Index(['California', 'Texas', 'New York', 'Florida', 'Illinois'], dtype='object')
Index(['population', 'area'], dtype='object')
```

# Pandas DataFrame Object (3)

❑ Can be considered as a generalization of a two-dimensional NumPy array

– A DataFrame has labels for the columns

❑ Can also think of a DataFrame as a specialization of a dictionary

– a DataFrame maps a column name to a Series of column data

UNLV

---

# Creating DataFrame Objects (1)

❑ Method 1: From a single Series object

```
1  population_dict = {'California': 38332521,'Texas': 26448193,'New York': 19651127,
2                     'Florida': 19552860,'Illinois': 12882135}
3  population = pd.Series(population_dict)
4
5  s = pd.DataFrame(population, columns=['population'])
6  print(s)
```

```
            population
California    38332521
Texas         26448193
New York      19651127
Florida       19552860
Illinois      12882135
```

UNLV

# Creating DataFrame Objects (2)

❑ Method 2: From a list of dictionaries

– Even if some keys in the dictionary are missing, Pandas will fill them in with NaN (i.e., "not a number") values

```
1  s = pd.DataFrame([{'a': 1, 'b': 2}, {'b': 3, 'c': 4}])
2  print(s)
```

```
     a  b    c
0  1.0  2  NaN
1  NaN  3  4.0
```

❑ Method 3: From a dictionary of Series objects

– See the example on slide 11&12

# Creating DataFrame Objects (3)

❑ Method 4: From a two-dimensional NumPy array

```
1  s = pd.DataFrame(np.random.rand(3, 2), # a 3x2 array of random numbers
2                   columns=['foo', 'bar'], # column labels
3                   index=['a', 'b', 'c']) # index
4  print(s)
```

```
        foo       bar
a  0.479797  0.440195
b  0.272726  0.658412
c  0.943578  0.228896
```

❑ Method 5: From a NumPy structured array

```
1  a = np.zeros(3, dtype=[('A', 'i8'), ('B', 'f8')])
2  s = pd.DataFrame(a)
3  print(s)
```

```
   A    B
0  0  0.0
1  0  0.0
2  0  0.0
```

# Read Data from File into DataFrame (1)

❑ The first step to any data science project is to import your data

❑ `read_csv()` function

  – File path is the argument
    • full file path which is prefixed by a / and includes the working directory in the specification,
    • or use the relative file path which doesn't.
  – The data will be read into a DataFrame

UNLV

# Read Data from File into DataFrame (2)

❑ The index will be automatically assigned.

❑ `set_index()`

  – Set the index to an existing column

UNLV

□ Lab

– Demographic Statistics
The program read the data from a CSV file
into a DataFrame and show the statistics

UNLV

# Data Selection in Series (1)

□ Series object acts in many ways like a
standard Python dictionary

– Use the key to access the value

– `in` and `not in` operator

– Assign new value

– Add a new key-value pair

UNLV

# Data Selection in Series (2)

```
1  data = pd.Series([0.25, 0.5, 0.75, 1.0],
2                   index=['a', 'b', 'c', 'd'])
3  print(data['b']) # use key to access the value
4  if 'a' in data: # the in operator
5      print('a is in data')
6
7  # update the value of b
8  data['b']=1.5
9
10 # assign a new key-value pair
11 data['e'] = 1.25
12
13 print(data)
```

```
0.5
a is in data
a    0.25
b    1.50
c    0.75
d    1.00
e    1.25
dtype: float64
```

**UNLV**

---

# Data Selection in Series (3)

□ Series object acts in many ways like a one-dimensional NumPy array

– Slices

- slicing with an explicit index (i.e., data['a':'c']), the final index is included in the slice
- slicing with an implicit index (i.e., data[0:2]), the final index is excluded from the slice

– Masking

– Fancy indexing

**UNLV**

# Data Selection in Series (4)

```
1  data = pd.Series([0.25, 0.5, 0.75, 1.0],
2                   index=['a', 'b', 'c', 'd'])
3  # slicing by explicit index
4  print(data['a':'c'])
5  # slicing by implicit integer index
6  print(data[0:2])
7  # masking
8  print(data[(data > 0.3) & (data < 0.8)])
9  # fancy indexing
10 print(data[['a', 'e']])
```

```
a    0.25
b    0.50
c    0.75
dtype: float64
a    0.25
b    0.50
dtype: float64
b    0.50
c    0.75
dtype: float64
a    0.25
e     NaN
dtype: float64
```

# Data Selection in DataFrame (1)

☐ DataFrame acts in many ways like a dictionary of Series structures sharing the same index

– Use the column label to retrieve the values

– Use the dictionary-style syntax to modify the object (e.g., add a new column)

# Data Selection in DataFrame (2)

```python
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})

# get the value for the key "area"
print(data['area'])

# add a new key-value pair.
# the value here is a Series, calculated based on pop and area
data['density'] = data['pop'] / data['area']

print(data)
```

```
California     423967
Texas          695662
New York       141297
Florida        170312
Illinois       149995
Name: area, dtype: int64
                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
```

25

**UNLV**

---

# Data Selection in DataFrame (3)

☐ DataFrame acts in many ways like a two-dimensional or structured array

– For array-style indexing, we need another convention: `loc`, `iloc`, and `ix`

• `loc:` index the data using the explicit index and column names

• `iloc:` index the array using the implicit Python-style index

– With `loc` and `iloc`,

• Can combine masking and fancy indexing

• Can update the value

26

**UNLV**

# Data Selection in DataFrame (4)

```
                area        pop
California    423967   38332521
Texas         695662   26448193
New York      141297   19651127
Florida       170312   19552860
Illinois      149995   12882135
```

```python
1  area = pd.Series({'California': 423967, 'Texas': 695662,
2                    'New York': 141297, 'Florida': 170312,
3                    'Illinois': 149995})
4  pop = pd.Series({'California': 38332521, 'Texas': 26448193,
5                   'New York': 19651127, 'Florida': 19552860,
6                   'Illinois': 12882135})
7  data = pd.DataFrame({'area':area, 'pop':pop})
8
9  # use the explicit index and column names
10 print(data.loc[:'Illinois', :'pop'])
11 # use the implicit Python-style index
12 print(data.iloc[:3, :2])
13
```

```
                area        pop
California    423967   38332521
Texas         695662   26448193
New York      141297   19651127
Florida       170312   19552860
Illinois      149995   12882135
                area        pop
California    423967   38332521
Texas         695662   26448193
New York      141297   19651127
```

27

# Data Selection in DataFrame (5)

```
                --
                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
```

```python
1  area = pd.Series({'California': 423967, 'Texas': 695662,
2                    'New York': 141297, 'Florida': 170312,
3                    'Illinois': 149995})
4  pop = pd.Series({'California': 38332521, 'Texas': 26448193,
5                   'New York': 19651127, 'Florida': 19552860,
6                   'Illinois': 12882135})
7  data = pd.DataFrame({'area':area, 'pop':pop})
8  data['density'] = data['pop'] / data['area']
9
10 # use the masking and fancy indexing
11 print(data.loc[data.density > 100, ['pop', 'density']])
12
```

```
                pop      density
New York   19651127   139.076746
Florida    19552860   114.806121
```

28

# Data Selection in DataFrame (6)

```
             area        pop
California  423967  38332521
Texas       695662  26448193
New York    141297  19651127
Florida     170312  19552860
Illinois    149995  12882135
```

```python
1  area = pd.Series({'California': 423967, 'Texas': 695662,
2                    'New York': 141297, 'Florida': 170312,
3                    'Illinois': 149995})
4  pop = pd.Series({'California': 38332521, 'Texas': 26448193,
5                   'New York': 19651127, 'Florida': 19552860,
6                   'Illinois': 12882135})
7  data = pd.DataFrame({'area':area, 'pop':pop})
8
9  data.iloc[0, 1] = 90
10 print(data)
11
```

```
             area        pop
California  423967        90
Texas       695662  26448193
New York    141297  19651127
Florida     170312  19552860
Illinois    149995  12882135
```

UNLV

29

---

# Data Selection in DataFrame (7)

☐ Note that *slicing* and *mask* refers to rows

```python
1  area = pd.Series({'California': 423967, 'Texas': 695662,
2                    'New York': 141297, 'Florida': 170312,
3                    'Illinois': 149995})
4  pop = pd.Series({'California': 38332521, 'Texas': 26448193,
5                   'New York': 19651127, 'Florida': 19552860,
6                   'Illinois': 12882135})
7  data = pd.DataFrame({'area':area, 'pop':pop})
8  data['density'] = data['pop'] / data['area']
9
10 print('slicing refers to rows. Using key')
11 print(data['Florida':'Illinois'])
12 print('\nslicing refers to rows. implicit index')
13 print(data[1:3])
14 print('\nMasking operations refers to rows')
15 print(data[data.density > 100])
```

```
slicing refers to rows. Using key
            area        pop     density
Florida    170312  19552860  114.806121
Illinois   149995  12882135   85.883763

slicing refers to rows. implicit index
            area        pop     density
Texas      695662  26448193   38.018740
New York   141297  19651127  139.076746

Masking operations refers to rows
            area        pop     density
New York   141297  19651127  139.076746
Florida    170312  19552860  114.806121
```

UNLV

30

□Lab

– Demographic Statistics (cont'd)
The program read the data from a CSV file
into a DataFrame and show the statistics

UNLV

# Operating on Data in Pandas (1)

□Pandas inherits much of this functionality
from NumPy, and the ufuncs

– Pandas is designed to work with NumPy,
any NumPy ufunc will work on Pandas
Series and DataFrame objects

UNLV

# Operating on Data in Pandas (2)

- For binary operations on two Series or DataFrame objects, Pandas will align indices in the process of performing the operation
  - The resulting array contains the *union* of indices of the two input arrays
  - Any item for which one or the other does not have an entry is marked with NaN, or "Not a Number,"

```python
1  area = pd.Series({'Alaska': 1723337, 'Texas': 695662,
2                    'California': 423967}, name='area')
3  population = pd.Series({'California': 38332521, 'Texas': 26448193,
4                    'New York': 19651127}, name='population')
5
6  density = population / area
7  print(density)
```

```
Alaska            NaN
California     90.413926
New York          NaN
Texas          38.018740
dtype: float64
```

33

# Operating on Data in Pandas (3)

```python
1  a = pd.Series([2, 4, 6], index=[0, 1, 2])
2  b = pd.Series([1, 3, 5], index=[1, 2, 3])
3
4  print(a+b)
```

```
0    NaN
1    5.0
2    9.0
3    NaN
dtype: float64
```

```python
1  a = pd.Series([2, 4, 6], index=[0, 1, 2])
2  b = pd.Series([1, 3, 5], index=[1, 2, 3])
3
4  # for NaN, fill in with 0
5  c = a.add(b, fill_value=0)
6  print(c)
```

```
0    2.0
1    5.0
2    9.0
3    5.0
dtype: float64
```

34

# Operating on Data in Pandas (4)

| Python Operator | Pandas Method(s) |
|---|---|
| + | add() |
| - | sub(), subtract() |
| * | mul(), multiply() |
| / | truediv(), div(), divide() |
| // | floordiv() |
| % | mod() |
| ** | pow() |

**UNLV**

# Operating on Data in Pandas (5)

- A similar type of alignment takes place for both columns and indices for operations on DataFrame

```
1  x = pd.DataFrame(np.random.randint(0, 20, (2, 2)),
2                   columns=list('AB'))
3  print("Content of x")
4  print(x)
5  y = pd.DataFrame(np.random.randint(0, 10, (3, 3)),
6                   columns=list('BAC'))
7  print("\nContent of y")
8  print(y)
9  print("\nContent of x+y")
10 print(x+y)
```

```
Content of x
    A   B
0  19  18
1   9   4

Content of y
   B  A  C
0  4  8  0
1  2  0  2
2  6  7  2

Content of x+y
      A     B    C
0  27.0  22.0  NaN
1   9.0   6.0  NaN
2   NaN   NaN  NaN
```

**UNLV**

# ❑Lab

– Demographic Statistics (cont'd)
The program read the data from a CSV file
into a DataFrame and show the statistics

UNLV

# Pandas (2)

UNLV

---

# Outline

- ☐ Missing Data

- ☐ Operating on Null Values

- ☐ MultiIndex

UNLV

# Missing Data (1)

☐ Real-world data is rarely clean and homogeneous

☐ Missing Data Conventions

– Using a *mask* that globally indicates missing values

– Choosing a *sentinel value* that indicates a missing entry

• Such as -9999 or some rare pattern.

UNLV

# Missing Data (2)

☐ In Pandas

– NaN (acronym for Not a Number) for missing numeric values

– None object for others missing values

☐ Note: an Empty string ("") is not equivalent to NaN or None

UNLV

# NaN: Missing Numerical Data (1)

☐ Acronym for Not a Number

☐ Special floating-point value

– Can be used to represent missing floating-point value

```
1  import numpy as np # import numpy
2  import pandas as pd # import pandas
3
4  # declare a numeric array, the second element is a missing value
5  exampNumArray = np.array([1, np.nan, 3, 4])
6  # print the array
7  print(exampNumArray)
8  # print the data type of the array
9  print(exampNumArray.dtype)
```

```
[ 1. nan  3.  4.]
float64
```

**UNLV**

---

# NaN: Missing Numerical Data (2)

☐ Data Virus: Regardless of the operation, the result of arithmetic with NaN will be another NaN

– Aggregates over the values are well defined (i.e., they don't result in an error) but not always useful:

```
1   import numpy as np
2   import pandas as pd
3
4   vals = np.array([1, np.nan, 3, 4])
5
6
7   print(vals.sum())
8   print(vals.max())
9   print(vals.min())
10  print(1+ vals)
11  print(0* vals)
12
```

```
nan
nan
nan
[ 2. nan  4.  5.]
[ 0. nan  0.  0.]
```

**UNLV**

# NaN: Missing Numerical Data (3)

☐ NumPy does provide some special aggregations that will ignore these missing values

- np.nansum()

- np.nanmin()

- np.nanmax()

```python
1   import numpy as np
2   import pandas as pd
3
4   vals = np.array([1, np.nan, 3, 4])
5
6
7   print(np.nansum(vals))
8   print(np.nanmin(vals))
9   print(np.nanmax(vals))
10
```

```
8.0
1.0
4.0
```

UNLV

---

# NaN and None

☐ NaN: Missing floating-point values

☐ None: Missing object value

- For example, string

☐ Pandas is built to handle the two of them nearly interchangeably

| Type | Conversion | Storing Value |
|------|-----------|---------------|
| floating | No change | np.nan |
| object/string | No change | None or np.nan |
| integer | Cast to float64 | np.nan |
| boolean | Cast to object | None or np.nan |

UNLV

# Operating on Null Values (1)

☐ Several useful methods for detecting, removing, and replacing null values in Pandas data structures

- `isnull()`: Generate a Boolean mask indicating missing values

- `notnull()`: Opposite of isnull()

- `dropna()`: Return a filtered version of the data

- `fillna()`: Return a copy of the data with missing values filled or imputed

UNLV

---

# Operating on Null Values (1): Example

```
1  data = pd.Series([1, np.nan, 'hello', None])
2  # function that returns Boolean mask
3  print(data.isnull())

0    False
1     True
2    False
3     True
dtype: bool
```

```
1  data = pd.Series([1, np.nan, 'hello', None])
2  # function that returns Boolean mask
3  print(data.notnull())

0     True
1    False
2     True
3    False
dtype: bool
```

```
1  # use the mask to filter data
2  print(data[data.notnull()])
3  # It is is equivalent to the following
4  print(data.dropna())

0        1
2    hello
dtype: object
0        1
2    hello
dtype: object
```

Detecting null values
The boolean array can be used for masking

Dropping null values

UNLV

# Operating on Null Values (2)

☐ **For a DataFrame, there are more options**

- We cannot drop single values from a DataFrame; we can only drop full rows or full columns

- By default, dropna() will drop all rows in which any null value is present

- Alternatively, we can drop NA values along a different axis

---

# Operating on Null Values (2): Example

```
1  df = pd.DataFrame([[1,      np.nan, 2],
2                     [2,      3,      5],
3                     [np.nan, 4,      6]])
4
5  print(df.dropna())
```

```
     0    1  2
1  2.0  3.0  5
```

Default action
Drop the rows containing null value

```
1  print(df.dropna(axis='columns'))
```

```
   2
0  2
1  5
2  6
```

Drop the columns containing null value

```
1  print(df.dropna(axis='rows'))
```

```
     0    1  2
1  2.0  3.0  5
```

Explicitly specification
Drop the rows containing null value

# Operating on Null Values (3)

❑You might rather be interested in dropping rows or columns with *all* NA values, or a majority of NA values

— how parameters

- `how = 'all'`: drop rows/columns that are *all* null values

- `how = 'any'`: drop rows/columns that contain *any* null values

— `thresh` parameter

- Specify a minimum number of non-null values for the row/column to be kept

UNLV

---

# Operating on Null Values (3): Example

```
1  df = pd.DataFrame([[1,      np.nan, 2, np.nan],
2                     [2,      3,      5, np.nan],
3                     [np.nan, 4,      6, np.nan],
4                     [np.nan, np.nan, 7, np.nan]])
5  # drop the columns with all na values
6  print(df.dropna(axis='columns', how='all'))
```

```
     0    1  2
0  1.0  NaN  2
1  2.0  3.0  5
2  NaN  4.0  6
3  NaN  NaN  7
```

Drop the column
if all the values in the column is null

```
1  # drop the columns with any na values
2  print(df.dropna(axis='columns', how='any'))
```

```
   2
0  2
1  5
2  6
3  7
```

Drop the column
if any value in the column is null

```
1  # drop the rows with less than 2 valid values
2  print(df.dropna(axis='rows', thresh=2))
```

```
     0    1  2    3
0  1.0  NaN  2  NaN
1  2.0  3.0  5  NaN
2  NaN  4.0  6  NaN
```

Drop the row
if it contains more than 2 non-null values

UNLV

❑ Lab

– Medical Tracking Data

– The program reads a csv file containing the data for a medical experiment. Several thousands of people received a treatment and came back for 37 tests. The test results were recorded. There are missing data in the file.

UNLV

# Filling Null Values

❑ Sometimes rather than dropping NA values, we'd rather replace them with a valid value.

– `fillna(value):` fill NA entries with a single value

– `fillna(method='ffill'):` forward-fill to propagate the previous value forward

– `fillna(method='bfill'):` back-fill to propagate the next values backward

UNLV

# Filling Null Values: Example

```
1  data = pd.Series([1, np.nan, 2,np.nan, 3],
2                   index=list('abcde'))
3  print(data.fillna(0))
```

```
a    1.0
b    0.0
c    2.0
d    0.0
e    3.0
dtype: float64
```

```
1  print(data.fillna(method='ffill'))
2  print(data.fillna(method='bfill'))
```

```
a    1.0
b    1.0
c    2.0
d    2.0
e    3.0
dtype: float64
a    1.0
b    2.0
c    2.0
d    3.0
e    3.0
dtype: float64
```

```
1  df = pd.DataFrame([[1,      np.nan, 2, np.nan],
2                     [2,      3,      5, np.nan],
3                     [np.nan, 4,      6, np.nan],
4                     [np.nan, np.nan, 7, np.nan]])
5  print(df.fillna(method='ffill', axis="columns"))
```

```
     0    1    2    3
0  1.0  1.0  2.0  2.0
1  2.0  3.0  5.0  5.0
2  NaN  4.0  6.0  6.0
3  NaN  NaN  7.0  7.0
```

The fills take place by taking
the value of previous columns

Note that if there isn't a previous column,
the null value remain unfilled.

---

☐Lab

– Vaccination Data
The program read a csv file containing
number of people received vaccinations (in
thousands).
There are missing data in the file.

□Exercise

– BMI for US Presidents
Write a program that reads
president_heights_updated.csv, and shows
the names of US presidents with the highest
and lowest BMI from the dataset.

UNLV

# MultiIndex (1)

□Datasets are not limited to one-
dimensional and two-dimensional

□MultiIndex data type contains multiple
levels of indexing

```
California  2000    33871648
            2010    37253956
New York    2000    18976457
            2010    19378102
Texas       2000    20851820
            2010    25145561
```

UNLV

# MultiIndex (1): Example

```python
1  import pandas as pd
2  import numpy as np
3
4  # stats and year
5  ind = [('California', 2000), ('California', 2010),
6         ('New York', 2000), ('New York', 2010),
7         ('Texas', 2000), ('Texas', 2010)]
8  # population for the states
9  populations = [33871648, 37253956,
10               18976457, 19378102,
11               20851820, 25145561]
12
13 # convert the index to a MultiIndex
14 ind = pd.MultiIndex.from_tuples(ind)
15 # use the MultiIndex as the index for the data
16 pop = pd.Series(populations, index=ind)
17
18 pop
```

```
California  2000    33871648
            2010    37253956
New York    2000    18976457
            2010    19378102
Texas       2000    20851820
            2010    25145561
dtype: int64
```

UNLV

---

# MultiIndex (2)

◻ We can add dimensions to the data frame

```python
1  import pandas as pd
2  import numpy as np
3
4  # stats and year
5  ind = [('California', 2000), ('California', 2010),
6         ('New York', 2000), ('New York', 2010),
7         ('Texas', 2000), ('Texas', 2010)]
8  # population for the states
9  populations = [33871648, 37253956,
10               18976457, 19378102,
11               20851820, 25145561]
12 # population for people under 18
13 minor= [9267089, 9284094,
14         4687374, 4318033,
15         5906301, 6879014]
16
17 # convert the index to a MultiIndex
18 ind = pd.MultiIndex.from_tuples(ind)
19 # use the MultiIndex as the index for the data
20 pop = pd.Series(populations, index=ind)
21 print(pop)
22 pop_df = pd.DataFrame({'total': pop,
23                        'under18': minor})
24 print(pop_df)
```

```
California  2000    33871648
            2010    37253956
New York    2000    18976457
            2010    19378102
Texas       2000    20851820
            2010    25145561
dtype: int64
                   total    under18
California  2000  33871648  9267089
            2010  37253956  9284094
New York    2000  18976457  4687374
            2010  19378102  4318033
Texas       2000  20851820  5906301
            2010  25145561  6879014
```

UNLV

# MultiIndex (3)

☐ all the ufuncs and other functionality for
data frame work with hierarchical indices
as well.

```
24  # apply ufunc for calculation
25  pop_df['ratio'] = pop_df['under18'] / pop_df['total']
26  print(pop_df)
```

```
                      total   under18     ratio
California 2000   33871648   9267089  0.273594
           2010   37253956   9284094  0.249211
New York   2000   18976457   4687374  0.247010
           2010   19378102   4318033  0.222831
Texas      2000   20851820   5906301  0.283251
           2010   25145561   6879014  0.273568
```

UNLV

23

# Pandas (3)

UNLV

---

# Outline

☐ **Combining Datasets**

- Concat

- Merge

☐ **Data Manipulation**

- Drop

- Replace

☐ **High-Performance Pandas: query()**

UNLV

# Combining Datasets

☐ Some of the most interesting studies of data come from combining different data sources

☐ Pandas includes functions and methods that make this sort of data wrangling fast and straightforward.

UNLV

---

# Concatenation with `pd.concat` (1)

☐ `pd.concat()`: By default, the concatenation takes place row-wise but allows specification of an axis

```python
1  df1=pd.DataFrame([['A1', 'B1'], ['A2', 'B2']], columns=list('AB'))
2  df2=pd.DataFrame([['C3', 'D3'], ['C4', 'D4']], columns=list('AB'))
3
4  print(df1)
5  print(df2)
6  print(pd.concat([df1, df2]))
7  print(pd.concat([df1, df2], axis=1))
```

```
     A    B
0   A1   B1
1   A2   B2
     A    B
0   C3   D3
1   C4   D4
     A    B
0   A1   B1
1   A2   B2
0   C3   D3
1   C4   D4
     A    B    A    B
0   A1   B1   C3   D3
1   A2   B2   C4   D4
```

UNLV

# Concatenation with `pd.concat` (2)

☐ `pd.concat()`: preserves indices, even if the result will have duplicate indices

– To ignore the index, use the `ignore_index` flag

```
1  df1=pd.DataFrame([['A1', 'B1'], ['A2', 'B2']], columns=list('AB'))
2  df2=pd.DataFrame([['C3', 'D3'], ['C4', 'D4']], columns=list('AB'))
3
4  print(df1)
5  print(df2)
6  print(pd.concat([df1, df2],ignore_index=True))
```

```
     A    B
0   A1   B1
1   A2   B2
     A    B
0   C3   D3
1   C4   D4
     A    B
0   A1   B1
1   A2   B2
2   C3   D3
3   C4   D4
```

UNLV

---

# Concatenation with Joins (1)

☐ In practice, data from different sources might have different sets of column names

– By default, the entries for which no data is available are filled with NaN values

```
1   df3=pd.DataFrame([['A1','B1','C1'], ['A2','B2','C2']]
2                  , columns=list('ABC')
3                  , index = list('01'))
4   df4=pd.DataFrame([['B3','C3','D3'], ['B4','C4','D4']]
5                  , columns=list('BCD')
6                  , index = list('34'))
7
8   print(df3)
9   print(df4)
10  print(pd.concat([df3, df4]))
```

```
     A    B    C
0   A1   B1   C1
1   A2   B2   C2
     B    C    D
3   B3   C3   D3
4   B4   C4   D4
     A    B    C     D
0   A1   B1   C1   NaN
1   A2   B2   C2   NaN
3  NaN   B3   C3    D3
4  NaN   B4   C4    D4
```

UNLV

# Concatenation with Joins (2)

◻ We can specify the options of `join` and `join_axes` parameters of the concatenate function.

– By default, the join is a union of the input columns (`join='outer'`)

– We can change this to an intersection of the columns using `join='inner'`

– We can also specify the index of the remaining columns using the `join_axes` argument

UNLV

---

# Concatenation with Joins: Example (1)

Get the intersection of the columns using `join='inner'`

```
 1  df3=pd.DataFrame([['A1','B1','C1'], ['A2','B2','C2']]
 2                  , columns=list('ABC')
 3                  , index = list('01'))
 4  df4=pd.DataFrame([['B3','C3','D3'], ['B4','C4','D4']]
 5                  , columns=list('BCD')
 6                  , index = list('34'))
 7
 8  print(df3)
 9  print(df4)
10  print(pd.concat([df3, df4], join='inner'))
```

```
    A   B   C
0  A1  B1  C1
1  A2  B2  C2
    B   C   D
3  B3  C3  D3
4  B4  C4  D4
    B   C
0  B1  C1
1  B2  C2
3  B3  C3
4  B4  C4
```

UNLV

# Concatenation with Joins: Example (2)

Specify the index of the remaining columns
using `join_axes`

```
1  df3=pd.DataFrame([['A1','B1','C1'], ['A2','B2','C2']]
2                   , columns=list('ABC')
3                   , index = list('01'))
4  df4=pd.DataFrame([['B3','C3','D3'], ['B4','C4','D4']]
5                   , columns=list('BCD')
6                   , index = list('34'))
7
8  print(df3)
9  print(df4)
10 print(pd.concat([df3, df4], join_axes=[df3.columns]))
```

```
     A    B    C
0   A1   B1   C1
1   A2   B2   C2
     B    C    D
3   B3   C3   D3
4   B4   C4   D4
     A    B    C
0   A1   B1   C1
1   A2   B2   C2
3  NaN   B3   C3
4  NaN   B4   C4
```

---

# Merge: One-to-One

❑ Merge two DataFrame objects

– Recognize the common column and use the column as a key to merge

– After merge, the order of entries in each column is not necessarily maintained

– The merge in general discards the index

# One-to-One Merge: Example

```
1  emp1 = pd.DataFrame({'employee': ['Bob', 'Jake', 'Lisa', 'Sue'],
2                       'group': ['Accounting', 'Engineering', 'Engineering', 'HR']})
3  emp2 = pd.DataFrame({'employee': ['Lisa', 'Bob', 'Jake', 'Sue'],
4                       'hire_date': [2004, 2008, 2012, 2014]})
5  print(emp1)
6  print(emp2)
7  print(pd.merge(emp1, emp2))
```

```
  employee        group
0      Bob   Accounting
1     Jake  Engineering
2     Lisa  Engineering
3      Sue           HR
  employee  hire_date
0     Lisa       2004
1      Bob       2008
2     Jake       2012
3      Sue       2014
  employee        group  hire_date
0      Bob   Accounting       2008
1     Jake  Engineering       2012
2     Lisa  Engineering       2004
3      Sue           HR       2014
```

UNLV

# Merge: Many-to-One

☐ Many-to-one merge happens when one of the two key columns contains duplicate entries.

– The resulting DataFrame will preserve those duplicate entries as appropriate.

UNLV

# Many-to-One Merge: Example

```
1  emp1 = pd.DataFrame({'employee': ['Bob', 'Jake', 'Lisa', 'Sue'],
2                       'group': ['Accounting', 'Engineering', 'Engineering', 'HR']})
3  emp2 = pd.DataFrame({'employee': ['Lisa', 'Bob', 'Jake', 'Sue'],
4                       'hire_date': [2004, 2008, 2012, 2014]})
5  emp3 = pd.merge(emp1, emp2)
6  emp4 = pd.DataFrame({'group': ['Accounting', 'Engineering', 'HR'],
7                       'supervisor': ['Carly', 'Guido', 'Steve']})
8  print(emp3)
9  print(emp4)
10 print(pd.merge(emp3, emp4))
```

```
  employee        group  hire_date
0      Bob   Accounting       2008
1     Jake  Engineering       2012
2     Lisa  Engineering       2004
3      Sue           HR       2014
        group supervisor
0  Accounting      Carly
1 Engineering      Guido
2          HR      Steve
  employee        group  hire_date supervisor
0      Bob   Accounting       2008      Carly
1     Jake  Engineering       2012      Guido
2     Lisa  Engineering       2004      Guido
3      Sue           HR       2014      Steve
```

"supervisor" information is repeated as needed

UNLV

13

---

# Merge: Many-to-Many

- If the key column in both the left and right DataFrame contains duplicates, then the result is a many-to-many merge

  – Some values are repeated as needed

UNLV

14

# Many-to-Many Merge: Example

```python
1  emp1 = pd.DataFrame({'employee': ['Bob', 'Jake', 'Lisa', 'Sue'],
2                       'group': ['Accounting', 'Engineering', 'Engineering', 'HR']})
3  emp5 = pd.DataFrame({'group': ['Accounting', 'Accounting',
4                       'Engineering', 'Engineering', 'HR', 'HR'],
5                       'skills': ['math', 'spreadsheets', 'coding', 'linux',
6                       'spreadsheets', 'organization']})
7  print(emp1)
8  print(emp5)
9  print(pd.merge(emp1, emp5))
```

```
  employee        group
0      Bob   Accounting
1     Jake  Engineering
2     Lisa  Engineering
3      Sue           HR
        group        skills
0  Accounting          math
1  Accounting  spreadsheets
2 Engineering        coding
3 Engineering         linux
4          HR  spreadsheets
5          HR  organization
  employee        group        skills
0      Bob   Accounting          math
1      Bob   Accounting  spreadsheets
2     Jake  Engineering        coding
3     Jake  Engineering         linux
4     Lisa  Engineering        coding
5     Lisa  Engineering         linux
6      Sue           HR  spreadsheets
7      Sue           HR  organization
```

15

UNLV

---

# Keyword on

☐ on: specify the name of the key column for merge

```python
1  emp1 = pd.DataFrame({'employee': ['Bob', 'Jake', 'Lisa', 'Sue'],
2                       'group': ['Accounting', 'Engineering', 'Engineering', 'HR']})
3  emp2 = pd.DataFrame({'employee': ['Lisa', 'Bob', 'Jake', 'Sue'],
4                       'hire_date': [2004, 2008, 2012, 2014]})
5  print(emp1)
6  print(emp2)
7  print(pd.merge(df1, df2, on='employee'))
```

```
  employee        group
0      Bob   Accounting
1     Jake  Engineering
2     Lisa  Engineering
3      Sue           HR
  employee  hire_date
0     Lisa       2004
1      Bob       2008
2     Jake       2012
3      Sue       2014
  employee        group  hire_date
0      Bob   Accounting       2008
1     Jake  Engineering       2012
2     Lisa  Engineering       2004
3      Sue           HR       2014
```

16

UNLV

# Keyword `left_on` and `right_on`

❑ `left_on` and `right_on`: merge two datasets with different column names

```
1  emp1 = pd.DataFrame({'employee': ['Bob', 'Jake', 'Lisa', 'Sue'],
2                       'group': ['Accounting', 'Engineering', 'Engineering', 'HR']})
3  emp3 = pd.DataFrame({'name': ['Bob', 'Jake', 'Lisa', 'Sue'],
4                       'salary': [70000, 80000, 120000, 90000]})
5  print(emp1)
6  print(emp3)
7  print(pd.merge(emp1, emp3, left_on="employee", right_on="name"))
```

```
  employee         group
0      Bob    Accounting
1     Jake   Engineering
2     Lisa   Engineering
3      Sue            HR
   name  salary
0   Bob   70000
1  Jake   80000
2  Lisa  120000
3   Sue   90000
  employee        group  name  salary
0      Bob   Accounting   Bob   70000
1     Jake  Engineering  Jake   80000
2     Lisa  Engineering  Lisa  120000
3      Sue           HR   Sue   90000
```

UNLV

---

# Keyword `left_index` and `right_index`

❑ Rather than merging on a column, we can merge on an index

```
1  emp1 = pd.DataFrame({'employee': ['Bob', 'Jake', 'Lisa', 'Sue'],
2                       'group': ['Accounting', 'Engineering', 'Engineering', 'HR']})
3  emp2 = pd.DataFrame({'employee': ['Lisa', 'Bob', 'Jake', 'Sue'],
4                       'hire_date': [2004, 2008, 2012, 2014]})
5  emp1 = emp1.set_index('employee')
6  emp2 = emp2.set_index('employee')
7  print(emp1)
8  print(emp2)
9  print(pd.merge(emp1, emp2, left_index=True, right_index=True))
```

```
                group
employee
Bob        Accounting
Jake      Engineering
Lisa      Engineering
Sue                HR
          hire_date
employee
Lisa           2004
Bob            2008
Jake           2012
Sue            2014
                group  hire_date
employee
Bob        Accounting       2008
Jake      Engineering       2012
Lisa      Engineering       2004
Sue                HR       2014
```

UNLV

# drop() method (1)

☐ Drop a column from the DataFrame

```
1  emp1 = pd.DataFrame({'employee': ['Bob', 'Jake', 'Lisa', 'Sue'],
2                       'group': ['Accounting', 'Engineering', 'Engineering', 'HR']})
3  emp3 = pd.DataFrame({'name': ['Bob', 'Jake', 'Lisa', 'Sue'],
4                       'salary': [70000, 80000, 120000, 90000]})
5  emp6 = pd.merge(emp1, emp3, left_on="employee", right_on="name").drop('name', axis="columns")
6  print(emp6)
```

```
  employee        group  salary
0      Bob   Accounting   70000
1     Jake  Engineering   80000
2     Lisa  Engineering  120000
3      Sue           HR   90000
```

---

# Inner and Outer Merge (1)

☐ When a value appears in one key column but not the other, we need to consider how to merge the data

```
1  guest1 = pd.DataFrame({'name': ['Peter', 'Paul', 'Mary'],
2                         'food': ['fish', 'beans', 'bread']},
3                        columns=['name', 'food'])
4  guest2 = pd.DataFrame({'name': ['Mary', 'Joseph'],
5                         'drink': ['wine', 'beer']},
6                        columns=['name', 'drink'])
7  print(guest1)
8  print(guest2)
9  print(pd.merge(guest1, guest2))
```

```
    name   food
0  Peter   fish
1   Paul  beans
2   Mary  bread
     name drink
0    Mary  wine
1  Joseph  beer
   name   food drink
0  Mary  bread  wine
```

*By default, the result contains the* intersection *of the two sets of inputs; this is what is known as an* inner join

# Inner and Outer Merge (2)

❑By using the keyword how, we can define outer, `left`, and `right` join

```
 1  guest1 = pd.DataFrame({'name': ['Peter', 'Paul', 'Mary'],
 2                         'food': ['fish', 'beans', 'bread']},
 3                         columns=['name', 'food']})
 4  guest2 = pd.DataFrame({'name': ['Mary', 'Joseph'],
 5                         'drink': ['wine', 'beer']},
 6                         columns=['name', 'drink']})
 7
 8  print(pd.merge(guest1, guest2, how = 'outer'))
 9  print(pd.merge(guest1, guest2, how = 'left'))
10  print(pd.merge(guest1, guest2, how = 'right'))
```

```
      name    food
0    Peter    fish
1     Paul   beans
2     Mary   bread
      name drink
0     Mary   wine
1   Joseph   beer
```

```
      name    food drink
0    Peter    fish   NaN
1     Paul   beans   NaN
2     Mary   bread  wine
3   Joseph    NaN   beer
      name    food drink
0    Peter    fish   NaN
1     Paul   beans   NaN
2     Mary   bread  wine
      name    food drink
0     Mary   bread  wine
1   Joseph    NaN   beer
```

21

**UNLV**

---

❑Lab

– US States Data
The program reads data from three csv files, representing data from different sources.
In the program, the data will be combined/merged into a single DataFrame.
The program will rank US states and territories by their 2010 population density.

22

**UNLV**

# drop() method: Remove Rows (1)

❑ `drop()` method can also be used to drop a row from a DataFrame

```
1  emp = pd.DataFrame({'employee':['Bob', 'Jake', 'Lisa','Sue','Ann','John'],
2                      'group':['Accounting', 'Engineering','Engineering','HR','RD','RD'],
3                      'salary':[70000, 80000, 120000, 90000,85000,85000]})
4
5  emp = emp.set_index('employee')
6  emp1 = emp.drop("Bob")
7  emp1
```

| employee | group | salary |
|---|---|---|
| Jake | Engineering | 80000 |
| Lisa | Engineering | 120000 |
| Sue | HR | 90000 |
| Ann | RD | 85000 |
| John | RD | 85000 |

| | employee | group | salary |
|---|---|---|---|
| 0 | Bob | Accounting | 70000 |
| 1 | Jake | Engineering | 80000 |
| 2 | Lisa | Engineering | 120000 |
| 3 | Sue | HR | 90000 |
| 4 | Ann | RD | 85000 |
| 5 | John | RD | 85000 |

23

---

# drop() method: Remove Rows (2)

❑ `drop()` method can be used to drop multiple rows

– Use fancy indexing (explicit index)

```
1  emp = pd.DataFrame({'employee':['Bob', 'Jake', 'Lisa','Sue','Ann','John'],
2                      'group':['Accounting', 'Engineering','Engineering','HR','RD','RD'],
3                      'salary':[70000, 80000, 120000, 90000,85000,85000]})
4
5  emp = emp.set_index('employee')
6  emp2 = emp.drop(["Bob", "Ann"])
7  emp2
```

| employee | group | salary |
|---|---|---|
| Jake | Engineering | 80000 |
| Lisa | Engineering | 120000 |
| Sue | HR | 90000 |
| John | RD | 85000 |

UNLV

24

# drop() method: Remove Rows (3)

□ `drop()` method can be used to drop multiple rows

– Use fancy indexing (implicit index)

```
 1  emp = pd.DataFrame({'employee':['Bob', 'Jake', 'Lisa','Sue','Ann','John'],
 2                      'group':['Accounting', 'Engineering','Engineering','HR','RD','RD'],
 3                      'salary':[70000, 80000, 120000, 90000,85000,85000]})
 4
 5  emp = emp.set_index('employee')
 6  emp3 = emp.drop(emp.index[[1, 4]])
 7  print(emp3)
 8
 9  emp4 = emp.drop(emp.index[:3])
10  print(emp4)
```

```
              group  salary
employee
Bob         Accounting   70000
Lisa        Engineering  120000
Sue                 HR    90000
John                RD    85000
              group  salary
employee
Sue          HR    90000
Ann          RD    85000
John         RD    85000
```

UNLV

25

---

# Other Ways to Drop Rows from a DataFrame

□ Use masking (on the value of a column)

```
 1  emp = pd.DataFrame({'employee':['Bob', 'Jake', 'Lisa','Sue','Ann','John'],
 2                      'group':['Accounting', 'Engineering','Engineering','HR','RD','RD'],
 3                      'salary':[70000, 80000, 120000, 90000,85000,85000]})
 4
 5
 6  emp5 = emp.loc[emp['employee'] != 'Bob']
 7  print(emp5)
 8  print('\n\n')
 9
10  emp6 = emp.loc[emp['salary'] < 100000]
11  print(emp6)
```

```
   employee       group  salary
1      Jake  Engineering   80000
2      Lisa  Engineering  120000
3       Sue          HR    90000
4       Ann          RD    85000
5      John          RD    85000



   employee       group  salary
0       Bob   Accounting   70000
1      Jake  Engineering   80000
3       Sue          HR    90000
4       Ann          RD    85000
5      John          RD    85000
```

UNLV

26

# Replace Values in a DataFrame

☐ `replace()` method can be used to update values in a DataFrame

```
1  emp = pd.DataFrame({'employee':['Bob', 'Jake', 'Lisa','Sue','Ann','John'],
2                      'group':['Accounting', 'Engineering','Engineering','HR','RD','RD'],
3                      'salary':[70000, 80000, 120000, 90000,85000,85000]})
4
5
6  emp = emp.replace('RD','R&D')
7  emp
```

|   | employee | group | salary |
|---|----------|-------|--------|
| 0 | Bob | Accounting | 70000 |
| 1 | Jake | Engineering | 80000 |
| 2 | Lisa | Engineering | 120000 |
| 3 | Sue | HR | 90000 |
| 4 | Ann | R&D | 85000 |
| 5 | John | R&D | 85000 |

UNLV

---

# Replace NaN Values

☐ For the whole DtaFrame

```
1  data = pd.read_csv('bonus.csv')
2
3  data = data.replace(np.nan, -99999)
```

☐ For one column

```
7  emp['bonus'] = emp['bonus'].replace(np.nan, 0)
8  emp
```

|   | employee | group | salary | bonus |
|---|----------|-------|--------|-------|
| 0 | Bob | Accounting | 70000 | NaN |
| 1 | Jake | Engineering | 80000 | 3000.0 |
| 2 | Lisa | Engineering | 120000 | 2000.0 |
| 3 | Sue | HR | 90000 | 15000.0 |
| 4 | Ann | RD | 85000 | NaN |
| 5 | John | RD | 0 | NaN |

|   | employee | group | salary | bonus |
|---|----------|-------|--------|-------|
| 0 | Bob | Accounting | 70000 | 0.0 |
| 1 | Jake | Engineering | 80000 | 3000.0 |
| 2 | Lisa | Engineering | 120000 | 2000.0 |
| 3 | Sue | HR | 90000 | 15000.0 |
| 4 | Ann | RD | 85000 | 0.0 |
| 5 | John | RD | 0 | 0.0 |

# Update a Cell in a DataFrame (1)

❑Specify the index and column

```
1  customer = pd.DataFrame({'name': ['Alan','Byona','Catherine','Dean','Franky'],
2              'age': [30, 69, 40, 18, 22],
3              'premium':[345,234,974,563,435]})
4
5  # update the first customer's age
6  customer.at[0,'age'] = 31
7  customer
```

|   | name | age | premium |
|---|------|-----|---------|
| 0 | Alan | 31 | 345 |
| 1 | Byona | 69 | 234 |
| 2 | Catherine | 40 | 974 |
| 3 | Dean | 18 | 563 |
| 4 | Franky | 22 | 435 |

UNLV

# Update a Cell in a DataFrame (2)

❑Specify the condition and column

```
1  customer = pd.DataFrame({'name': ['Alan','Byona','Catherine','Dean','Franky'],
2              'age': [30, 69, 40, 18, 22],
3              'premium':[345,234,974,563,435]})
4
5  # update Alan's age
6  customer.at[customer['name'] == 'Alan','age'] = 31
7  customer
```

|   | name | age | premium |
|---|------|-----|---------|
| 0 | Alan | 31 | 345 |
| 1 | Byona | 69 | 234 |
| 2 | Catherine | 40 | 974 |
| 3 | Dean | 18 | 563 |
| 4 | Franky | 22 | 435 |

UNLV

□ Lab

  – Player Salary Data
    The program reads data from a csv files,
    representing data of NBA players' salaries in
    different years. This program can show the
    salary for a certain year, and allows the user
    to update the salary.

UNLV

---

# DataFrame.sort_value()

□ sort_value(by=column,
  asending=Boolean)

```
1  df = pd.DataFrame({'col1': ['A', 'A', 'B', np.nan, 'D', 'C'],
2                     'col2': [2, 1, 9, 8, 7, 4],
3                     'col3': [0, 1, 9, 4, 2, 3]})
4  print(df)
5  print(df.sort_values(by='col1'))
```

```
   col1  col2  col3
0    A     2     0
1    A     1     1
2    B     9     9
3   NaN    8     4
4    D     7     2
5    C     4     3
   col1  col2  col3
0    A     2     0
1    A     1     1
2    B     9     9
5    C     4     3
4    D     7     2
3   NaN    8     4
```

UNLV

# `DataFrame.sort_value()`: Example

```
     col1  col2  col3
0     A     2     0
1     A     1     1
2     B     9     9
3   NaN     8     4
4     D     7     2
5     C     4     3
```

```
1  df = pd.DataFrame({'col1': ['A', 'A', 'B', np.nan, 'D', 'C'],
2                     'col2': [2, 1, 9, 8, 7, 4],
3                     'col3': [0, 1, 9, 4, 2, 3]})
4  print(df.sort_values(by=['col1', 'col2']))
```

```
     col1  col2  col3
1     A     1     1
0     A     2     0
2     B     9     9
5     C     4     3
4     D     7     2
3   NaN     8     4
```

**Sort by two columns**

```
1  df = pd.DataFrame({'col1': ['A', 'A', 'B', np.nan, 'D', 'C'],
2                     'col2': [2, 1, 9, 8, 7, 4],
3                     'col3': [0, 1, 9, 4, 2, 3]})
4  print(df.sort_values('col1', ascending=False))
```

```
     col1  col2  col3
4     D     7     2
5     C     4     3
2     B     9     9
0     A     2     0
1     A     1     1
3   NaN     8     4
```

**Sort by col1, in descending order**

33

UNLV

---

# `DataFrame.query()` Method

□ Can be used to filter data with conditions on multiple columns

– A more efficient computation

– Compared to the masking expression this is much easier to read and understand

34

UNLV

❑Lab

– US States Data (cont.)
The program reads data from three csv files,
representing data from different sources.
In the program, the data will be
combined/merged into a single DataFrame.
The program will rank US states and
territories by their 2010 population density.

UNLV

# Visualization with Matplotlib (1)

UNLV

---

# Outline

❑ Introduction to Matplotlib

❑ Simple Line Plots

❑ Simple Scatter Plots

❑ Histograms, Binnings, and Density

❑ Customizing Plot Legends and Colorbars

UNLV

# Introduction to Matplotlib (1)

- ☐ Data visualization library built on NumPy arrays
  - Allows visual access to huge amounts of data in easily digestible visuals
  - Large user base and an active developer base
  - Predated Pandas by more than a decade, and thus is not designed for use with Pandas DataFrames
- ☐ Seaborn
  - Provides an API on top of Matplotlib that offers choices for plot style and color defaults
  - Defines simple high-level functions for common statistical plot types
  - Integrates with the functionality provided by Pandas DataFrames

UNLV

# Introduction to Matplotlib (2)

- ☐ Importing Matplotlib

```
1  # import matplotlib, set the alias as mpl
2  import matplotlib as mpl
3  # import the pyplot module, set the alias as plt
4  import matplotlib.pyplot as plt
```

  - Pyplot is the most used module of Matplotlib
    - Provides an interface like MATLAB but instead, it uses Python and it is open source
- ☐ IPython is built to work well with Matplotlib

UNLV

# Introduction to Matplotlib (3)

☐ `plt.show()`

– If we run the .py file from the shell, the `show()` function is needed to opens a window that display the figure

```
1   # import numpy, set the alias as np
2   import numpy as np
3   # import the pyplot module, set the alias as plt
4   import matplotlib.pyplot as plt
5
6   # an array of 100 numbers evenly distributed between 0 an
7   x = np.linspace(0, 10, 100)
8
9   plt.plot(x, np.sin(x))
10  plt.plot(x, np.cos(x))
11
12  plt.show()
```

# Introduction to Matplotlib (4)

☐ IPython is built to work well with Matplotlib if we specify Matplotlib mode

– `%matplotlib inline` will lead to static images of the plot embedded in the notebook

– Creating a plot will embed a PNG image of the resulting graphic

– It needs to be done only once per kernel/session

UNLV

# Introduction to Matplotlib (5)

```
1   # import numpy, set the alias as np
2   import numpy as np
3   # import the pyplot module, set the alias as plt
4   import matplotlib.pyplot as plt
5
6   # an array of 100 numbers evenly distributed between 0 and 100
7   x = np.linspace(0, 10, 100)
8
9   %matplotlib inline
10  plt.plot(x, np.sin(x))
11  plt.plot(x, np.cos(x))
```

[<matplotlib.lines.Line2D at 0x1763436eac8>]

7

---

# Simple Line Plots (1)

❑ For all Matplotlib plots, we start by creating a figure and an axes

```
1   %matplotlib inline
2   import matplotlib.pyplot as plt
3   import numpy as np
4
5   # assign the figure object to a variable
6   fig = plt.figure()
7   # assign the axes object to a variable
8   ax = plt.axes()
```



figure is a single container that contains all the objects representing axes, graphics, text, and labels.

axes is a bounding box with ticks and labels, which will eventually contain the plot elements that make up our visualization

8

# Simple Line Plots (1)

☐ Use the `plot()` function to draw the plot

```
1   # import numpy, set the alias as np
2   import numpy as np
3   # import the pyplot module, set the alias as plt
4   import matplotlib.pyplot as plt
5
6   %matplotlib inline
7
8   # an array of 100 numbers evenly distributed between 0 and 100
9   a = np.linspace(0, 10, 100)
10  # an array of exponential values of -a
11  b = np.exp(-a)
12
13  # use a as x axis and b as y axis for a line plot
14  plt.plot(a, b)
```

UNLV

---

# Simple Line Plots (3)

☐ To create a single figure with multiple lines, just simply call the plot function multiple times

```
1   # import numpy, set the alias as np
2   import numpy as np
3   # import the pyplot module, set the alias as plt
4   import matplotlib.pyplot as plt
5
6   # an array of 100 numbers evenly distributed between 0 and 100
7   a = np.linspace(0, 10, 100)
8   # an array of exponential values of -a
9   b = np.exp(-a)
10  c = np.exp(-a*a)
11  # use a as x axis and b as y axis for a line plot
12  plt.plot(a, b)
13  # call plot() again to add another line
14  plt.plot(a, c)
```

UNLV

# Adjusting the Plot: Axes Limits

☐ To adjust axis limits is to use the `plt.xlim()` and `plt.ylim()` methods

```python
1  # import numpy, set the alias as np
2  import numpy as np
3  # import the pyplot module, set the alias as plt
4  import matplotlib.pyplot as plt
5
6  # an array of 100 numbers evenly distributed between 0 and 100
7  x = np.linspace(0, 10, 100)
8
9  plt.plot(x, np.sin(x))
10
11 plt.xlim(-1, 11) # set the x axis to -1 and 11
12 plt.ylim(-1.5, 1.5) # set the y axis to -1.5 and 1.5
```



11

# Line Colors and Styles (1)

☐ `color` keyword: accepts a string argument representing virtually any imaginable color

– If no color is specified, Matplotlib will automatically cycle through a set of default colors for multiple lines

```python
1  x = np.linspace(0, 10, 1000)
2  plt.plot(x, np.sin(x - 0), color='blue')        # specify color by name
3  plt.plot(x, np.sin(x - 1), color='g')           # short color code (rgbcmyk)
4  plt.plot(x, np.sin(x - 2), color='0.75')        # Grayscale between 0 and 1
5  plt.plot(x, np.sin(x - 3), color='#FFDD44')     # Hex code (RRGGBB from 00 to FF)
6  plt.plot(x, np.sin(x - 4), color=(1.0,0.2,0.3)) # RGB tuple, values 0 to 1
7  plt.plot(x, np.sin(x - 5), color='chartreuse'); # all HTML color names supported
```



12

UNLV

# Line Colors and Styles (2)

☐ `linestyle` keyword: Specify the line style

```python
1  x = np.linspace(0, 10, 1000)
2  plt.plot(x, x + 0, linestyle='solid')
3  plt.plot(x, x + 1, linestyle='dashed')
4  plt.plot(x, x + 2, linestyle='dashdot')
5  plt.plot(x, x + 3, linestyle='dotted');
6
7  # For short, you can use the following codes:
8  plt.plot(x, x + 4, linestyle='-')   # solid
9  plt.plot(x, x + 5, linestyle='--')  # dashed
10 plt.plot(x, x + 6, linestyle='-.')  # dashdot
11 plt.plot(x, x + 7, linestyle=':');  # dotted
```



13

# Line Colors and Styles (3)

☐ `linestyle` and `color` codes can be combined into a single non-keyword argument

```python
1  x = np.linspace(0, 10, 1000)
2  plt.plot(x, x + 0, '-g')  # solid green
3  plt.plot(x, x + 1, '--c') # dashed cyan
4  plt.plot(x, x + 2, '-.k') # dashdot black
5  plt.plot(x, x + 3, ':r');  # dotted red
```



14

# Labeling Plots

□ `title()`, `xlabel()`, and `ylabel()`

  – Set the title and axis labels

□ `legend()` and the keyword `label` in `plot()`

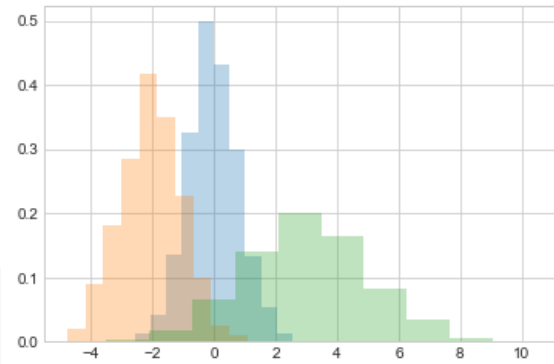  – Set the plot legend

UNLV

---

# Labeling Plots: Example



```
1   %matplotlib inline
2   # import numpy, set the alias as np
3   import numpy as np
4   # import the pyplot module, set the alias as
5   import matplotlib.pyplot as plt
6
7   # an array of 100 numbers evenly distributed
8   x = np.linspace(0, 10, 100)
9
10  # set the title of the plot
11  plt.title("Sine and Cosine Curves")
12  # set the label of x axis
13  plt.xlabel("x")
14  # set the label of y axis
15  plt.ylabel("sin(x) and cos(x)");
16
17  # set the plot legend for each line on the plot
18  plt.plot(x, np.sin(x), label='sin(x)')
19  plt.plot(x, np.cos(x), label='cos(x)')
20  plt.legend() # show the legend
```

UNLV

# Simple Scatter Plots

☐ `scatter()` function

```
1   %matplotlib inline
2   # import numpy, set the alias as np
3   import numpy as np
4   # import the pyplot module, set the alias as plt
5   import matplotlib.pyplot as plt
6
7   # an array of 30 numbers evenly distributed between 0 and 100
8   x = np.linspace(0, 10, 30)
9   # an array of 30 random integers between 0 and 9
10  y = np.random.randint(0, 10, 30)
11
12  plt.scatter(x, y)
13
```

# Scatter Plots with Color and Size

☐ c keyword: assign different colors to the dots

  – array or list of colors or color

☐ s keyword: assign different size to the dot

```
1   x = np.random.randn(100)
2   y = np.random.randn(100)
3   colors = np.random.rand(100)
4   sizes = 1000 * colors
5
6   # c is defined by the numbers in colors
7   # s is defined by the numbers in size
8   # cmap specifies the color map
9   plt.scatter(x, y, c=colors, s=sizes, alpha=0.3,
10             cmap='viridis')
11  plt.colorbar()  # show color scale
```

UNLV

# Colormaps

☐ Matplotlib has a number of built-in colormaps

– They are  accessible via at
  `matplotlib.colormaps`

```python
1  import matplotlib.pyplot as plt
2
3
4  plt.colormaps()
```
```
['Accent',
 'Accent_r',
 'Blues',
 'Blues_r',
 'BrBG',
 'BrBG_r',
 'BuGn',
 'BuGn_r',
 'BuPu',
 'BuPu_r',
 'CMRmap',
 'CMRmap_r',
 'Dark2',
 'Dark2_r',
 'GnBu',
```

UNLV

---

☐ Lab

– Height and Weight by Age Group
  This program reads the data from a csv file
  and then plot the relationships between
  height and weight

UNLV

# Histograms

☐ A simple histogram can be a great first step in understanding a dataset

☐ `hist()` function

- `bins` keyword: specify the number of bins

- `histtype` keyword:
  - 'bar' is a traditional bar-type histogram. If multiple data are given the bars are arranged side by side.
  - 'barstacked' is a bar-type histogram where multiple data are stacked on top of each other.
  - 'step' generates a lineplot that is by default unfilled.
  - 'stepfilled' generates a lineplot that is by default filled.

- `alpha` keyword: set the opacity

- `density` keyword: The area (or integral) under the histogram will sum to 1

UNLV

---

# Histograms: Example (1)

```
1  # import numpy, set the alias as np
2  import numpy as np
3  # import the pyplot module, set the alias as plt
4  import matplotlib.pyplot as plt
5  %matplotlib inline
6  # 1000 random numbers with mean=0, sd = 0.8
7  x1 = np.random.normal(0, 0.8, 1000)
8  plt.hist(x1, histtype='stepfilled', bins=10)
9  plt.show()
```



UNLV

# Histograms: Example (2)



```python
1   # import numpy, set the alias as np
2   import numpy as np
3   # import the pyplot module, set the alias as
4   import matplotlib.pyplot as plt
5   %matplotlib inline
6
7   # three random number arrays, with 1000 numbres each
8   x1 = np.random.normal(0, 0.8, 1000)
9   x2 = np.random.normal(-2, 1, 1000)
10  x3 = np.random.normal(3, 2, 1000)
11
12  # a shared set of keywords
13  kwargs = dict(histtype='stepfilled', alpha=0.3, density=True, bins=10)
14
15  # draw the three subsets of data, with the same set of keywords
16  plt.hist(x1, **kwargs)
17  plt.hist(x2, **kwargs)
18  plt.hist(x3, **kwargs)
```

**UNLV**

---

# Two-Dimensional Histograms

☐ We can also create histograms in two-dimensions by dividing points among two-dimensional bins

– Create a heat map of the data

☐ `hist2d()` function

– `bins` keyword: specify the number of bins for the two dimensions

– `cmap` keyword: color theme

**UNLV**

# Two-Dimensional Histograms: Example

```
21  # draw the 2-d scatter plot
22  # 30 bins on each dimension. Color map as blue
23  plt.hist2d(weight, height, bins=30, cmap='Blues')
24  # show the colorar on the side with labels
25  plt.colorbar().set_label('counts in bin')
```

---

# Customizing Plot Style

☐ `plt.style.use()`

– Specify the style you would like to apply

– You can use `plt.style.available` to see all the available styles

  • ['seaborn-ticks', 'ggplot', 'dark_background', 'bmh', 'seaborn-poster', 'seaborn-notebook', 'fast', 'seaborn', 'classic', 'Solarize_Light2', 'seaborn-dark', 'seaborn-pastel', 'seaborn-muted', '_classic_test', 'seaborn-paper', 'seaborn-colorblind', 'seaborn-bright', 'seaborn-talk', 'seaborn-dark-palette', 'tableau-colorblind10', 'seaborn-darkgrid', 'seaborn-whitegrid', 'fivethirtyeight', 'grayscale', 'seaborn-white', 'seaborn-deep']

# Customizing Plot Style: Example

```
6  plt.style.use('seaborn-muted')
7
8  x = np.linspace(0, 10, 1000)
9  # Get the fiure and axes from the plot
10 fig, ax = plt.subplots()
11 # draw the two lines with respective style and label
12 ax.plot(x, np.sin(x), '-b', label='Sine')
13 ax.plot(x, np.cos(x), '--r', label='Cosine')
14
15 # adjust plots with equal axis ratios
16 ax.axis('equal')
17 # show the legend
18 ax.legend()
```

<matplotlib.legend.Legend at 0x1763a39ba08>

```
6  plt.style.use('ggplot')
7
8  x = np.linspace(0, 10, 1000)
9  # Get the fiure and axes from the plot
10 fig, ax = plt.subplots()
11 # draw the two lines with respective style and labels
12 ax.plot(x, np.sin(x), '-b', label='Sine')
13 ax.plot(x, np.cos(x), '--r', label='Cosine')
14
15 # adjust plots with equal axis ratios
16 ax.axis('equal')
17 # show the legend
18 ax.legend()
```

<matplotlib.legend.Legend at 0x276d1c53b48>

---

# Customizing Plot Legends

- ☐ `loc` keyword: specify the location

- ☐ `frameon` keyword: turn on or off the frame

- ☐ `ncol` keyword: specify the number of columns

- ☐ `fancybox` keyword: use a rounded box or not

- ☐ `shadow` keyword: add a shadow

# Customizing Plot Legends: Example

```
 6  plt.style.use('seaborn-white')
 7
 8  x = np.linspace(0, 10, 1000)
 9  # Get the fiure and axes from the plot
10  fig, ax = plt.subplots()
11  # draw the two lines with respective style and labels
12  ax.plot(x, np.sin(x), '-b', label='Sine')
13  ax.plot(x, np.cos(x), '--r', label='Cosine')
14
15  # adjust plots with equal axis ratios
16  ax.axis('equal')
17  # show the legend to use fancybox, turn the frame on, add shadow
18  # make the location as lower center, two columns
19  ax.legend(fancybox=True, frameon=True, shadow=True, loc='lower center', ncol=2)
```

<matplotlib.legend.Legend at 0x276d2fd0f08>

UNLV

---

## ☐Lab

– California Cities
The program reads data from a csv file and plot the California cites. The size of the dots represents the area, and the color shows the population

UNLV

□Exercise

– Height Weight & BodyFat
Please use the data bodyData.csv to
visualize the height, weight, and percentage
body fat data. For example, create a figure
as shown below

# Visualization with Matplotlib (2)

UNLV

---

# Outline

- ☐ Pie Chart

- ☐ Multiple Subplots

- ☐ Three-Dimensional Plotting

- ☐ Visualization with Seaborn

UNLV

# Pie Chart

☐ `pie()` function to display the pie chart

    – Values and labels are lists

    – `autopct`: auto-labeling the percentage

        • Formatting string

    – explode: offsetting a slice

---

```python
import matplotlib.pyplot as plt

# Labels a list
categories = ['Utilities', 'Food', 'Entertainment', 'Clothing', 'Misc.']
# Values to be charted as a list
amounts = [312,658, 230, 498, 123]

# offsetting a slice; only "explode" the 3rd slice
offset = (0, 0, 0.1, 0, 0)

# draw the pie chart, showing percentage, offset the "Entertainment slice"
plt.pie(amounts, labels = categories, autopct ='%1.2f%%', explode=offset)
# add title to the chart
plt.title('Expenses')
plt.show()
```

# Multiple Subplots

□ Subplots

 – Compare different views of data side by side

 – Groups of smaller axes that can exist together within a single figure

□ Different ways to create subplots

 – Subplots by Hand

 – Grids of Subplots

UNLV

---

# Subplots by Hand (1)

□ `plt.axes ()` function

 – by default this creates a standard axes object that fills the entire figure.

 – It also takes a four numbers in the figure coordinate system to represent [left, bottom, width, height]

```
1   # import numpy, set the alias as np
2   import numpy as np
3   # import the pyplot module, set the alias as plt
4   import matplotlib.pyplot as plt
5
6   # standard axes
7   ax1 = plt.axes()
8   # create an inset axes by setting the x and y position to 0.65
9   # and the x and y extents to 0.2
10  ax2 = plt.axes([0.65, 0.65, 0.2, 0.2])
```

# Subplots by Hand (2)

◻ `plt.axes ()` function

    – After creating the new axes, plot the second line

```
 6  # an array of 100 numbers evenly distributed between 0 and 100
 7  a = np.linspace(0, 10, 100)
 8  # an array of exponential values of -a
 9  c = np.exp(-a)
10  # standard axes
11  ax1 = plt.axes()
12  # use a as x axis and b as y axis for a line plot
13  plt.plot(a, c)
14
15  # an array of exponential values of -a*a
16  d = np.exp(-a*a)
17  # create an inset axes by setting the x and y position to 0.65
18  # and the x and y extents to 0.2
19  ax2 = plt.axes([0.65, 0.65, 0.2, 0.2])
20  # call plot() again to add another line
21  plt.plot(a, d, color = "green")
```

---

# Simple Grids of Subplots (1)

◻ `plt.subplots(rows, columns)`

    – Can be used to create subplots in rows and columns

    – The first optional arguments define the number of rows and columns of the subplot grid

    – The returned axes is a NumPy array containing the list of created Axes

        • When stacked in one direction, axes is a one-dimensional array

        • When a grid is created, axes is a two-dimensional array

**UNLV**

# Simple Grids of Subplots (2)

– The fig variable saves the figure object returned by subplots()

– The axs variable saves the axes array returned by subplots()

```
1  # subplot in two rows.
2  # save the figure and axes to varaibles
3  fig, axs = plt.subplots(2)
4  # use the fig variable to set the title
5  fig.suptitle('Vertically stacked subplots')
```

```
1  # subplot in two rows.
2  # save the figure and axes to varaibles
3  fig, axs = plt.subplots(1,2)
4  # use the fig variable to set the title
5  fig.suptitle('side-by-side subplots')
```



9

---

# Simple Grids of Subplots (3)

– Then, use index to specify the axes to plot

```
2   x = np.linspace(0, 10, 100)
3   # sine function of x squared
4   y = np.sin(x ** 2)
5
6   # subplot in two rows.
7   # save the figure and axes to varaibles
8   fig, axs = plt.subplots(2)
9   # use the fig variable to set the title
10  fig.suptitle('Vertically stacked subplots')
11  # Plot the figure at the first row
12  axs[0].plot(x, y)
13  # Plot the figure at the second row
14  axs[1].plot(x, -y)
```

```
2   x = np.linspace(0, 10, 100)
3   # sine function of x squared
4   y = np.sin(x ** 2)
5
6   # subplot in two rows.
7   # save the figure and axes to varaibles
8   fig, axs = plt.subplots(1,2)
9   # use the fig variable to set the title
10  fig.suptitle('side-by-side subplots')
11  # Plot the figure at the first column
12  axs[0].plot(x, y)
13  # Plot the figure at the second column
14  axs[1].plot(x, -y)
```



10

# Simple Grids of Subplots (4)

☐ When more than one column or rows are specified, use the index of two-dimensional array to specify the subplot

```
1  year = ['2018','2019']
2  quarter = ['Q1', 'Q2', 'Q3','Q4']
3
4  fig, axs = plt.subplots(len(year), len(quarter))
5  fig.subplots_adjust(hspace=0.4, wspace=0.8)
6
7  for y in range(len(year)):
8      for q in range(len(quarter)):
9          axs[y, q].set_title('Y'+year[y]+quarter[q])
```
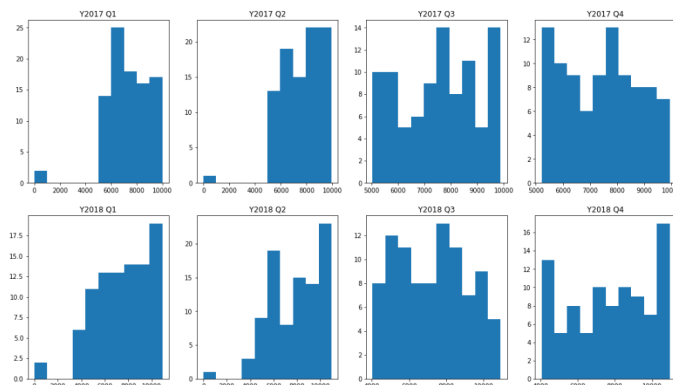
UNLV

---

☐ Lab

– Sales History Comparison
The program reads a file containing the quarterly sales data and plot histogram for comparison.

UNLV

# Three-Dimensional Plotting (1)

☐ It is enabled by the `mplot3d` toolkit

```
1  import matplotlib.pyplot as plt
2  |
3  # to enable the 3d plotting
4  from mpl_toolkits import mplot3d
5
6  %matplotlib inline
7
8  ax = plt.axes(projection='3d')
9
```

UNLV

---

# Three-Dimensional Plotting (2)

☐ `ax.plot3D()` and `ax.scatter3D()`

– The most basic three-dimensional plot is a line or collection of scatter plot created from sets of (x, y, z) triples

– The call signature for these is nearly identical to that of their two-dimensional counterparts

UNLV

# Three-Dimensional Plotting: Example

```python
1  ax = plt.axes(projection='3d')
2
3  # Data for a three-dimensional line
4  zline = np.linspace(0, 15, 1000)
5  xline = np.sin(zline)
6  yline = np.cos(zline)
7  # 3d line plot
8  ax.plot3D(xline, yline, zline, 'gray')
9
10  # Data for three-dimensional scattered points
11  zdata = 15 * np.random.random(100)
12  xdata = np.sin(zdata) + 0.1 * np.random.randn(100)
13  ydata = np.cos(zdata) + 0.1 * np.random.randn(100)
14  # 3d scatter plot
15  ax.scatter3D(xdata, ydata, zdata, c=zdata, cmap='Greens')
```

UNLV

---

☐ Lab

– Three-Dimensional Body Data
  This program reads the data from a csv file
  and then plot the data points

UNLV

□Exercise

– Revise Three-Dimensional Body Data
program
Please create three data sets: people less
than or equal 30 years old, people between
30 and 60, and people over 60. Plot the
three datasets on the 3D scatter plot.

UNLV

---

# Visualization with Seaborn

□Provides an API on top of Matplotlib that offers
choices for plot style and color defaults,
defines simple high-level functions for common
statistical plot types

□Integrates with the functionality provided by
Pandas DataFrames

□Build in some dataset that can be used for
learning purposes

– Get all the available dataset by using
`sns.get_dataset_names()`

UNLV

# Seaborn Style

```
1   import numpy as np
2   import matplotlib.pyplot as plt
3
4   %matplotlib inline
5
6   # Create some data
7   rng = np.random.RandomState(0)
8   x = np.linspace(0, 10, 500)
9   y = np.cumsum(rng.randn(500, 6), 0)
10
11  # The old matplotlib style
12  plt.plot(x, y)
13  plt.legend('ABCDEF', ncol=2, loc='upper left')
```

```
1   import numpy as np
2   import matplotlib.pyplot as plt
3   import seaborn as sns
4
5   %matplotlib inline
6
7   # Create some data
8   rng = np.random.RandomState(0)
9   x = np.linspace(0, 10, 500)
10  y = np.cumsum(rng.randn(500, 6), 0)
11
12  # use the seaborn style
13  sns.set()
14
15  plt.plot(x, y)
16  plt.legend('ABCDEF', ncol=2, loc='upper left')
```

# Distribution Plot

☐ `sns.kdeplot()`

– Instead of histogram, get a smooth estimate of the distribution using a kernel density estimation

```
1   import numpy as np
2   import pandas as pd
3   import matplotlib.pyplot as plt
4   %matplotlib inline
5
6   data = np.random.multivariate_normal([0, 0], [[5, 2], [2, 2]], size=2000)
7   data = pd.DataFrame(data, columns=['x', 'y'])
8
9   for col in 'xy':
10      plt.hist(data[col], density=True, alpha=0.5)
```

```
1   for col in 'xy':
2       sns.kdeplot(data[col], shade=True)
```
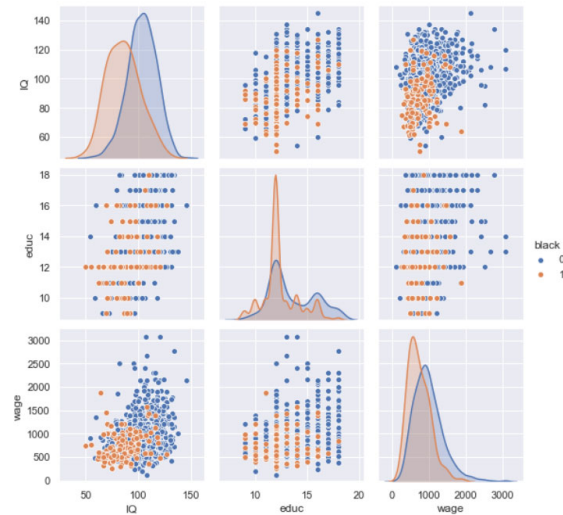
# Exploring Data (1)

☐ Pair Plots: `sns.pairplot()`

– Useful for exploring correlations between multidimensional data
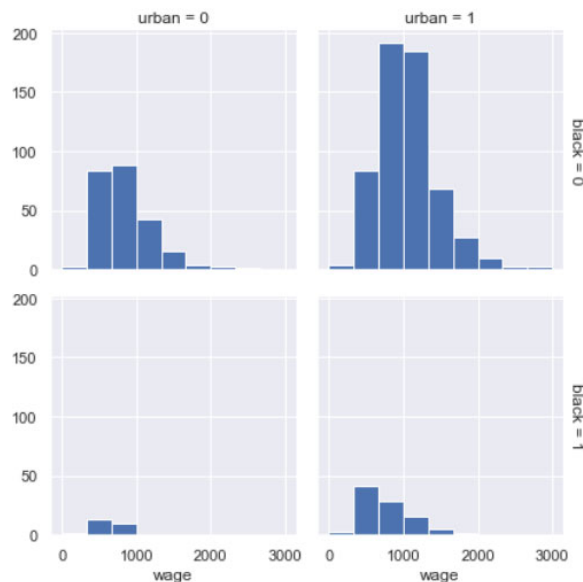
– plot all pairs of values against each other

UNLV

---

# Exploring Data (2)

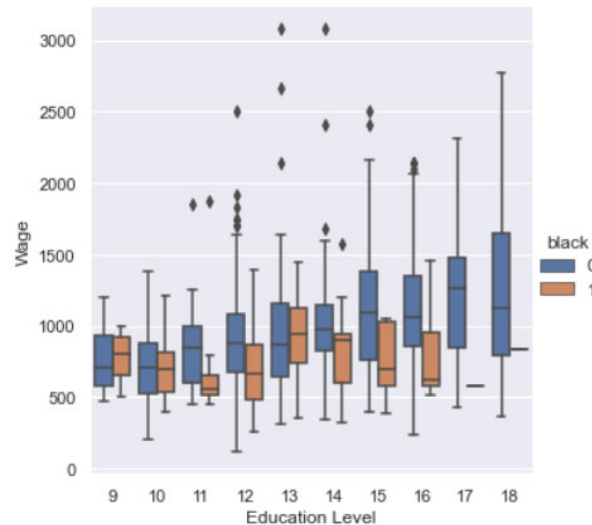☐ Faceted histograms: `sns.FacetGrid()`

– histograms of subsets

UNLV

# Exploring Data (3)

☐ Factor plots: `sns.catplot()`

   – view the distribution of a parameter within bins defined by any other parameter
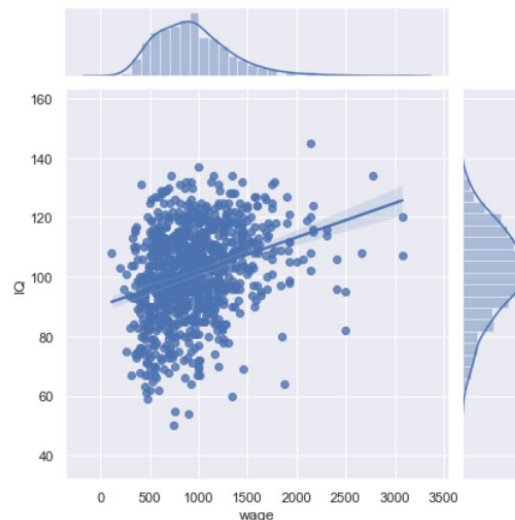
UNLV

---

# Exploring Data (4)

☐ Joint distributions: `sns.jointplot()`

   – show the joint distribution between different datasets, along with the associated marginal distributions

UNLV

❑Lab

– Tip Distribution
The program is a test of many useful functions in Seaborn, including pair plot, faceted histograms, factor plots, and joint distribution

UNLV

❑Exercise

– Wage in 1980
Write a program that read the data "wage.csv" and create the following diagrams:

• pair plot for four variables: IQ, education, black, and wage.

• faceted histograms for wage, categorized by black and urban (10 bins between 0 and 3000)

• factor plot: for each education level, compare the wage by black or not

• joint distribution on (1) IQ and wage and (2) experience and wage

– Data source: Introductory Econometrics: A Modern Approach, 6e by Jeffrey M. Wooldridge.

UNLV