Programming Principles in Python (CSCI 503)

Data

Dr. David Koop





CPU-Bound vs. I/O-Bound



Threading

- Threading address the I/O waits by letting separate pieces of a program run at the same time
- Threads run in the same process
- Threads share the same memory (and global variables)
- Operating system schedules threads; it can manage when each thread runs, e.g. round-robin scheduling
- When blocking for I/O, other threads can run













Python Threading Speed

- If I/O bound, threads work great be used by other threads
- Threads do not run simultaneously in standard Python, i.e. they cannot take advantage of multiple cores
- Use threads when code is I/O bound, otherwise no real speed-up plus some overhead for using threads

• If I/O bound, threads work great because time spent waiting can now be





Python and the GIL

- Solution for reference counting (used for garbage collection)
 Could add locking to every value/data structure, but with multiple locks
- Could add locking to every value/date
 comes possible deadlock
- Python instead has a Global Interpreter Lock (GIL) that must be acquired to execute any Python code
- This effectively makes Python single-threaded (faster execution)
- Python requires threads to give up GIL after certain amount of time
- Python 3 improved allocation of GIL to threads by not allowing a single CPUbound thread to hog it

D. Koop, CSCI 503/490, Fall 2021





5

Multiprocessing

- most cases
- Big win: can take advantage of multiple cores!

D. Koop, CSCI 503/490, Fall 2021

Multiple processes do not need to share the same memory, interact less Python makes the difference between processes and threads minimal in









Multiprocessing using concurrent.futures

• import concurrent.futures import multiprocessing as mp import time

def dummy(num): time.sleep(5) return num ** 2

results = executor.map(dummy, range(10))

• mp.get context('fork') changes from 'spawn' used by default in MacOS, works in notebook

with concurrent.futures.ProcessPoolExecutor(max workers=5, mp context=mp.get context('fork')) as executor:





asyncio

- Single event loop that controls when each task is run
- Tasks can be ready or waiting
- Tasks are **not interrupted** like they are with threading - Task controls when control goes back to the main event loop

 - Either waiting or complete
- Event loop keeps track of whether tasks are ready or waiting - Re-checks to see if new tasks are now ready - Picks the task that has been waiting the longest
- async and await keywords
- Requires support from libraries (e.g. aiohttp)









When to use threading, asyncio, or multiprocessing?

- If your code has a lot of I/O or Network usage:
 - If there is library support, use asyncio
 - Otherwise, multithreading is your best bet (lower overhead)
- If you have a GUI
 - Multithreading so your UI thread doesn't get locked up
- If your code is CPU bound:
 - You should use multiprocessing (if your machine has multiple cores)











Concurrency Comparison

Concurrency Type	Switching Decision	Number o Processor
Pre-emptive multitasking (threading)	The operating system decides when to switch tasks external to Python.	
Cooperative multitasking (asyncio)	The tasks decide when to give up control.	
Multiprocessing (multiprocessing)	The processes all run at the same time on different processors.	Man













pandas

- Contains high-level data structures and manipulation tools designed to make data analysis fast and easy in Python
- Built on top of NumPy
- Built with the following requirements:
 - Data structures with labeled axes (aligning data)
 - Support time series data
 - Do arithmetic operations that include metadata (labels)
 - Handle missing data
 - Add merge and relational operations

D. Koop, CSCI 503/490, Fall 2021



11



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 503/490, Fall 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
- Arithmetic operations **align**:

In [28]:	obj3	In [29]: obj	4	In [30]: obj	3 + obj4
Out[20]. 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
				dtype: float	:64
				[VV. N	AcKinney, Pyt









- A dictionary of Series (labels for each series) A spreadsheet with row keys (the index) and 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





15

DataFrame Constructor Inputs

Type

2D ndarray dict of arrays, lists, or tuples NumPy structured/record array dict of Series

dict of dicts

list of dicts or Series

List of lists or tuples Another DataFrame

NumPy MaskedArray

Notes

Treated as the "dict of arrays" case

Series" case.

DataFrame's column labels

Treated as the "2D ndarray" case

D. Koop, CSCI 503/490, Fall 2021

- A matrix of data, passing optional row and column labels
- Each sequence becomes a column in the DataFrame. All sequences must be the same length.
- Each value becomes a column. Indexes from each Series are unioned together to form the result's row index if no explicit index is passed.
- Each inner dict becomes a column. Keys are unioned to form the row index as in the "dict of
- Each item becomes a row in the DataFrame. Union of dict keys or Series indexes become the
- The DataFrame's indexes are used unless different ones are passed
- Like the "2D ndarray" case except masked values become NA/missing in the DataFrame result

[W. McKinney, Python for Data Analysis]











DataFrame Access and Manipulation

- df.values \rightarrow 2D NumPy array
- Accessing a column:
 - df["<column>"]
 - df.<column>
 - Both return Series
 - Dot syntax only works when the column is a valid identifier
- Assigning to a column:
 - df["<column>"] = <scalar> # all cells set to same value
 - df["<column>"] = <array> # values set in order
 - df["<column>"] = <series> # values set according to match between df and series indexes





df =	f = 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	= pd.read_csv('penguins_lter.csv')									
Column N	James	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	('penguins_l	ter.csv')							
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
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





	<pre>df = 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
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 503/490, Fall 2021







344 rows × 17 columns

D. Koop, CSCI 503/490, Fall 2021

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 503/490, Fall 2021

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 503/490, Fall 2021

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 [
							""
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







DataFrame Index

- Similar to index for Series
- Immutable
- Can be shared with multiple structures (DataFrames or Series)
- in operator works with: 'Ohio' in df.index
- \bullet Can choose new index column(s) with <code>set_index()</code>
- reindex creates a new object with the data conformed to new index
 - obj2 = obj.reindex(['a', 'b', 'c', 'd', 'e'])
 - can fill in missing values in different ways





Reading & Writing Data in Pandas

Format	Data Description
text	<u>CSV</u>
text	Fixed-Width Text File
text	<u>JSON</u>
text	HTML
text	Local clipboard
	MS Excel
binary	<u>OpenDocument</u>
binary	HDF5 Format
binary	Feather Format
binary	Parquet Format
binary	ORC Format
binary	<u>Msgpack</u>
binary	<u>Stata</u>
binary	<u>SAS</u>
binary	<u>SPSS</u>
binary	Python Pickle Format
SQL	SQL
SQL	Google BigQuery

D. Koop, CSCI 503/490, Fall 2021

Reader	Writer
read_csv	to_csv
read_fwf	
read_json	to_json
read_html	to_html
read_clipboard	to_clipboard
read_excel	to_excel
read_excel	
read_hdf	to_hdf
read_feather	to_feather
read_parquet	to_parquet
read_orc	
read_msgpack	to_msgpack
read_stata	to_stata
read_sas	
read_spss	
read_pickle	to_pickle
read_sql	to_sql
read_gbq	to_gbq

[https://pandas.pydata.org/pandas-docs/stable/user_guide/io.html]

read_csv

- Convenient method to read csv files
- Lots of different options to help get data into the desired format
- **Basic:** df = pd.read csv(fname)
- Parameters:

 - path: where to read the data from - sep (Or delimiter): the delimiter $(', ', '', '', ' \setminus t', ' \setminus s+')$
 - header: if None, no header
 - index col: which column to use as the row index - names: list of header names (e.g. if the file has no header)

 - skiprows: number of list of lines to skip

Writing CSV data with pandas

- Basic: df.to csv(<fname>)
- Change delimiter with sep kwarg:
 - df.to csv('example.dsv', sep='|')
- Change missing value representation - df.to csv('example.dsv', na rep='NULL')
- Don't write row or column labels:
 - df.to csv('example.csv', index=False, header=False)
- Series may also be written to csv

Documentation

- pandas <u>documentation</u> is pretty good

D. Koop, CSCI 503/490, Fall 2021

Lots of recipes on stackoverflow for particular data manipulations/queries

Food Inspections Example

