Programming Principles in Python (CSCI 503/490)

Data

Dr. David Koop





:	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	onIslandStageIndividual IDCersTorgersenAdult, 1 Egg StageN1A1ersTorgersenAdult, 1 Egg StageN1A2ersTorgersenAdult, 1 Egg StageN2A1ersTorgersenAdult, 1 Egg StageN2A2ersTorgersenAdult, 1 Egg StageN3A1ersTorgersenAdult, 1 Egg StageN3A1ersBiscoeAdult, 1 Egg StageN39A1ersBiscoeAdult, 1 Egg StageN39A1ersBiscoeAdult, 1 Egg StageN43A1ersBiscoeAdult, 1 Egg StageN43A1	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	Image: Line Date Egg Culmen Length (mm) Yes 11/11/07 39.1 Yes 11/11/07 39.5 Yes 11/16/07 40.3 Yes 11/16/07 NaN Yes 11/16/07 NaN Yes 11/16/07 36.7 Mo 12/1/09 NaN Yes 11/22/09 46.8 Yes 11/22/09 50.4 Yes 11/22/09 45.2 Yes 11/22/09 49.9	
		340	PAL0910	121	Gentoo penguin (Pygoscelis papua)	Anvers	Biscoe	ndStageIndividual IDClutch CompletionDate EggCulme CompletionienAdult, 1 Egg StageN1A1Yes11/11/07ienAdult, 1 Egg StageN1A2Yes11/11/07ienAdult, 1 Egg StageN2A1Yes11/16/07ienAdult, 1 Egg StageN2A2Yes11/16/07ienAdult, 1 Egg StageN2A2Yes11/16/07ienAdult, 1 Egg StageN3A1Yes11/16/07ienAdult, 1 Egg StageN3A1Yes11/16/07ienAdult, 1 Egg StageN3A2No12/1/09ienAdult, 1 Egg StageN39A1Yes11/22/09ienAdult, 1 Egg StageN39A2Yes11/22/09ienAdult, 1 Egg StageN39A2Yes11/22/09ienAdult, 1 Egg StageN43A1Yes11/22/09	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 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	If = pd.read_csv('penguins_lter.csv')SstudyNameSample NumberSpeciesRegionIslandStageIndividual IDClutch Completion0PAL07081Adelie Penguin (Pygoscelis adeliae)AnversTorgersenAdult, 1 Egg StageN1A1Yes11/1PAL07082Adelie Penguin (Pygoscelis adeliae)AnversTorgersenAdult, 1 Egg StageN1A2Yes11/2PAL07083Adelie Penguin (Pygoscelis adeliae)AnversTorgersenAdult, 1 Egg StageN2A1Yes11/3PAL07084Adelie Penguin (Pygoscelis adeliae)AnversTorgersenAdult, 1 Egg StageN2A2Yes11/4PAL07085Adelie Penguin (Pygoscelis adeliae)AnversTorgersenAdult, 1 Egg StageN2A2Yes11/4PAL07085Adelie Penguin (Pygoscelis adeliae)AnversTorgersenAdult, 1 Egg StageN3A1Yes11/4PAL07085Adelie Penguin (Pygoscelis papua)AnversTorgersenAdult, 1 Egg StageN3A1Yes11/5Mutricold12Gentoo penguin (Pygoscelis papua)AnversBiscoeAdult, 1 Egg StageN3A2No124PAL0910121Gentoo penguin (Pygoscelis papua)AnversBiscoeAdult, 1 Egg StageN3A2Yes11/2344PAL0910122Gentoo penguin (Py	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





	df =	pd.read_csv	<pre>('penguins_l</pre>	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 503/490, Spring 2023







344 rows × 17 columns

D. Koop, CSCI 503/490, Spring 2023

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, Spring 2023

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, Spring 2023

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









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

D. Koop, CSCI 503/490, Spring 2023





3

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

D. Koop, CSCI 503/490, Spring 2023





3

Reading and Writing Data

- Reading:
 - df = pd.read_csv(fname)
- Writing
 - df.to_csv(fname)
- Many options also possible on both
 - sep: the separator (defaults to comma)
 - skiprows: when reading, number of list of lines to skip
 - index: set to None when writing if unimportant
- Also methods for other formats (json, parquet, sql)
- Methods are read_* and to_*

ן





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











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









Avoid inplace



D. Koop, CSCI 503/490, Spring 2023

47 / 48



Northern Illinois University









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 503/490, Spring 2023

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

cul_length cul_depth







Split-Apply-Combine



D. Koop, CSCI 503/490, Spring 2023

[W. McKinney, Python for Data Analysis]









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 Smith	1		2	_			
Jane Doe		16	11				
Mary John	son	3	1		name	trt	result
				-	John Smith	\mathbf{a}	
	Initia	Data			Jane Doe	a	16
					Mary Johnson	a	3
					John Smith	b	2
					Jane Doe	b	11
Joh	n Smith	Jane Doe	Mary Joh	nson	Mary Johnson	b	1
treatmenta		16		3	Tidvr	$) \rightarrow + \rightarrow$	
treatmentb	2	11		1	TIQY L	Jala	
	Ŧ						

Iranspose







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 503/490, Spring 2023

(b) Tidy data

[H. Wickham, 2014]



<u>Assignment 7</u>

- Musical Artists Datasets
- Downloading and uncompressing files
- Finding files using OS libraries
- Load per-artist numpy arrays, each saved in the .npy format
- Store per-month dataframes, each in a csv file
- Issue with r.e.m..npy





Food Inspections Example





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-overlappi
endswith	Returns True if string ends with suf
startswith	Returns True if string starts with pr
join	Use string as delimiter for concatena
index	Return position of first character in s
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,	Trim whitespace, including newlines,
rstrip,	for each element.
lstrip	
split	Break string into list of substrings us
lower	Convert alphabet characters to lower
иррег	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 503/490, Spring 2023

ing 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.

```
s; 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']





Method chaining in pandas

- Tom Augspurger's <u>post</u>
- <u>Effective Pandas</u> book by Matt Harrison
- Functions written for chaining, and pipe allows custom functions

```
• def read(fp):
  df = (pd.read csv(fp)
          .rename(columns=str.lower)
          .drop('unnamed: 36', axis=1)
          .pipe(extract city name)
          .pipe(time to datetime, ['dep time', 'arr time',
          .assign(fl date=lambda x: pd.to datetime(x['fl date']),
                  dest=lambda x: pd.Categorical(x['dest']),
```

'crs arr time', 'crs dep time']) origin=lambda x: pd.Categorical(x['origin']), tail num=lambda x: pd.Categorical(x['tail num']), unique carrier=lambda x: pd.Categorical(x['unique carrier']), cancellation code=lambda x: pd.Categorical(x['cancellation_code'])))





Example: Inspect Intermediate Results

• def csnap(df, fn=lambda x: x.shape, msg=None):

** ** **

if msg: print(msg) display(fn(df)) return df

- wine.pipe(csnap) # display data frame
 - .rename(columns={"color intensity": "ci"})

 - .pipe(csnap) # display data frame

D. Koop, CSCI 503/490, Spring 2023

...

""" Custom Help function to print things in method chaining. Returns back the df to further use in chaining.

.assign(color filter=lambda x: np.where(x.hue > 1, 1, 0))





