Advanced Data Management (CSCI 490/680)

Data Cleaning

Dr. David Koop





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 490/680, Spring 2020

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]









Types of arguments for readers

- Indexing: choose a column to index the data, get column names from file or user
- Type inference and data conversion: automatic or user-defined
- Datetime parsing: can combine information from multiple columns
- Iterating: deal with very large files
- Unclean Data: skip rows (e.g. comments) or deal with formatted numbers (e.g. 1,000,345)









JavaScript Object Notation (JSON)

- A format for web data
- Looks very similar to python dictionaries and lists
- Example:
 - { "name": "Wes", "places lived": ["United States", "Spain", "Germany"], "pet": null, "siblings": [{"name": "Scott", "age": 25, "pet": "Zuko"}, {"name": "Katie", "age": 33, "pet": "Cisco"}] }
- Only contains literals (no variables) but allows null
- Values: strings, arrays, dictionaries, numbers, booleans, or null
 - Dictionary keys must be strings
 - Quotation marks help differentiate string or numeric values





Binary Formats

- CSV, JSON, and XML are all text formats
- What is a binary format?
- Pickle: Python's built-in serialization
- HDF5: Library for storing large scientific data
 - Hierarchical Data Format, supports compression
 - Interfaces in C, Java, MATLAB, etc.
 - Use pd. HDFStore to access
 - Shortcuts: read hdf/to hdf, need to specify object
- Excel: need to specify sheet when a spreadsheet has multiple sheets
 - pd.ExcelFile Of pd.read excel









Databases

- Similar syntax from other database
 Oracle, etc.)
- SQLAlchemy: Python package that abstracts away differences between different database systems
- SQLAIchemy gives support for reading queries to data frame:
 - import sqlalchemy as sqla db = sqla.create_engine('sqlite:///mydata.sqlite') pd.read_sql('select * from test', db)

• Similar syntax from other database systems (MySQL, Microsoft SQL Server,





6

Dirty Data





Geolocation Errors

- address
- world of trouble" [Washington Post, 2016]



D. Koop, CSCI 490/680, Spring 2020

Maxmind helps companies determine where users are located based on IP

"How a quiet Kansas home wound up with 600 million IP addresses and a









Numeric Outliers

12 13 14 21 22 26 33 35 36 37 39 42 45 47 54 57 61 68 450 ages of employees (US)

400

300



D. Koop, CSCI 490/680, Spring 2020

median 37 * mean 58.52632 * variance 9252.041 *





Northern Illinois University







Dirty Data: Data Scientist's View

- Combination of:
 - Statistician's View: data has non-ideal samples for model
 - Database Expert's View: missing data, corrupted data
 - Domain Expert's View: data doesn't pass the smell test
- All of the views present problems with the data
- The goal may dictate the solutions:
 - Median value: don't worry too much about crazy outliers - Generally, aggregation is less susceptible by numeric errors - Be careful, the data may be correct...

D. Koop, CSCI 490/680, Spring 2020







10

Be careful how you detect dirty data

- The appearance of a hole in the earth's ozone layer over Antarctica, first malfunctioning.
 - National Center for Atmospheric Research



D. Koop, CSCI 490/680, Spring 2020

detected in 1976, was so unexpected that scientists didn't pay attention to what their instruments were telling them; they thought their instruments were











Data Wrangling

- better analyzed
- Data cleaning: getting rid of inaccurate data
- Data transformations: changing the data from one representation to another • Data reshaping: reorganizing the data
- Data merging: combining two datasets

• Data wrangling: transform raw data to a more meaningful format that can be





12

<u>Assignment 2</u>

- Similar to Assignment 1, now with pandas
- Part 5:
 - CS 680 \rightarrow Required
 - CS 490 → Extra Credit
- Due Friday, Feb. 7

D. Koop, CSCI 490/680, Spring 2020

m	onth	1	2	3	4	5	6	7	8	9	10	11	12	
	year													
	1851	0	0	0	0	0	1	2	1	1	1	0	0	
	1852	0	0	0	0	0	0	0	1	3	1	0	0	
	1853	0	0	0	0	0	0	0	3	4	1	0	0	
	1854	0	0	0	0	0	1	0	1	2	1	0	0	
	1855	0	0	0	0	0	0	0	4	1	0	0	0	
	2015	0	0	0	0	1	1	1	3	5	0	1	0	
	2016	1	0	0	0	1	2	0	5	5	1	1	0	
	2017	0	0	0	1	0	2	3	4	4	3	1	0	
	2018	0	0	0	0	1	0	2	3	7	3	0	0	
	All	4	1	1	6	38	125	171	448	623	353	89	14	1

169 rows × 13 columns







Wrangler: Interactive Visual Specification of Data Transformation Scripts

S. Kandel, A. Paepcke, J. Hellerstein, J. Heer





Wrangler

- Data cleaning takes a lot of time and human effort
- "Tedium is the message"
- Repeating this process on multiple data sets is even worse!
- Solution:
 - interactive interface (mixed-initiative)
 - transformation language with natural language "translations"
 - suggestions + "programming by demonstration"





Previous Work: Potter's Wheel

- V. Raman and J. Hellerstein, 2001
- Defines structure extractions for identifying fields
- Defines transformations on the data
- Allows user interaction





Potter's Wheel: Structure Extraction



# Structures	Final Structure Chason
	I'mai Suuciule Chosen
Enumerated	(Punc = Punctuation)
5	Integer
5	IspellWord
12	AllCapsWord
9	Int Punc(/) Int Punc(/) Int
5	Capitalized Word
5	Int(len 2) Punc(:) Int(len 2)
9	Double Punc('.') Double
5	<i>Punc(") IspellWord Punc(\)</i>
4	ξ^*
3	AllCapsWord(HTTP)
6	Punc(/) Double(1.0)
	[V. Raman and J. Hellerste
	Northern Illinois Univ









Potter's Wheel: Transforms

Transform			
Format	$\phi(R,i,f)$	=	$\{(a_1,\ldots,a_{i-1})\}$
Add	$\alpha(R,x)$	=	$\{(a_1,\ldots,a_n,x)\}$
Drop	$\pi(R,i)$	=	$\{(a_1,\ldots,a_{i-1})\}$
Сору	$\kappa((a_1,\ldots,a_n),i)$	=	$\{(a_1,\ldots,a_n,a_n)\}$
Merge	$\mu((a_1,\ldots,a_n),i,j,glue)$	=	$\{(a_1,\ldots,a_{i-1})\}$
Split	$\omega((a_1,\ldots,a_n),i,\text{splitter})$) =	$\{(a_1,\ldots,a_{i-1})\}$
Divide	$\delta((a_1,\ldots,a_n),i,pred)$	=	$\{(a_1,\ldots,a_{i-1})\}$
			$\{(a_1,\ldots,a_{i-1})\}$
Fold	$\lambda(R, i_1, i_2, \dots i_k)$	—	$\{(a_1,\ldots,a_{i_1}-$
			(a_1,\ldots,a_n) (
Select	$\sigma(R, \text{pred})$	=	$\{(a_1,\ldots,a_n)\mid$

Notation: R is a relation with n columns. i, j are column indices and a_i represents the value of a column in a row. x and glue are values. f is a function mapping values to values. $x \oplus y$ concatenates x and y. splitter is a position in a string or a regular expression, left(x, splitter) is the left part of x after splitting by splitter. pred is a function returning a boolean.

D. Koop, CSCI 490/680, Spring 2020

Definition $, a_{i+1}, \ldots, a_n, f(a_i)) \mid (a_1, \ldots, a_n) \in R \}$ $x) \mid (a_1, \ldots, a_n) \in R \}$ $\{a_{i+1}, \ldots, a_n\} \mid (a_1, \ldots, a_n) \in R\}$ $a_i) \mid (a_1, \ldots, a_n) \in R \}$ $\{a_{i+1}, \ldots, a_{j-1}, a_{j+1}, \ldots, a_n, a_i \oplus \text{glue} \oplus a_j) \mid (a_1, \ldots, a_n) \in R\}$ $a_{i+1}, \ldots, a_n, \text{left}(a_i, \text{splitter}), \text{right}(a_i, \text{splitter})) \mid (a_1, \ldots, a_n) \in R$ $\{a_i, a_{i+1}, \ldots, a_n, a_i, \text{null}\} \mid (a_1, \ldots, a_n) \in R \land \text{pred}(a_i)\} \cup$ $\{a_{i+1},\ldots,a_n, \text{ null},a_i\} \mid (a_1,\ldots,a_n) \in R \land \neg \text{pred}(a_i)\}$ $a_{i_1+1}, a_{i_1+1}, \ldots, a_{i_2-1}, a_{i_2+1}, \ldots, a_{i_k-1}, a_{i_k+1}, \ldots, a_n, a_{i_l})$ $\in R \land 1 \leq l \leq k$ $(a_1,\ldots,a_n) \in R \wedge \operatorname{pred}((a_1,\ldots,a_n))$

[V. Raman and J. Hellerstein, 2001]





18

Potter's Wheel: Example

		Stewart,Bob	
Anna	Davis		'(.*),(
		Dole,Jerry	
Joan	Marsh		

Bob	Stewart	2 Mer
Anna	Davis	
Jerry	Dole	
Joan	Marsh	

D. Koop, CSCI 490/680, Spring 2020



[V. Raman and J. Hellerstein, 2001]











Potter's Wheel: Inferring Structure from Examples

Example Values Split By User	Inferred Structure	Comments
(is user specified split position)		
Taylor, Jane , \$52,072 Blair, John , \$73,238 Tony Smith , \$1,00,533	$(<\xi^* > < `,` Money >)$	Parsing is doable despite no good de limiter. A <i>regular expression</i> domain can infer a structure of \$[0-9,]* for last component.
MAA to SIN JFK to SFO LAX - ORD SEA / OAK	$(< len 3 identifier > < \xi^* > < len 3 identifier >)$	Parsing is possible despite multiple delimiters.
321 Blake #7 , Berkeley , CA 94720 719 MLK Road , Fremont , CA 95743	(<number <math="">\xi^* > < ',' word> <',' (2 letter word) (5 letter integer)>)</number>	Parsing is easy because of consistent delimiter.

D. Koop, CSCI 490/680, Spring 2020

[V. Raman and J. Hellerstein, 2001]











Wrangler Transformation Language

- Based on Potter's Wheel
- Map: Delete, Extract, Cut, Split, Update
- Lookup/join: Use external data (e.g. from zipcode \rightarrow state) Reshape: Fold and Unfold (aka pivot)
- Positional: Fill and lag
- Sorting, aggregation, key generation, schema transforms

D. Koop, CSCI 490/680, Spring 2020





21

Interface

- Automated Transformation Suggestions
- Editable Natural Langua
 - Transform Script Import Export Fill Bangladesh by copying value Split data repeatedly on **newline** into **rows** above

Fill Bangladesh by ✓ copying interpolating values from above

- Fill Bangladesh by averaging t values from above
- Visual Transformation Pl
- Transformation History



split	# split1	split2	# split3	🗰 split4
	2004	2004	2004	2003
ATE	Participation Rate 2004	Mean SAT I Verbal	Mean SAT I Math	Participation Ro
v York	87	497	510	82
nnecticut	85	515	515	84
ssachusetts	85	518	523	82
v Jersey	83	501	514	85
v Hampshire	80	522	521	75
- - •	77	489	476	77
ine	76	505	501	70
ınsylvania	74	501	502	73
Laware	73	500	499	73
orgia	73	494	493	66
split	# fold	Abc fold1	# value	
v York	2004	Participation Rate 2004	87	
v York	2004	Mean SAT I Verbal	497	
v York	2004	Mean SAT I Math	510	
v York	2003	Participation Rate 2003	82	
v York	2003	Mean SAT I Verbal	496	
v York	2003	Mean SAT I Math	510	
nnecticut	2004	Participation Rate 2004	85	
nnecticut	2004	Mean SAT I Verbal	515	
nnecticut	2004	Mean SAT I Math	515	
nnecticut	2003	Participation Rate 2003	84	
nnecticut	2003	Mean SAT I Verbal	512	
	2002	NALL CAT T MALL	F 1 A	











Automation from past actions

- Infer parameter sets from user interaction
- Generating transforms
- Ranking and ordering transformations:
 - Based on user preferences, difficulty, and corpus frequency
 - Sort transforms by type and diversify suggestions

Reported crime in Alabama

(a)

(b)	<i>before:</i> <i>selection:</i> <i>after</i> :	{ 'in', ' '} { 'Alabama' } Ø	'Alabama' \rightarrow {'Alabama', word 'in' \rightarrow {'in', word, lowercase} '' \rightarrow {''}
(c)	<i>before:</i> <i>selection:</i> <i>after</i> :	{(' '), ('in', ' ') {('Alabama'), (Ø	<pre>(word, ``), (lowercase, ``)} (word)}</pre>
(d)	$\{(), (`Alabata) \\ \{(`, '), (), ()\} \\ \{(`, '), (`Alabata) \\ \{(`, '), (`, '), (`, '), (`Alabata) \\ \{(`, '), (`$	ma'),()} bama'),()} { },()} {),()} {'Alabama'),()}	<pre>{(),(word),()} {(word, ``),(),()} {(word, ``),('Alabama'),()} {(word, ``),(word),()} {(lowercase, ``),(),()} {(lowercase, ``),('Alabama'),() {(lowercase, ``),(word),()}</pre>
(\mathbf{n})	[(lowarcas	a ('Alabama	()) (10 z (10 hom))

), (Alaballia), () \rightarrow /[a-L]+ (Alaballia)/ (\mathbf{U}) *illowercuse*,













Evaluation

- Compare with Excel
- Tests:
 - Extract text from a single string entry
 - Fill in missing values with estimates
 - Reshape tables
- Allowed users to ask questions about Excel, not Wrangler
- helpful
- Complaint: No manual fallback, make implications of user choices more obvious for users

Found significant effect of tool and users found previews and suggestions









Task Completion Times



D. Koop, CSCI 490/680, Spring 2020





25

IM

TR.

ех

SU

ext

ext

ext

Sou

abc

6 Ca

31 adt

32 adt

33 adt

34 adt

hts in Prediction

Partially underlined Figure 12 qualified retrieval

TYPE	ITEM	COLOR	SIZ
	P. I <u>KE</u>	GREEN	

D. Koc

- equality operators: \neq , >, >=, <, <=. If no inequality c used as a prefix, equality is implied. The symbol $\neq c$ placed by \neg or $\neg =$.
- Partially underlined qualified retrieval. Print the green start with the letter I. This is found in Figure 12. The not underlined, and it is a constant. Therefore, the sys all the green items that start with the letter I. The use tially underline at the beginning, middle or end of a wo tence, or a paragraph, as in the example, XPAY, whi find a word, a sentence or a paragraph such that som that sentence or paragraph there exist the letters PA. example element can be blank, then a word, a sente

paragraph that starts or ends with the letters PA also qu

The partial underline feature is useful if an entry is a se text and the user wishes to search to find all examples tain a special word or root. If, for example, the query entries with the word Texas, the formulation of this qu TEXAS Y.

Update suggestions when given more information

Qualified retrieval using links. Print all the green iter the toy department. This is shown in Figure 13.2015 this user displays both the TYPE table and the SALES table









Data Wrangling Tasks

- Unboxing: Discovery & Assessment: What's in there? (types, distribution) • Structuring: Restructure data (table, nested data, pivot tables)
- Cleaning: does data match expectations (often involves user)
- Enriching & Blending: Adding new data
- Optimizing & Publishing: Structure for storage or visualization









Differences with Extract-Transform-Load (ETL)

- ETL:
 - Who: IT Professionals
 - Why: Create static data pipeline
 - What: Structured data
 - Where: Data centers
- "Modern Data Preparation":
 - Who: Analysts
 - Why: Solve problems by designing recipes to use data
 - What: Original, custom data blended with other data
 - Where: Cloud, desktop





Northern Illinois University





Trifacta Wrangler







Paper Critique

- Foofah: Transforming Data By Example, Z. Jin et al., 2017
- Due Tuesday **before** class, submit via Blackboard
- Read the paper
- Look up references if necessary
- Keep track of things you are confused by or that seem problematic
- Write a few sentences summarizing the paper's contribution
- Write more sentences discussing the paper and what you think the paper does well or doesn't do well at
- For this response, compare/contrast with Wrangler/Trifacta
- Length: 1/2-1 page







Data Cleaning in pandas









Handling Missing Data

- Filtering out missing data:
 - Can choose rows or columns
- Filling in missing data:
 - with a default value
 - with an interpolated value
- In pandas:

Argument	Description
dropna	Filter axis labels based on whether values for much missing data to tolerate.
fillna	Fill in missing data with some value or using
isnull	Return boolean values indicating which value
notnull	Negation of isnull.

D. Koop, CSCI 490/680, Spring 2020

r each label have missing data, with varying thresholds for how

an interpolation method such as 'ffill' or 'bfill'. es are missing/NA.

[W. McKinney, Python for Data Analysis]









Filling in missing data

• fillna arguments:

Argument	Description
value	Scalar value or dict-like object
method	Interpolation; by default 'ff
axis	Axis to fill on; default axis=
inplace	Modify the calling object wit
limit	For forward and backward fil

D. Koop, CSCI 490/680, Spring 2020

- ct to use to fill missing values
- fill' if function called with no other arguments
- =0
- hout producing a copy
- lling, maximum number of consecutive periods to fill

[W. McKinney, Python for Data Analysis]











Filtering and Cleaning Data

- Find duplicates
 - first instance is **not marked** as a duplicate
- Remove duplicates:
 - drop duplicates: drops all rows where duplicated is True - keep: which value to keep (first or last)
- Can pass specific columns to check for duplicates, e.g. check only key column

- duplicated: returns boolean Series indicating whether row is a duplicate -





Changing Data

- Convert strings to upper/lower case
- Convert Fahrenheit temperatures to Celsius
- Create a new column based on another column

```
In [56]: lowercased
Out[56]:
                                meat_to_animal = {
0
           bacon
     pulled pork
                                  'bacon': 'pig',
1
2
           bacon
                                  'pulled pork': 'pig',
3
        pastrami
                                  'pastrami': 'cow',
     corned beef
4
                                  'corned beef': 'cow',
           bacon
5
                                  'honey ham': 'pig',
        pastrami
6
                                  'nova lox': 'salmon'
       honey ham
        nova lox
Name: food, dtype: object
```

D. Koop, CSCI 490/680, Spring 2020

<pre>In [57]: data['animal'] = lowercased.map(meat_to_animal)</pre>				
In [58]: data Out[58]:				
	food	ounces	animal	
0	bacon	4.0	pig	
1	pulled pork	3.0	pig	
2	bacon	12.0	pig	
3	Pastrami	6.0	COW	
4	corned beef	7.5	COW	
5	Bacon	8.0	pig	
6	pastrami	3.0	COW	
7	honey ham	5.0	pig	
8	nova lox	6.0	salmon	

[W. McKinney, Python for Data Analysis]









Replacing Values

- fillna is a special case
- What if –999 in our dataset was identified as a missing value?

In [In [62]	
Out[61]:		Out [62]
Ο	1.0	Ο
1	-999.0	1
2	2.0	2
3	-999.0	3
4	-1000.0	4 -100
5	3.0	5
dtyp	e: float64	dtype:

Can pass list of values or dictionary to change different values

```
: data.replace(-999, np.nan)
```

```
1.0
```

```
NaN
```

```
2.0
```

```
NaN
```

```
00.0
```

```
3.0
```

```
float64
```







Clamping Values

In [93]: data.descr	<pre>ibe()</pre>		
Out[93]:			
	0	1	2	
count	1000.000000	1000.000000	1000.000000	1000.00000
mean	0.049091	0.026112	-0.002544	-0.051827
std	0.996947	1.007458	0.995232	0.998311
min	-3.645860	-3.184377	-3.745356	-3.428254
25%	-0.599807	-0.612162	-0.687373	-0.747478
50 %	0.047101	-0.013609	-0.022158	-0.088274
75 %	0.756646	0.695298	0.699046	0.623333
max	2.653656	3.525865	2.735527	3.366626

D. Koop, CSCI 490/680, Spring 2020

Values above or below a specified thresholds are set to a max/min value

```
In [97]: data[np.abs(data) > 3] = np.sign(data) * 3
In [98]: data.describe()
Out[98]:
                                            2
                                                          3
                 0
                               1
       1000.000000
                    1000.000000
                                  1000.000000
                                               1000.000000
count
                                                 -0.051765
          0.050286
                       0.025567
                                    -0.001399
mean
          0.992920
                                     0.991414
                                                  0.995761
std
                       1.004214
         -3.000000
                                    -3.000000
min
                       -3.000000
                                                 -3.000000
25%
         -0.599807
                       -0.612162
                                    -0.687373
                                                 -0.747478
          0.047101
                       -0.013609
                                    -0.022158
                                                  -0.088274
50%
                       0.695298
75%
                                     0.699046
          0.756646
                                                  0.623331
          2.653656
                       3.000000
                                     2.735527
                                                  3.000000
max
```







Computing Indicator Values

- Useful for machine learning
- Want to take possible values and map them to 0-1 indicators
- Example:

•

a b c

0 0 1 0

1 0 1 0

2 1 0 0

3 0 0 1

4 1 0 0

5 0 1 0

```
In [110]: pd.get_dummies(df['key'])
Out[110]:
```

```
In [109]: df = pd.DataFrame({'key': ['b', 'b', 'a', 'c', 'a', 'b'],
                            'data1': range(6)})
```







String Transformation

- One of the reasons for Python's popularity is string/text processing
- split (<delimiter>): break a string into pieces:
 - -s = "12, 13, 14"slist = s.split(',') # ["12", "13", " 14"]
- <delimiter>.join([<str>]): join several strings by a delimiter
 - ":".join(slist) # "12:13: 14"
- strip(): remove leading and trailing whitespace
 - [p.strip() for p in slist] # ["12", "13", "14"]







String Transformation

replace (<from>, <to>): change substrings to another substring

- s.replace(',', ':') # "12:13: 14"

- upper()/lower(): casing
 - "AbCd".upper () # "ABCD"
 - "AbCd".lower() # "abcd"

D. Koop, CSCI 490/680, Spring 2020

substrings to another substring





String Transformations

- index (<str>): find where a substring first occurs (Error if not found)
 - -s = "12, 13, 14"
 - s.index(',') # 2
 - s.index(':') # ValueError raised
- find (<str>): same as index but -1 if not found
 - s.find(',') # 2 s.find(':') # -1
- startswith()/endswith(): boolean checks for string occurrence
 - s.startswith("1") # True s.endswith("5") # False





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 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 490/680, Spring 2020

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.







Regular Expressions

- AKA regex
- A syntax to better specify how to decompose strings
- Look for patterns rather than specific characters
- "31" in "The last day of December is 12/31/2020."
- May work for some questions but now suppose I have other lines like: "The last day of September is 9/30/2020."
- ...and I want to find dates that look like:
- <numbers>/<numbers>/<numbers>
- Cannot search for every combination!
- d+/d+/d+







Regular Expressions

- Character classes:
 - $\ \ d = digits$
 - $\ s = spaces$
 - $\w = word character [a-zA-Z0-9]$
 - [a-z] = lowercase letters (square brackets indicate a set of chars)
- Repeating characters or patterns
 - + = one or more (any number)
 - * = zero or more (any number)
 - -? = zero or one
 - $\{ < number > \} = a \text{ specific number (or range) of occurrences}$





Regular Expressions in Python

- import re
- re.search(<pattern>, <str_to_check>)
 - Returns None if no match, information about the match otherwise
- Capturing information about what is in a string \rightarrow parentheses
- $(\d+)/\d+/\d+$ will **capture** information about the month
- match = re.search('(\d+)/\d+/\d+','12/31/2016')
 if match:
 match.group() # 12
- re.findall(<pattern>, <str_to_check>)
 - Finds all matches in the string, search only finds the first match
- Can pass in flags to alter methods: e.g. re.IGNORECASE





Pandas String Methods

- to the entire series
- Fast (vectorized) on whole columns or datasets
- USe .str.<method name>
- .str is important!
 - data = pd.Series({'Dave': 'dave@google.com',

 - 'Wes': np.nan})

data.str.contains('gmail') data.str.split('@').str[1] data.str[-3:]

• Any column or series can have the string methods (e.g. replace, split) applied

```
'Steve': 'steve@gmail.com',
'Rob': 'rob@gmail.com',
```





