### Advanced Data Management (CSCI 490/680)

#### 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.





2

## 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
- Expressing conditional logic as array expressions instead of loops with ifelif-else branches
- Group-wise data manipulations (aggregation, transformation, function) application).





Northern Illinois University









## NumPy Arrays

- data1 = [6, 7.5, 8, 0, 1]arr1 = np.array(data1)
- Zeros: np.zeros(10), Ones: np.ones((4,5)), **Empty:** np.empty((2,2))
- Types: Each array has a fixed type unlike other variables in python

#### D. Koop, CSCI 680/490, Spring 2021

# • # of dimensions: arr2.ndim, Shape: arr2.shape, Type: arr2.dtype





## 2D Array Slicing

#### How to obtain the blue slice from array arr?





D. Koop, CSCI 680/490, Spring 2021



[W. McKinney, Python for Data Analysis]



Northern Illinois University







## 2D Array Slicing

#### How to obtain the blue slice from array arr?





D. Koop, CSCI 680/490, Spring 2021



[W. McKinney, Python for Data Analysis]



Northern Illinois University







## 2D Array Slicing

#### How to obtain the blue slice from array arr?





D. Koop, CSCI 680/490, Spring 2021



[W. McKinney, Python for Data Analysis]









### 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









## <u>Assignment 1</u>

- Due Monday, Feb. 1 at 11:59pm
- Using Python for data analysis on Info Wanted ads Provided a1.ipynb file (right-click and download) Use basic python for now to demonstrate language knowledge
- No pandas (for now)
- Use Anaconda or hosted Python environment
- Turn .ipynb file in via Blackboard
- Notes:
  - Bug in URL (https instead of http),
  - Bug in Problem 1 solution





## Other Operations

- Fancy Indexing: arr[[1,2,3]]
- Transposing arrays: arr.T
- Reshaping arrays: arr.reshape((3,5))
- Unary universal functions (ufuncs): np.sqrt, np.exp
- Binary universal functions: np.add, np.maximum







## Unary Universal Functions

Description
Compute the absolute value
Compute the square root of
Compute the square of each
Compute the exponent e <sup>x</sup> of
Natural logarithm (base <i>e</i> ),
Compute the sign of each el
Compute the ceiling of each number)
Compute the floor of each e
Round elements to the near
Return fractional and integr
Return boolean array indica
Return boolean array indica respectively
Regular and hyperbolic trigo
Inverse trigonometric functi
Compute truth value of not

#### D. Koop, CSCI 680/490, Spring 2021

e element-wise for integer, floating-point, or complex values each element (equivalent to arr \*\* 0.5) element (equivalent to arr \*\* 2) f each element log base 10, log base 2, and log(1 + x), respectively

lement: 1 (positive), 0 (zero), or –1 (negative) element (i.e., the smallest integer greater than or equal to that

element (i.e., the largest integer less than or equal to each element)

- rest integer, preserving the dtype
- ral parts of array as a separate array
- ating whether each value is NaN (Not a Number)
- nting whether each element is finite (non-inf, non-NaN) or infinite,

onometric functions

ions

t  $\times$  element-wise (equivalent to  $\sim arr$ ).

#### [W. McKinney, Python for Data Analysis]









## **Binary Universal Functions**

Function	Description
add	Add corresponding
subtract	Subtract elements
multiply	Multiply array elen
divide, floor_divide	Divide or floor divi
рожег	Raise elements in f
maximum, fmax	Element-wise max
minimum, fmin	Element-wise mini
mod	Element-wise mod
copysign	Copy sign of values
greater, greater_equal,	Perform element-v
less, less_equal,	operators >, >=,
equal, not_equal	
logical_and,	Compute element-
logical_or, logical_xor	&   , ^)

#### D. Koop, CSCI 680/490, Spring 2021

- elements in arrays
- in second array from first array
- nents
- de (truncating the remainder)
- first array to powers indicated in second array
- (imum; fmax ignores NaN
- imum; fmin ignores NaN
- lulus (remainder of division)
- s in second argument to values in first argument
- wise comparison, yielding boolean array (equivalent to infix

<, <=, ==, !=)

-wise truth value of logical operation (equivalent to infix operators

[W. McKinney, Python for Data Analysis]









## Statistical Methods

Method	Description
SUM	Sum of all the elements in the arra
mean	Arithmetic mean; zero-length array
std, var	Standard deviation and variance, reddenominator n
min, max	Minimum and maximum
argmin, argmax	Indices of minimum and maximum
CUMSUM	Cumulative sum of elements starting
cumprod	Cumulative product of elements sta

#### D. Koop, CSCI 680/490, Spring 2021

- ay or along an axis; zero-length arrays have sum 0
- ys have NaN mean
- espectively, with optional degrees of freedom adjustment (default

- elements, respectively
- ng from 0
- arting from 1

[W. McKinney, Python for Data Analysis]





/SIS] 11

## More

- Other methods:
  - any and all
  - sort
  - unique
- Linear Algebra (numpy.linalg)
- Pseudorandom Number Generation (numpy.random)





### Data

- What is data?
  - Types
  - Semantics
- How is data structured?
  - Tables (Data Frames)
  - Databases
  - Data Cubes
- What formats is data stored in?
- Raw versus derived data





### Data

• What is this data?

R011	42ND STREET & 8TH AVENUE	00228985	00008471	00000441	00001455	00000134	00033341	00071255
R170	14TH STREET-UNION SQUARE	00224603	00011051	00000827	00003026	00000660	00089367	00199841
R046	42ND STREET & GRAND CENTRAL	00207758	00007908	00000323	00001183	00003001	00040759	00096613

- Semantics: real-world meaning of the data
- Type: structural or mathematical interpretation
- Both often require metadata
  - Sometimes we can infer some of this information
  - Line between data and metadata isn't always clear

this information isn't always clear





### Data

	REMOTE	STATION	FF 1
1	R011	42ND STREET & 8TH AVENUE	00228985
2	R170	14TH STREET-UNION SQUARE	00224603
3	R046	42ND STREET & GRAND CENTRAL	00207758
4	R012	34TH STREET & 8TH AVENUE	00188311
5	R293	34TH STREET – PENN STATION	00168768
6	R033	42ND STREET/TIMES SQUARE	00159382
7	R022	34TH STREET & 6TH AVENUE	00156008
8	R084	59TH STREET/COLUMBUS CIRCLE	00155262
9	R020	47-50 STREETS/ROCKEFELLER	00143500
10	R179	86TH STREET-LEXINGTON AVE	00142169
11	R023	34TH STREET & 6TH AVENUE	00134052
12	R029	PARK PLACE	00121614
13	R047	42ND STREET & GRAND CENTRAL	00100742
14	R031	34TH STREET & 7TH AVENUE	00095076
15	R017	LEXINGTON AVENUE	00094655
16	R175	8TH AVENUE-14TH STREET	00094313
17	R057	BARCLAYS CENTER	00093804
18	R138	WEST 4TH ST-WASHINGTON SO	00093562

							_
/	SEN/DIS	7-D AFAS UNL	D AFAS/RMF L	JOINT RR TKT	7-D UNL	30-D UNL	
5	00008471	00000441	00001455	00000134	00033341	00071255	
3	00011051	00000827	00003026	00000660	00089367	00199841	
3	00007908	00000323	00001183	00003001	00040759	00096613	
L	00006490	00000498	00001279	00003622	00035527	00067483	
3	00006155	00000523	00001065	00005031	00030645	00054376	
2	00005945	00000378	00001205	00000690	00058931	00078644	
3	00006276	00000487	00001543	00000712	00058910	00110466	
2	00009484	00000589	00002071	00000542	00053397	00113966	
)	00006402	00000384	00001159	00000723	00037978	00090745	
9	00010367	00000470	00001839	00000271	00050328	00125250	
2	00005005	00000348	00001112	00000649	00031531	00075040	
1	00004311	00000287	00000931	00000792	00025404	00065362	
2	00004273	00000185	00000704	00001241	00022808	00068216	
5	00003990	00000232	00000727	00001459	00024284	00038671	
5	00004688	00000190	00000833	00000754	00020018	00055066	
3	00003907	00000286	00001144	00000256	00038272	00074661	
1	00004204	00000454	00001386	00001491	00039113	00068119	
,	00004677	00000251	00000965	00000127	00031628	00074458	





## Dataset Types

#### → Tables

→ Networks





 $\rightarrow$  Multidimensional Table



→ Trees



D. Koop, CSCI 680/490, Spring 2021

→ Geometry (Spatial)







Northern Illinois University





## Data Terminology

- Items
  - An **item** is an individual discrete entity
  - e.g., a row in a table
- Attributes
  - logged
  - a.k.a. variable, (data) dimension
  - e.g., a column in a table

#### D. Koop, CSCI 680/490, Spring 2021

#### - An attribute is some specific property that can be measured, observed, or





### Tables

Α	В	С	S	Т	U
Order ID	Order Date	Order Priority	Product Container	Product Base Margin	Ship Date
3	10/14/06	5-Low	Large Box	0.8	10/21/06
6	2/21/08	4-Not Specified	Small Pack	0.55	2/22/08
32	7/16/07	2-High	Small Pack	0.79	7/17/07
32	7/16/07	2-High	Jumbo Box	•1 .	7/17/07
32	7/16/07	2-High	Medium Box	attribute	7/18/07
32	7/16/07	2-High	Medium Box	0.05	7/18/07
35	10/23/07	4-Not Specified	Wrap Bag	0.52	10/24/07
35	10/23/07	4-Not Specified	Small Box	0.58	10/25/07
36	11/3/07	1-Urgent	Small Box	0.55	11/3/07
65	3/18/07	1-Urgent	Small Pack	0.49	3/19/07
66	1 (20 (05	5-Low	Wrap Bag	0.56	1/20/05
69	item <sup>5</sup>	4-Not Specified	Small Pack	0.44	6/6/05
69	5	4-Not Specified	Wrap Bag	0.6	6/6/05
70	12/18/06	5-Low	Small Box	0.59	12/23/06
70	12/18/06	5-Low	Wrap Bag	0.82	12/23/06
96	4/17/05	2-High	Small Box	0.55	4/19/05
97	1/29/06	3-Medium	Small Box	0.38	1/30/06
129	11/19/08	5-Low	Small Box	0.37	11/28/08
130	5/8/08	2-High	Small Box	0.37	5/9/08
130	5/8/08	2-High	Medium Box	0.38	5/10/08
130	5/8/08	2-High	Small Box	0.6	5/11/08
132	6/11/06	3-Medium	Medium Box	0.6	6/12/06
132	6/11/06	3-Medium	Jumbo Box	0.69	6/14/06
134	5/1/08	4-Not Specified	Large Box	0.82	5/3/08
135	10/21/07	4-Not Specified	Small Pack	0.64	10/23/07
166	9/12/07	2-High	Small Box	0.55	9/14/07
193	8/8/06	1-Urgent	Medium Box	0.57	8/10/06
194	4/5/08	3-Medium	Wrap Bag	0.42	4/7/08





### Tables





- - row ~ item (usually)
- column ~ attribute
- label labe
- Key: identifies each item (row)
  - Usually unique
  - Allows join of data from 2+ tables
  - Compound key: key split among multiple columns, e.g. (state, year) for population
- Multidimensional:
  - Split compound key

D. Koop, CSCI 680/490, Spring 2021

#### Data organized by rows & columns







### Attribute Types

# Categorical

D. Koop, CSCI 680/490, Spring 2021







Northern Illinois University





## Categorial, Ordinal, and Quantitative

Α	В	(	2	S	Т	U
Order ID	Order Date	Order Priorit	ty	Product Container	Product Base Margin	Ship Date
3	10/14/06	5-Low		Large Box	0.8	10/21/06
6	2/21/08	4-Not Speci	fied	Small Pack	0.55	2/22/08
32	7/16/07	2-High		Small Pack	0.79	7/17/07
32	7/16/07	2-High		Jumbo Box	0.72	7/17/07
32	7/16/07	2-High		Medium Box	0.6	7/18/07
32	7/16/07	2-High		Medium Box	0.65	7/18/07
35	10/23/07	4-Not Speci	fied	Wrap Bag	0.52	10/24/07
35	10/23/07	4-Not Speci	fied	Small Box	0.58	10/25/07
36	11/3/07	1-Urgent		Small Box	0.55	11/3/07
65	3/18/07	1-Urgent		Small Pack	0.49	3/19/07
66	1/20/05	5-Low		Wrap Bag	0.56	1/20/05
69	6/4/05	4-Not Speci	fied	Small Dack	0.44	6/6/05
69	6/4/05	4-Not Spec	anar	ntitativo	0.6	6/6/05
70	12/18/06	5-Low	yuai	illative	0.59	12/23/06
70	12/18/06	5-Low	ordi	nal	0.82	12/23/06
96	4/17/05	2-High	UI UI		0.55	4/19/05
97	1/29/06	3-Medium	cate	porical	0.38	1/30/06
129	11/19/08	5-Low	cute	5011041	0.37	11/28/08
130	5/8/08	2-High		Small Box	0.37	5/9/08
130	5/8/08	2-High		Medium Box	0.38	5/10/08
130	5/8/08	2-High		Small Box	0.6	5/11/08
132	6/11/06	3-Medium		Medium Box	0.6	6/12/06
132	6/11/06	3-Medium		Jumbo Box	0.69	6/14/06
134	5/1/08	4-Not Speci	fied	Large Box	0.82	5/3/08
135	10/21/07	4-Not Speci	fied	Small Pack	0.64	10/23/07
166	9/12/07	2-High		Small Box	0.55	9/14/07
193	8/8/06	1-Urgent		Medium Box	0.57	8/10/06
194	4/5/08	3-Medium		Wrap Bag	0.42	4/7/08

#### D. Koop, CSCI 680/490, Spring 2021





21

## Categorial, Ordinal, and Quantitative

Α	В	(	C	S	Т	U
Order ID	Order Date	Order Priori	ty	Product Container	Product Base Margin	Ship Date
3	10/14/06	5-Low		Large Box	0.8	10/21/06
6	2/21/08	4-Not Speci	fied	Small Pack	0.55	2/22/08
32	7/16/07	2-High		Small Pack	0.79	7/17/07
32	7/16/07	2-High		Jumbo Box	0.72	7/17/07
32	7/16/07	2-High		Medium Box	0.6	7/18/07
32	7/16/07	2-High		Medium Box	0.65	7/18/07
35	10/23/07	4-Not Speci	fied	Wrap Bag	0.52	10/24/07
35	10/23/07	4-Not Speci	fied	Small Box	0.58	10/25/07
36	11/3/07	1-Urgent		Small Box	0.55	11/3/07
65	3/18/07	1-Urgent		Small Pack	0.49	3/19/07
66	1/20/05	5-Low		Wrap Bag	0.56	1/20/05
69	6/4/05	4-Not Spec	fied	Small Pack	0.44	6/6/05
69	6/4/05	4-Not Spec	ana	atitativo	0.6	6/6/05
70	12/18/06	5-Low	yuai	illative	0.59	12/23/06
70	12/18/06	5-Low	ordi	nal	0.82	12/23/06
96	4/17/05	2-High		1101	0.55	4/19/05
97	1/29/06	3-Medium	cate	gorical	0.38	1/30/06
129	11/19/08	5-Low	cute	5011041	0.37	11/28/08
130	5/8/08	2-High		Small Box	0.37	5/9/08
130	5/8/08	2-High		Medium Box	0.38	5/10/08
130	5/8/08	2-High		Small Box	0.6	5/11/08
132	6/11/06	3-Medium		Medium Box	0.6	6/12/06
132	6/11/06	3-Medium		Jumbo Box	0.69	6/14/06
134	5/1/08	4-Not Speci	fied	Large Box	0.82	5/3/08
135	10/21/07	4-Not Specified		Small Pack	0.64	10/23/07
166	9/12/07	2-High		Small Box	0.55	9/14/07
193	8/8/06	1-Urgent		Medium Box	0.57	8/10/06
194	4/5/08	3-Medium		Wrap Bag	0.42	4/7/08
	1 / 1 / 0 0	a				1 (11) (0.0)









## Attribute Types

- May be further specified for computational storage/processing - Categorical: string, boolean, blood type
- - Ordered: enumeration, t-shirt size
  - Quantitative: integer, float, fixed decimal, datetime
- Sometimes, types can be inferred from the data
  - e.g. numbers and none have decimal points  $\rightarrow$  integer
  - could be incorrect (data doesn't have floats, but could be)









### Ordering Direction

#### Sequential



D. Koop, CSCI 680/490, Spring 2021

# Diverging







Northern Illinois University





## Sequential and Diverging Data

- Sequential: homogenous range from a minimum to a maximum
  - Examples: Land elevations, ocean depths
- Diverging: can be deconstructed into two sequences pointing in opposite directions
  - Has a **zero point** (not necessary 0)
  - Example: Map of both land elevation and ocean depth

1000 500 -500 -1000-1500 -2000 -2500 -3000 -3500













## Sequential and Diverging Data

- Sequential: homogenous range from a minimum to a maximum
  - Examples: Land elevations, ocean depths
- Diverging: can be deconstructed into two sequences pointing in opposite directions
  - Has a **zero point** (not necessary 0)
  - Example: Map of both land elevation and ocean depth

















### Cyclic Data













### Cyclic Data



D. Koop, CSCI 680/490, Spring 2021

month







- The meaning of the data
- Example: 94023, 90210, 02747, 60115









- The meaning of the data
- Example: 94023, 90210, 02747, 60115
  - Attendance at college football games?









- The meaning of the data
- Example: 94023, 90210, 02747, 60115
  - Attendance at college football games?
  - Salaries?









- The meaning of the data
- Example: 94023, 90210, 02747, 60115
  - Attendance at college football games?
  - Salaries?
  - Zip codes?
- Cannot always infer based on what the data looks like • Often require semantics to better understand data, column names help May also include rules about data: a zip code is part of an address that
- uniquely identifies a residence
- Useful for asking good questions about the data









## Data Model vs. Conceptual Model

- Data Model: raw data that has a specific data type (e.g. floats):
  - Temperature Example: [32.5, 54.0, -17.3] (floats)
- Conceptual Model: how we think about the data
  - Includes semantics, reasoning
  - Temperature Example:
    - Quantitative: [32.50, 54.00, -17.30]











## Data Model vs. Conceptual Model

- Data Model: raw data that has a specific data type (e.g. floats):
  - Temperature Example: [32.5, 54.0, -17.3] (floats)
- Conceptual Model: how we think about the data
  - Includes semantics, reasoning
  - Temperature Example:
    - Quantitative: [32.50, 54.00, -17.30]
    - Ordered: [warm, hot, cold]











## Data Model vs. Conceptual Model

- Data Model: raw data that has a specific data type (e.g. floats):
  - Temperature Example: [32.5, 54.0, -17.3] (floats)
- Conceptual Model: how we think about the data
  - Includes semantics, reasoning
  - Temperature Example:
    - Quantitative: [32.50, 54.00, -17.30]
    - Ordered: [warm, hot, cold]
    - Categorical: [not burned, burned, not burned]



















- Often, data in its original form isn't as useful as we would like
- Examples: Data about a basketball team's games









- Often, data in its original form isn't as useful as we would like
- Examples: Data about a basketball team's games
- Example 1: 1stHalfPoints, 2ndHalfPoints
  - More useful to know total number of points
  - Points = 1stHalfPoints + 2ndHalfPoints









- Often, data in its original form isn't as useful as we would like
- Examples: Data about a basketball team's games
- Example 1: 1stHalfPoints, 2ndHalfPoints
  - More useful to know total number of points
  - Points = 1stHalfPoints + 2ndHalfPoints
- Example 2: Points, OpponentPoints
  - Want to have a column indicating win/loss
  - Win = True if (Points > OpponentPoints) else False









- Often, data in its original form isn't as useful as we would like
- Examples: Data about a basketball team's games
- Example 1: 1stHalfPoints, 2ndHalfPoints
  - More useful to know total number of points
  - Points = 1stHalfPoints + 2ndHalfPoints
- Example 2: Points, OpponentPoints
  - Want to have a column indicating win/loss
  - Win = True if (Points > OpponentPoints) else False
- Example 3: Points
  - Want to have a column indicating how that point total ranks
  - Rank = index in sorted list of all Point values









#### 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 680/490, Spring 2021

# 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 680/490, Spring 2021

'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]:	obj3	In [29]: obj	4	In [30]: obj	3 + obj4
Out[28]:		Out[29]:		Out[30]:	
Ohio	35000	California	NaN	California	NaN
Oregon	16000	Ohio	35000	Ohio	70000
Texas	71000	Oregon	16000	Oregon	32000
Utah	5000	Texas	71000	Texas	142000
dtype: i	nt64	dtype: float	64	Utah	NaN
71	•	71	•	dtype: float	64
				[VV. N	AcKinney, Pyt











## Data Frame

- 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









## Chicago Food Inspections Exploration

- Based on David Beazley's PyData Chicago talk
- Data
- YouTube video: <u>https://www.youtube.com/watch?v=j6VSAsKAj98</u>
- Our in-class exploration:
  - Python can give answers fairly quickly
  - Data analysis questions:
    - What is information is available
    - Questions are interesting about this dataset
    - How to decide on good follow-up questions
    - What the computations mean







## Chicago Food Inspections Exploration







