Programming Principles in Python (CSCI 503)

Data

Dr. David Koop





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









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, Spring 2021

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









- 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





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









	studyName	Sample	Species	Region	Island	Stage	Individual	Clutch	Date	Culmen Length
		Number		-			ID	Completion	Egg	(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_cs	v('penguins_l	ter.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











344 rows × 17 columns

D. Koop, CSCI 503, Spring 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









	df =	<pre>pd.read_csv</pre>	('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



D. Koop, CSCI 503, Spring 2021











344 rows × 17 columns

D. Koop, CSCI 503, Spring 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, Spring 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, Spring 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 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









Indexing

- Same as with NumPy arrays but can use Series's index labels
- Slicing with labels: NumPy is **exclusive**, Pandas is **inclusive**!
 - s = Series(np.arange(4)) s[0:2] # gives two values like numpy
 - s = Series(np.arange(4), index=['a', 'b', 'c', 'd'])s['a':'c'] # gives three values, not two!
- Obtaining data subsets
 - []: get columns by label
 - loc: get rows/cols by label
 - iloc: get rows/cols by position (integer index)
- For single cells (scalars), also have at and iat





Filtering

- Same as with numpy arrays but allows use of column-based criteria
 - data [data < 5] = 0
 - data[data['three'] > 5]
- Multiple criteria, use &, \mid , and \sim ; remember parentheses!
 - data[(data['three'] > 5) & (data['two'] < 10)]

• data < 5 \rightarrow boolean data frame, can be used to select specific elements







Filtering

df[df['Culmen Length (mm)'] > 40]

	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









Filtering

df[df['Culmen Length (mm)'] > 40]

	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









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, Spring 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, Spring 2021

Lots of recipes on stackoverflow for particular data manipulations/queries





<u>Assignment 7</u>

- Downloading and unarchiving files
- File system manipulation
- Threading
- Basic Data Manipulation
- Due Friday





Derived Data

- Create new columns from existing columns
 r["PctFail"] = r['Fail'] / r['Total']
- Note that operations are computed in a vectorized manner
- Similarities to functional paradigm (map/filter):
 - specify the operation once
 - no loops
 - interpreted as an operation on the entire column

columns r['Total'] in a vectorized manne map/filter):





Aggregation

- Descriptive statistics
 - df['Culmen Length (mm)'].mean()
 - .median()
 - .describe()
 - .count()
 - .min(), .max()
- Also general methods
 - .sum()
 - .product()





Split-Apply-Combine



D. Koop, CSCI 503, Spring 2021

[W. McKinney, Python for Data Analysis]











Split-Apply-Combine

- Similar to Map (split+apply) Reduce (combine) paradigm
- The Pattern:
 - 1. Split the data by some grouping variable
 - 2. Apply some function to each group independently
 - 3. Combine the data into some output dataset
- The apply step is usually one of:
 - Aggregate
 - Transform
 - Filter







In Pandas

- groupby method creates a GroupBy object
- groupby doesn't actually compute anything until there is an apply/aggregate step or we wish to examine the groups
- Choose keys (columns) to group by
- size() is the count of each group
- Other aggregates also work







Split-Apply-Combine

- df.groupby('Island')[['Culmen Length (mm)',
- df.groupby('Island').agg({'Culmen Length (mm)': 'mean',
- df.groupby('Island').agg(cul length=('Culmen Length (mm)', 'mean'), cul depth=('Culmen Depth (mm)', 'mean'))

Island		
Biscoe	45.257485	15.874850
Dream	44.167742	18.344355
Torgersen	38.950980	18.429412

D. Koop, CSCI 503, Spring 2021

```
'Culmen Depth (mm)']].mean()
 'Culmen Depth (mm) ': 'mean'})
```

cul_length cul_depth







Different Data Layouts

	treatm	ienta t	reatmentb	-			
John Smith			2	-			
Jane Doe		16	11			4 4	
Mary Johnson		3	1		name	trt	r
		_		-	John Smith	a	
	nitial D)ata			Jane Doe	a	
					Mary Johnson	a	
					John Smith	b	
					Jane Doe	b	
John Sn	nith Ja	ane Doe	Mary Joh	nson	Mary Johnson	b	
nenta		16		3		$) \rightarrow + \rightarrow$	
nenth	2	11		1	I IQY L	Jala	

	trea	atmenta	treatmentb					
John S	Smith		2					
Jane I	Doe	16	11		-			
Mary	Johnson	3	1		_	name	trt	result
						John Smith	\mathbf{a}	
	Initia	l Data				Jane Doe	a	16
						Mary Johnson	a	3
						John Smith	b	2
						Jane Doe	b	11
	John Smith	Jane Doe	Mary Joh	nson	-	Mary Johnson	b	1
treatmenta		16		3		Tidv F	$) \rightarrow + \rightarrow$	
treatmentb	2	11		1		hay L	າລເລ	

Transpose









Problem: Variables stored in both rows & columns

Mexico Weather, Global Historical Climatology Network

id	year	month	element	d1	d2	d3	d4	d5	d6	d7	d8
MX17004	2010	1	tmax								
MX17004	2010	1	tmin								
MX17004	2010	2	tmax		27.3	24.1					
MX17004	2010	2	tmin		14.4	14.4					
MX17004	2010	3	tmax					32.1			
MX17004	2010	3	tmin					14.2			
MX17004	2010	4	tmax								
MX17004	2010	4	tmin								
MX17004	2010	5	tmax								
MX17004	2010	5	tmin								

D. Koop, CSCI 503, Spring 2021







22

Problem: Variables stored in both rows & columns

Mexico Weather, Global Historical Climatology Network

id	year	month	element	d1	d2	d3	d4	d5	d6	d7	d8
MX17004	2010	1	tmax								
MX17004	2010	1	tmin								
MX17004	2010	2	tmax		27.3	24.1					
MX17004	2010	2	tmin		14.4	14.4					
MX17004	2010	3	tmax					32.1			
MX17004	2010	3	tmin					14.2			
MX17004	2010	4	tmax								
MX17004	2010	4	tmin								
MX17004	2010	5	tmax								
MX17004	2010	5	tmin								

Variable in columns: day; Variable in rows: tmax/tmin











Solution: Melting + Pivot

id	date	element	value	id	date	tmax	tmin
MX17004	2010-01-30	tmax	27.8	MX17004	2010-01-30	27.8	14.5
MX17004	2010-01-30	tmin	14.5	MX17004	2010-02-02	27.3	14.4
MX17004	2010-02-02	tmax	27.3	MX17004	2010-02-03	24.1	14.4
MX17004	2010-02-02	tmin	14.4	MX17004	2010-02-11	29.7	13.4
MX17004	2010-02-03	tmax	24.1	MX17004	2010-02-23	29.9	10.7
MX17004	2010-02-03	tmin	14.4	MX17004	2010-03-05	32.1	14.2
MX17004	2010-02-11	tmax	29.7	MX17004	2010-03-10	34.5	16.8
MX17004	2010-02-11	tmin	13.4	MX17004	2010-03-16	31.1	17.6
MX17004	2010-02-23	tmax	29.9	MX17004	2010-04-27	36.3	16.7
MX17004	2010-02-23	tmin	10.7	MX17004	2010-05-27	33.2	18.2

(a) Molten data

D. Koop, CSCI 503, Spring 2021

(b) Tidy data

[H. Wickham, 2014]









Melt

Want to keep each observation separate (tidy), aka pivot_longer

	location	Temperature	Jan-2010	Feb-2010	Ma
0	CityA	Predict	30	45	
1	CityB	Actual	32	43	

df.melt(id vars=["location", "Temperature"], var name="Date", value name="Value")

	location	Temperature	Date	Value
0	CityA	Predict	Jan-2010	30
1	CityB	Actual	Jan-2010	32
2	CityA	Predict	Feb-2010	45
3	CityB	Actual	Feb-2010	43
4	CityA	Predict	Mar-2010	24
5	CityB	Actual	Mar-2010	22











Pivot

- "wide" format (aka pivot_wider)
- Long format: column names are data values...
- Wide format: more like spreadsheet format
- Example:

	date	item	value	
0	1959-03-31	realgdp	2710.349	
1	1959-03-31	infl	0.000	
2	1959-03-31	unemp	5.800	
3	1959-06-30	realgdp	2778.801	
4	1959-06-30	infl	2.340	
5	1959-06-30	unemp	5.100	
6	1959-09-30	realgdp	2775.488	
7	1959-09-30	infl	2.740	
8	1959-09-30	unemp	5.300	
9	1959-12-31	realgdp	2785.204	

Sometimes, we have data that is given in "long" format and we would like

```
.pivot('date', 'item', 'value')
```

item	infl	realgdp	unemp
date			
1959-03-31	0.00	2710.349	5.8
1959-06-30	2.34	2778.801	5.1
1959-09-30	2.74	2775.488	5.3
1959-12-31	0.27	2785.204	5.6
1960-03-31	2.31	2847.699	5.2

[W. McKinney, Python for Data Analysis]









Reshaping Data

- Reshape/pivoting are fundamental operations
- Can have a nested index in pandas
- 3rd) and associated representative rankings
- Could write this in different ways:

number	one	two	three
state			
Ohio	0	1	2
Colorado	3	4	5

stat Ohio

Coloi

D. Koop, CSCI 503, Spring 2021

• Example: Congressional Districts (Ohio's 1st, 2nd, 3rd, Colorado's 1st, 2nd,

state	Ohio	Colorado
number		
one	0	3
two	1	4
three	2	5

е	number	
	one	0
	two	1
	three	2
rado	one	3
	two	4
	three	5









Reshaping Data

- Reshape/pivoting are fundamental operations
- Can have a nested index in pandas
- 3rd) and associated representative rankings
- Could write this in different ways:

number	one	two	three
state			
Ohio	0	1	2
Colorado	3	4	5

stat Ohio

```
MultiIndex
```

Colo

D. Koop, CSCI 503, Spring 2021

• Example: Congressional Districts (Ohio's 1st, 2nd, 3rd, Colorado's 1st, 2nd,

state	Ohio	Colorado
number		
one	0	3
two	1	4
three	2	5

е	number	
	one	0
	two	1
	three	2
rado	one	3
	two	4
	three	5









Stack and Unstack

- stack: pivots from the columns into rows (may produce a Series!)
- unstack: pivots from rows into columns
- unstacking may add missing data
- stacking filters out missing data (unless dropna=False)
- level one two three number



Color

D. Koop, CSCI 503, Spring 2021

• can unstack at a different level by passing it (e.g. 0), defaults to innermost

	Т		state number	Ohio	Colorado		
			one	0	3		
			two	1	4		
	1_		three	2	5		
ac	K		/				
j	number						
	one	0	IING	tack	$-(\bigcirc)$		
	two	1	und				
	three	2					
cado	one	3					
	two	4					
	three	5		[\	V. McKinn	ey, Python for Data A	na
						Northern Illinois Univers	sity







String Methods

- Can do many of the same methods used for single strings on entire columns • Requires .str prefix before calling the method
- violations.value.str.strip().str.split(' Comments:') Also helps when extracting from a list - comments.str[1]









String Methods

Argument	Description
count	Return the number of non-overlapping
endswith	Returns True if string ends with suf
startswith	Returns True if string starts with pr
join	Use string as delimiter for concatenation
index	Return position of first character in su
find	Return position of first character of <i>fi</i> if not found.
rfind	Return position of first character of <i>la</i>
replace	Replace occurrences of string with ar
strip, rstrip, lstrip	Trim whitespace, including newlines; for each element.
solit	Break string into list of substrings usi
lower	Convert alphabet characters to lower
upper	Convert alphabet characters to upper
casefold	Convert characters to lowercase, and common comparable form.
ljust,	Left justify or right justify, respective
rjust	character) to return a string with a m

D. Koop, CSCI 503, Spring 2021

ng occurrences of substring in the string.

ffix.

refix.

ating a sequence of other strings.

ubstring if found in the string; raises ValueError if not found.

first occurrence of substring in the string; like index, but returns -1

last occurrence of substring in the string; returns –1 if not found.

nother string.

```
; equivalent to x.strip() (and rstrip, lstrip, respectively)
```

ing passed delimiter.

rcase.

rcase.

convert any region-specific variable character combinations to a

ely; pad opposite side of string with spaces (or some other fill ninimum width.











Support for Datetime

- Python has datetime library to support dates and times pandas has a Timestamp data type that functions somewhat similarly
- Pandas can convert timestamps
 - pd.to_datetime: versatile, can often guess format
- Like string methods, also a . dt accessor for datetime methods/properties
- With a timestamp, filtering based on datetimes becomes easier
 - df[df['Inspection Date'] > '2021']







Food Inspections Example





