Advanced Data Management (CSCI 490/680)

Data Wrangling

Dr. David Koop





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 490/680, Spring 2020

Contains high-level data structures and manipulation tools designed to make









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 490/680, Spring 2020

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





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





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









Indexing Data Frames

- Brackets can be ambiguous:
 - df['Address']
 - df[0:4]
- - df.loc[:, 'Address']
 - df.iloc[0:4,:]
- Putting them together:
 - df.iloc[0:4,:].loc[:,'Address']
 - df.loc[df.index[0:4], 'Address']

D. Koop, CSCI 490/680, Spring 2020

• .loc and .iloc require more code (always row and column), but are clearer









Sorting by Value (sort_values)

- sort values method on series - obj.sort values()
- first)
- sort values on DataFrame:
 - df.sort values (<list-of-columns>)
 - df.sort values(by=['a', 'b'])
 - Can also use axis=1 to sort by index labels

D. Koop, CSCI 490/680, Spring 2020



• Missing values (NaN) are at the end by default (na position controls, can be





Unique Values and Value Counts

- unique returns an array with only the unique values (no index) - s = Series(['c','a','d','a','a','b','b','c','c']) s.unique() # array(['c', 'a', 'd', 'b'])
- Data Frames use drop duplicates
- value counts returns a Series with index frequencies:
 - s.value counts() # Series({'c': 3,'a': 3,'b': 2,'d': 1})









Statistics

- sum: column sums (axis=1 gives sums over rows)
- missing values are excluded unless the whole slice is NaN
- idxmax, idxmin are like argmax, argmin (return index)
- describe: shortcut for easy stats!

In [204]:	<pre>df.describe()</pre>	In [2
Out[204]:			

	one	two	In [2
count	3.000000	2.000000	Out[2
mean	3.083333	-2.900000	count
std	3.493685	2.262742	uniqu
min	0.750000	-4.500000	top
25%	1.075000	-3.700000	freq
50%	1.400000	-2.900000	dtype
75%	4.250000	-2.100000	
max	7.100000	-1.300000	

```
205]: obj = Series(['a', 'a', 'b', 'c'] * 4)
In [206]: obj.describe()
Out[206]:
count
         16
unique
           3
           а
           8
dtype: object
```







<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	nonth	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







D. Koop, CSCI 490/680, Spring 2020

Data Formats





Comma-separated values (CSV) Format

- Comma is a field separator, newlines denote records
 - a,b,c,d,message 1,2,3,4,hello 5, 6, 7, 8, world 9,10,11,12,foo
- May have a header (a, b, c, d, message), but not required
- No type information: we do not know what the columns are (numbers, strings, floating point, etc.)
 - Default: just keep everything as a string
- Type inference: Figure out the type to make each column based on values What about commas in a value? \rightarrow double quotes





Delimiter-separated Values

- Comma is a **delimiter**, specifies boundary between fields
- Could be a tab, pipe (1), or perhaps spaces instead
- All of these follow similar styles to CSV





Fixed-width Format

- Old school
- Each field gets a certain number of spots in the file
- Example:

- id8141	360.242940	149
id1594	444.953632	166
id1849	364.136849	183
id1230	413.836124	184
id1948	502.953953	173

• Specify exact character ranges for each field, e.g. 0-6 is the id

- .910199
- .985655
- .628767
- .375703
- .237159

- 11950.7
- 11788.4
- 11806.2
- 11916.8
- 12468.3





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)











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 (', ', '', '\t', '\t', '\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





More read_csv/read_tables arguments

Argument	Description
skiprows	Number of rows at beginning of fi
na_values	Sequence of values to replace with
comment	Character(s) to split comments off
parse_dates	Attempt to parse data to dateti Otherwise can specify a list of colu combine multiple columns togethe
keep_date_col	If joining columns to parse date, k
converters	Dict containing column number of function f to all values in the 'fo
dayfirst	When parsing potentially ambigue 2012); False by default.
date_parser	Function to use to parse dates.
ΠΓΟ₩S	Number of rows to read from begi
iterator	Return a TextParser object for
chunksize	For iteration, size of file chunks.

D. Koop, CSCI 490/680, Spring 2020

ile to ignore or list of row numbers (starting from 0) to skip. h NA.

- the end of lines.
- ime; False by default. If True, will attempt to parse all columns. Imn numbers or name to parse. If element of list is tuple or list, will er and parse to date (e.g., if date/time split across two columns).
- keep the joined columns; False by default.
- name mapping to functions (e.g., { 'foo': f} would apply the oo' column).
- ous dates, treat as international format (e.g., 7/6/2012 -> June 7,

inning of file.

^r reading file piecemeal.

[W. McKinney, Python for Data Analysis]









Chunked Reads

- With very large files, we may not want to read the entire file
- Why?
 - Time
 - Want to understand part of data before processing all of it
- Reading only a few rows:
 - df = pd.read csv('example.d
- Reading chunks:
 - Get an iterator that returns the next chunk of the file
 - chunker = pd.read csv('example.csv', chunksize=1000)
 - for piece in chunker: process data (piece)





Python csv module

- Also, can read csv files outside of pandas using csv module
 - import csv with open ('persons of concern.csv', 'r') as f: for i in range(3): next(f) reader = csv.reader(f)
 - records = [r for r in reader] # r is a list

• Or

- import csv with open('persons of concern.csv', 'r') as f: for i in range(3): next(f) reader = csv.DictReader(f) records = [r for r in reader] # r is a dict









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

D. Koop, CSCI 490/680, Spring 2020





eXtensible Markup Language (XML)

- Older, self-describing format with nesting; each field has tags
- Example:
 - <INDICATOR> <INDICATOR SEQ>373889</INDICATOR SEQ> <PARENT SEQ></PARENT SEQ> <AGENCY NAME>Metro-North Railroad</AGENCY NAME> <INDICATOR NAME>Escalator Avail.</INDICATOR NAME> <PERIOD YEAR>2011/PERIOD YEAR> <PERIOD MONTH>12/PERIOD MONTH> <CATEGORY>Service Indicators</CATEGORY> <FREQUENCY>M</FREQUENCY> <YTD TARGET>97.00</YTD TARGET> </INDICATOR>
- Top element is the **root**









XML

- No built-in method
- Use lxml library (also can use ElementTree)
- from lxml import objectify path = 'datasets/mta perf/Performance MNR.xml' parsed = objectify.parse(open(path)) root = parsed.getroot() data = []skip fields = ['PARENT SEQ', 'INDICATOR SEQ', 'DESIRED CHANGE', 'DECIMAL PLACES'] for elt in root.INDICATOR: el data = $\{\}$ for child in elt.getchildren(): if child.tag in skip fields: continue el data[child.tag] = child.pyval data.append(el data) perf = pd.DataFrame(data)

D. Koop, CSCI 490/680, Spring 2020

[W. McKinney, Python for Data Analysis]



Northern Illinois University







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









What is the problem with reading this data?

```
• [{"name": "Wes",
   "places lived": ["United States", "Spain", "Germany"],
   "pet": null,
   "siblings":
      {"name": "Scott", "age": 25, "pet": "Zuko"},
      {"name": "Katie", "age": 33, "pet": "Cisco"}]
  }
  {"name": "Nia",
   "address": {"street": "143 Main",
               "city": "New York",
                "state": "New York"},
   "pet": "Fido",
   "siblings":
      {"name": "Jacques", "age": 15, "pet": "Fido"}]
  },
 ...
```









Reading JSON data

- Python has a built-in json module
 - with open ('example.json') as f: data = json.load(f)
 - Can also load/dump to strings:
 - json.loads, json.dumps
- Pandas has read json, to json methods









JSON Orientation

- produced by to json() with a corresponding orient value. The set of possible orients is:
 - split: dict like {index -> [index], columns -> [columns], data \rightarrow [values] }
 - records: list like [{column -> value}, ..., {column -> value}]
 - index: dict like {index -> {column -> value}}
 - columns: dict like {column -> {index -> value}}
 - values: just the values array

Indication of expected JSON string format. Compatible JSON strings can be











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
 - Interfaces in C, Java, MATLAB, etc.
 - Supports compression
 - 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









Databases

Dir	n_Date				
8	Id	~0 1	1		
	Date				
	Day				
	Day_of_Week		Ι.		
	Month			Fa	ct_Sa
	Month_Name		<u> </u>		Date_
	Quarter				Store_
	Quarter_Name				Produc
	Year				Units_













Databases

- features
 - links between tables via foreign keys
 - SQL to create, store, and query data
- sqlite3 is a simple database with built-in support in python
- Python has a database API which lets you access most database systems through a common API.

Relational databases are similar to multiple data frames but have many more









Python DBAPI Example

import sqlite3

con = sqlite3.connect('mydata.sqlite') con.execute(query) con.commit() # Insert some data data = [('Atlanta', 'Georgia', 1.25, 6), ('Tallahassee', 'Florida', 2.6, 3), ('Sacramento', 'California', 1.7, 5)] stmt = "INSERT INTO test VALUES(?, ?, ?, ?)" con.executemany(stmt, data) con.commit()

query = """CREATE TABLE test(a VARCHAR(20), b VARCHAR(20), c REAL, d INTEGER);"""

[W. McKinney, Python for Data Analysis]









Databases

- Oracle, etc.)
- SQLAIchemy: 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,







What if data isn't correct/trustworthy/in the right format?







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









This takes a lot of time!



D. Koop, CSCI 490/680, Spring 2020

What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets; 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%

[CrowdFlower Data Science Report, 2016]







...and it isn't the most fun thing to do



D. Koop, CSCI 490/680, Spring 2020

What's the least enjoyable part of data science?

- Building training sets: 10%
- Cleaning and organizing data: 57%
- Collecting data sets: 21%
- Mining data for patterns: 3%
- Refining algorithms: 4%
- Other: 5%

[CrowdFlower Data Science Report, 2016]









Dirty Data: Statistician's View

- Some process produces the data
- Want a model but have non-ideal samples:
 - Distortion: some samples corrupted by a process
 - Selection bias: likelihood of a sample depends on its value
 - Left and right censorship: users come and go from scrutiny
 - Dependence: samples are not independent (e.g. social networks)
- You can add/augment models for different problems, but cannot model everything
- Trade-off between accuracy and simplicity











Dirty Data: Database Expert's View

- Got a dataset
- Some values are missing, corrupted, wrong, duplicated
- Results are absolute (relational model)
- Better answers come from improving the quality of values in the dataset

D. Koop, CSCI 490/680, Spring 2020







Dirty Data: Domain Expert's View

- Data doesn't look right
- Answer doesn't look right
- What happened?
- Domain experts carry an implicit model of the data they test against You don't always need to be a domain expert to do this
- - Can a person run 50 miles an hour?
 - Can a mountain on Earth be 50,000 feet above sea level?
 - Use common sense

D. Koop, CSCI 490/680, Spring 2020







Dirty Data: Data Scientist's View

- Combination of the previous three views
- 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...







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











Where does dirty data originate?

- Source data is bad, e.g. person entered it incorrectly
- Transformations corrupt the data, e.g. certain values processed incorrectly due to a software bug
- Integration of different datasets causes problems
- Error propagation: one error is magnified









Types of Dirty Data Problems

- Separator Issues: e.g. CSV without respecting double quotes -12, 13, "Doe, John", 45
- Naming Conventions: NYC VS. New York
- Missing required fields, e.g. key
- Different representations: 2 vs. two
- Redundant records: may be exactly the same or have some overlap
- Formatting issues: 2017-11-07 vs. 07/11/2017 vs. 11/07/2017

D. Koop, CSCI 490/680, Spring 2020

• Truncated data: "Janice Keihanaikukauakahihuliheekahaunaele" becomes "Janice Keihanaikukauakahihuliheek" on Hawaii license









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







Data Cleaning







Wrangler: Interactive Visual Specification of Data Transformation Scripts

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





Data Wrangler Demo

<u>http://vis.stanford.edu/wrangler/app/</u>

Transform Script	Impo	rt Export
Split data repeatedly or rows	on newline into	
Split split repeatedly of	on ','	
Promote row 0 to head	ler	
Text Columns Rows	Table	Clear
Delete row 7		
Delete empty rows		1
		1
Fill row / by copying v	alues from above	1
		1
		1-
D. Koop, CSCI 490/680, Spring 2020)	1



	Tear Year	Property_crime_rate
0	Reported crime in Alabama	
1		
2	2004	4029.3
3	2005	3900
4	2006	3937
5	2007	3974.9
5	2008	4081.9
7		
8	Reported crime in Alaska	
9		
)	2004	3370.9
1	2005	3615
2	2006	3582
3	2007	3373.9
4	2008	2928.3



