Advanced Data Management (CSCI 640/490)

Data & Pandas

Dr. David Koop





Arrays

- Usually a fixed size—lists are meant to change size
- Are mutable—tuples are not
- Store only one type of data—lists and tuples can store anything • Are faster to access and manipulate than lists or tuples
- Can be multidimensional:

 - Can have list of lists or tuple of tuples but no guarantee on shape - Multidimensional arrays are rectangles, cubes, etc.









Why NumPy?

- Fast vectorized array operations for data munging and cleaning, subsetting and filtering, transformation, and any other kinds of computations
- Common array algorithms like sorting, unique, and set operations
- Efficient descriptive statistics and aggregating/summarizing data
- Data alignment and relational data manipulations for merging and joining together heterogeneous data sets
- elif-else branches
- Group-wise data manipulations (aggregation, transformation, function) application).





Northern Illinois University

• Expressing conditional logic as array expressions instead of loops with if









Speed Benefits

- Compare random number generation in pure Python versus numpy
- Python:
 - import random %timeit rolls list = [random.randrange(1,7)
- With NumPy:
 - %timeit rolls array = np.random.randint(1, 7, 60 000)
- Significant speedup (80x+)

D. Koop, CSCI 640/490, Spring 2024

for i in range(0, 60 000)]





Assignment 2

- Assignment 1 Questions with pandas, DuckDB, and polars
- CS 640 students do all, CS 490 do pandas & DuckDB (polars is EC)
- Can work by framework or by query
- Most questions can be answered with a single statement... but that statement can take a while to write
 - Read documentation
 - Check hints









Array Shape

- Our normal way of checking the size of a collection is... len
- How does this work for arrays?
- arr1 = np.array([1,2,3,6,9]) len(arr1) # 5
- arr2 = np.array([[1.5,2,3,4],[5,6,7,8]])len(arr2) # 2
- All dimension lengths \rightarrow shape: arr2.shape # (2,4)
- Number of dimensions: arr2.ndim # 2
- Can also reshape an array:
 - arr2.reshape(4,2)
 - arr2.reshape(-1,2) # what happens here?









Array Programming

- Lists:
 - c = []
 for i in range(len(a)):
 c.append(a[i] + b[i])
- How to improve this?

D. Koop, CSCI 640/490, Spring 2024





7

Array Programming

- Lists:
 - C = | | for i in range(len(a)): c.append(a[i] + b[i])
 - -c = [aa + bb for aa, bb in zip(a,b)]
- NumPy arrays:
 - -c = a + b
- More functional-style than imperative
- Internal iteration instead of external









Operations

- a = np.array([1, 2, 3])b = np.array([6, 4, 3])
- (Array, Array) Operations (**Element-wise**)
 - Addition, Subtraction, Multiplication
 - -a + b # array([7, 6, 6])
- (Scalar, Array) Operations (**Broadcasting**):
 - Addition, Subtraction, Multiplication, Division, Exponentiation
 - a ** 2 # array([1, 4, 9])
 - -b + 3 # array([9, 7, 6])









More on Array Creation

- Zeros: np.zeros(10)
- Ones: np.ones((4,5)) # shape
- Empty: np.empty((2,2))
- _like versions: pass an existing array and matches shape with specified contents
- Range: np.arange(15) # constructs an array, not iterator!





Indexing

- Same as with lists plus shorthand for 2D+
 - $\operatorname{arr1} = \operatorname{np.array}([6, 7, 8, 0, 1])$
 - arr1[1]
 - arr1[-1]
- What about two dimensions?
 - $\operatorname{arr2} = \operatorname{np.array}([[1.5, 2, 3, 4], [5, 6, 7, 8]])$
 - arr[1][1]
 - arr[1,1] # shorthand





2D Indexing



D. Koop, CSCI 640/490, Spring 2024

	axis 1	
0	1	2
, 0	0, 1	0, 2
, 0	1, 1	1, 2
, 0	2, 1	2, 2

[W. McKinney, Python for Data Analysis]









Slicing

- 1D: Similar to lists
 - arr1 = np.array([6, 7, 8, 0, 1])
 - arr1[2:5] # np.array([8,0,1]), sort of
- Can **mutate** original array:
 - arr1[2:5] = 3 # supports assignment
 - arr1 # the original array changed
- Slicing returns views (copy the array if original array shouldn't change)
 - arr1[2:5] # a view
 - arr1[2:5].copy() # a new array





Slicing

- 2D+: comma separated indices as shorthand:
 - arr2 = np.array([[1.5,2,3,4],[5,6,7,8]])
 - a[1:3,1:3]
 - a[1:3,:] # works like in single-dimensional lists
- Can combine index and slice in different dimensions
 - a[1,:] # gives a row
 - a[:,1] # gives a column

D. Koop, CSCI 640/490, Spring 2024

shorthand: 4], [5, 6, 7, 8]])

ingle-dimensional lists erent dimensions









How to obtain the blue slice from array arr?

D. Koop, CSCI 640/490, Spring 2024

[W. McKinney, Python for Data Analysis]



Northern Illinois University













How to obtain the blue slice from array arr?

D. Koop, CSCI 640/490, Spring 2024



[W. McKinney, Python for Data Analysis]















How to obtain the blue slice from array arr?

D. Koop, CSCI 640/490, Spring 2024

Expression	Shape	
arr[:2, 1:]	(2, 2)	
arr[2]	(3,)	
arr[2, :]	(3,)	

(1, 3)

arr[2:, :]

















How to obtain the blue slice from array arr?

D. Koop, CSCI 640/490, Spring 2024

Expression	Shape
arr[:2, 1:]	(2, 2)
arr[2]	(3,)
arr[2, :]	(3,)
arr[2:, :]	(1, 3)
[· · · · · ·]	(2 2)
arr[:, :2]	(3, 2)

[W. McKinney, Python for Data Analysis]















How to obtain the blue slice from array arr?

D. Koop, CSCI 640/490, Spring 2024

Expression	Shape
arr[:2, 1:]	(2, 2)
arr[2]	(3,)
$\operatorname{arr}[2, :]$	(3,)
arr[z:, :]	(1,3)
arr[:, :2]	(3, 2)
arr[1, :2]	(2,)
arr[1:2, :2]	(1, 2)

[W. McKinney, Python for Data Analysis]











Reshaping

- reshape:
 - arr2.reshape(4,2) # returns new view
- resize:
 - arr2.resize(4,2) # no return, modifies arr2 in place
- flatten:
 - arr2.flatten() # array([1.5,2.,3.,4.,5.,6.,7.,8.])
- ravel:
 - arr2.ravel() # array([1.5,2.,3.,4.,5.,6.,7.,8.])
- flatten and ravel look the same, but ravel is a view





Boolean Indexing

- names == 'Bob' gives back booleans that represent the element-wise comparison with the array names
- Boolean arrays can be used to index into another array:
 - data[names == 'Bob']
- Can even mix and match with integer slicing
- Can do boolean operations (&, |) between arrays (just like addition, subtraction)
 - data[(names == 'Bob') | (names == 'Will')]
- Note: or and and do not work with arrays
- We can set values too! data [data < 0] = 0





Array Transformations

- Transpose
 - arr2.T # flip rows and columns
- Stacking: take iterable of arrays and stack them horizontally/vertically
 - $\operatorname{arrh1} = \operatorname{np.arange}(3)$
 - $\operatorname{arrh2} = \operatorname{np.arange}(3, 6)$
 - np.vstack([arrh1, arrh2])
 - np.hstack([arr1.T, arr2.T]) # ???





numpy Functions

- Unary: abs, sqrt, log, ceil, sin, cos, tan, arccos, arcsin, ...
- Binary: add, subtract, multiple, divide, $\ldots <$, >, >=, <=, ==, !=
- Statistics: sum, mean, std, min, max, argmin, argmax
- Boolean: any, all
- Others: sort, unique
- Linear Algebra (numpy.linalg)
- Pseudorandom Number Generation (numpy.random)





pandas

- data analysis fast and easy in Python
- Built on top of NumPy
- Requirements:
 - Data structures with labeled axes (aligning data)
 - Time series data
 - Arithmetic operations that include metadata (labels)
 - Handle missing data
 - Merge and relational operations

D. Koop, CSCI 640/490, Spring 2024

Contains high-level data structures and manipulation tools designed to make









Pandas Code Conventions

- Universal:
 - import pandas as pd
- Also used:
 - from pandas import Series, DataFrame







Series

- A one-dimensional array (with a type) with an **index**
- Index defaults to numbers but can also be text (like a dictionary)
- Allows easier reference to specific items
- obj = pd.Series([7,14,-2,1])
- Basically two arrays: obj.values and obj.index
- Can specify the index explicitly and use strings
- obj2 = pd.Series([4, 7, -5, 3])index=['d', 'b', 'a', 'c'])
- Kind of like fixed-length, ordered dictionary + can create from a dictionary
- obj3 = pd.Series({'Ohio': 35000, 'Texas': 71000,

D. Koop, CSCI 640/490, Spring 2024

'Oregon': 16000, 'Utah': 5000})









Series

- Indexing: s[1] Or s['Oregon']
- Can check for missing data: pd.isnull(s) Or pd.notnull(s)
- Both index and values can have an associated name:
 - s.name = 'population'; s.index.name = 'state'
- Addition and NumPy ops work as expected and preserve the index-value link
- These operations **align**:

In [28]: Out[28]:	obj3	In [29]: obj Out[29]:	4	In [30]: obj Out[30]:	3 + obj4
Ohio Oregon Texas Utah dtype: i	35000 16000 71000 5000 nt64	California Ohio Oregon Texas dtype: float	NaN 35000 16000 71000 64	California Ohio Oregon Texas Utah dtype: float [W. N	NaN 70000 32000 142000 NaN 64 4CKinney, Pythor
Spring 202	24				Northern









- A dictionary of Series (labels for each series)
- A spreadsheet with column headers
- Has an index shared with each series
- Allows easy reference to any cell
- df = DataFrame({'state': ['Ohio', 'Ohio', 'Ohio', 'Nevada'], 'year': [2000, 2001, 2002, 2001], 'pop': [1.5, 1.7, 3.6, 2.4]})
- Index is automatically assigned just as with a series but can be passed in as well via index kwarg
- Can reassign column names by passing columns kwarg









df =	df = pd.read_csv('penguins_lter.csv')									
:	studyName	Sample Number	Species	Region	Island	Stage	Individual ID	Clutch Completion	Date Egg	Culmen Length (mm)
0	PAL0708	1	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A1	Yes	11/11/07	39.1
1	PAL0708	2	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A2	Yes	11/11/07	39.5
2	PAL0708	3	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N2A1	Yes	11/16/07	40.3
3	PAL0708	4	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N2A2	Yes	11/16/07	NaN
4	PAL0708	5	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N3A1	Yes	11/16/07	36.7
339	PAL0910	120	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	Adult, 1 Egg Stage	N38A2	No	12/1/09	NaN
340	PAL0910	121	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	Adult, 1 Egg Stage	N39A1	Yes	11/22/09	46.8
341	PAL0910	122	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	Adult, 1 Egg Stage	N39A2	Yes	11/22/09	50.4
342	PAL0910	123	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	Adult, 1 Egg Stage	N43A1	Yes	11/22/09	45.2
343	PAL0910	124	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	Adult, 1 Egg Stage	N43A2	Yes	11/22/09	49.9

344 rows × 17 columns









		df =	<pre>pd.read_csv('penguins_lter.csv')</pre>									
Column	Name	es	studyName	Sample Number	Species	Region	Island	Stage	Individual ID	Clutch Completion	Date Egg	Culmen Length (mm)
	-	0	PAL0708	1	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A1	Yes	11/11/07	39.1
		1	PAL0708	2	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A2	Yes	11/11/07	39.5
		2	PAL0708	3	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N2A1	Yes	11/16/07	40.3
		3	PAL0708	4	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N2A2	Yes	11/16/07	NaN
		4	PAL0708	5	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N3A1	Yes	11/16/07	36.7
		339	PAL0910	120	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	Adult, 1 Egg Stage	N38A2	No	12/1/09	NaN
		340	PAL0910	121	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	Adult, 1 Egg Stage	N39A1	Yes	11/22/09	46.8
		341	PAL0910	122	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	Adult, 1 Egg Stage	N39A2	Yes	11/22/09	50.4
		342	PAL0910	123	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	Adult, 1 Egg Stage	N43A1	Yes	11/22/09	45.2
		343	PAL0910	124	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	Adult, 1 Egg Stage	N43A2	Yes	11/22/09	49.9

344 rows × 17 columns









	df =	pd.read_csv	<pre>/('penguins_l'</pre>	ter.csv')							
Column Names		studyName	Sample Number	Species	Region	Island	Stage	Individual ID	Clutch Completion	Date Egg	Culmen Length (mm)
	0	PAL0708	1	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A1	Yes	11/11/07	39.1
	1	PAL0708	2	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A2	Yes	11/11/07	39.5
	2	PAL0708	3	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N2A1	Yes	11/16/07	40.3
	3	PAL0708	4	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N2A2	Yes	11/16/07	NaN
	4	PAL0708	5	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N3A1	Yes	11/16/07	36.7
Index											
	339	PAL0910	120	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	Adult, 1 Egg Stage	N38A2	No	12/1/09	NaN
	340	PAL0910	121	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	Adult, 1 Egg Stage	N39A1	Yes	11/22/09	46.8
	341	PAL0910	122	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	Adult, 1 Egg Stage	N39A2	Yes	11/22/09	50.4
	342	PAL0910	123	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	Adult, 1 Egg Stage	N43A1	Yes	11/22/09	45.2
	343	PAL0910	124	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	Adult, 1 Egg Stage	N43A2	Yes	11/22/09	49.9

344 rows × 17 columns

D. Koop, CSCI 640/490, Spring 2024









	df = pd.read_csv('penguins_lter.csv') M Names siudyName Sample Number Species Region Island Stage Individual ID Clutch Bag Date Egg Culmen Length (mm) 0 PAL0708 1 Adelie Penguin (Pygoscelis adeliae) Anvers Torgersen Adult, 1 Egg Stage N1A1 Yes 11/11/07 39.5 1 PAL0708 2 Adelie Penguin (Pygoscelis adeliae) Anvers Torgersen Adult, 1 Egg Stage N1A2 Yes 11/11/07 39.5 2 PAL0708 3 Adelie Penguin (Pygoscelis adeliae) Anvers Torgersen Adult, 1 Egg Stage N2A1 Yes 11/107 39.5 2 PAL0708 3 Adelie Penguin (Pygoscelis adeliae) Anvers Torgersen Adult, 1 Egg N2A1 Yes 11/16/07 40.3										
Column Names		studyName	Sample Number	Species	Region	Island	Stage	Individual ID	Clutch Completion	Date Egg	Culmen Length (mm)
	0	PAL0708	1	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A1	Yes	11/11/07	39.1
	1	PAL0708	2	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A2	Yes	11/11/07	39.5
	2	PAL0708	3	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N2A1	Yes	11/16/07	40.3
	3	PAL0708	4	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N2A2	Yes	11/16/07	NaN
	4	PAL0708	5	Adelie Penguin (Pygoscelis adeliae)	Anvers	Torgersen	Adult, 1 Egg Stage	N3A1	Yes	11/16/07	36.7
Index											
	339	PAL0910	120	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	Adult, 1 Egg Stage	N38A2	No	12/1/09	NaN
	340	PAL0910	121	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	Adult, 1 Egg Stage	N39A1	Yes	11/22/09	46.8
	341	PAL0910	122	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	Adult, 1 Egg Stage	N39A2	Yes	11/22/09	50.4
	342	PAL0910	123	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	Adult, 1 Egg Stage	N43A1	Yes	11/22/09	45.2
	343	PAL0910	124	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	Adult, 1 Egg Stage	N43A2	Yes	11/22/09	49.9

344 rows × 17 columns

D. Koop, CSCI 640/490, Spring 2024











344 rows × 17 columns

D. Koop, CSCI 640/490, Spring 2024

ies	Region	Island	Stage	Individual ID	Clutch Completion	Date Egg	Culmen Length (mm)
elis ae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A1	Yes	11/11/07	39.1
elis ae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A2	Yes	11/11/07	39.5
elis ae)	Anvers	Torgersen	Adult, 1 Egg Stage	N2A1	Yes	11/16/07	40.3
elis ae)	Anvers	Torgersen	Adult, 1 Egg Stage	N2A2	Yes	11/16/07	NaN
elis ae)	Anvers	Torgersen	Adult, 1 Egg Stage	N3A1	Yes	11/16/07	36.7
elis ua)	Anvers	Biscoe	Adult, 1 Egg Stage	N38A2	No	12/1/09	NaN
elis ua)	Anvers	Biscoe	Adult, 1 Egg Stage	N39A1	Yes	11/22/09	46.8
elis ua)	Anvers	Biscoe	Adult, 1 Egg Stage	N39A2	Yes	11/22/09	50.4
elis ua)	Anvers	Biscoe	Adult, 1 Egg Stage	N43A1	Yes	11/22/09	45.2
elis ua)	Anvers	Biscoe	Adult, 1 Egg Stage	N43A2	Yes	11/22/09	49.9











D. Koop, CSCI 640/490, Spring 2024

ies	Region	Island	Stage	Individual ID	Clutch Completion	Date Egg	Culmen Length (mm)
elis ae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A1	Yes	11/11/07	39.1
elis ae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A2	Yes	11/11/07	39.5
elis ae)	Anvers	Torgersen	Adult, 1 Egg Stage	N2A1	Yes	11/16/07	40.3
elis ae)	Anvers	Torgersen	Adult, 1 Egg Stage	N2A2	Yes	11/16/07	NaN
elis ae)	Anvers	Torgersen	Adult, 1 Egg Stage	N3A1	Yes	11/16/07	36.7
elis ua)	Anvers	Biscoe	Adult, 1 Egg Stage	N38A2	No	12/1/09	NaN
elis ua)	Anvers	Biscoe	Adult, 1 Egg Stage	N39A1	Yes	11/22/09	46.8
elis ua)	Anvers	Biscoe	Adult, 1 Egg Stage	N39A2	Yes	11/22/09	50.4
elis ua)	Anvers	Biscoe	Adult, 1 Egg Stage	N43A1	Yes	11/22/09	45.2
elis ua)	Anvers	Biscoe	Adult, 1 Egg Stage	N43A2	Yes	11/22/09	49.9











D. Koop, CSCI 640/490, Spring 2024

ies	Region	Island	Stage	Individual ID	Clutch Completion	Date Egg	Culmen Length (mm)
elis ae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A1	Yes	11/11/07	39.1
elis ae)	Anvers	Torgersen	Adult, 1 Egg Stage	N1A2	Yes	11/11/07	39.5
elis ae)	Anvers	Torgersen	Adult, 1 Egg Stage	N2A1	Yes	11/16/07	40.3
elis ae)	Anvers	Torgersen	Adult, 1 Egg Stage	N2A2	Yes	11/16/07	NaN
elis ae)	Anvers	Torgersen	Adult, 1 Egg Stage	N3A1	Yes	11/16/07	Missina F
							""
elis ua)	Anvers	Biscoe	Adult, 1 Egg Stage	N38A2	No	12/1/09	NaN
elis ua)	Anvers	Biscoe	Adult, 1 Egg Stage	N39A1	Yes	11/22/09	46.8
elis ua)	Anvers	Biscoe	Adult, 1 Egg Stage	N39A2	Yes	11/22/09	50.4
elis ua)	Anvers	Biscoe	Adult, 1 Egg Stage	N43A1	Yes	11/22/09	45.2
elis ua)	Anvers	Biscoe	Adult, 1 Egg Stage	N43A2	Yes	11/22/09	49.9











Chicago Food Inspections Example

- Use pandas to analyze food inspection data
 - (see notebook)





