Advanced Data Management (CSCI 640/490)

Data Wrangling

Dr. David Koop





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





2

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]

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





Arithmetic

- Add, subtract, multiply, and divide are element-wise like numpy
- ...but use labels to align
- ...and missing labels lead to NaN (not a number) values

| In [28]: obj3 | | In [29]: obj | In [29]: obj4 | | 3 + obj4 |
|---------------|-------|--------------|---------------|--------------|----------|
| Out[28] | • | Out[29]: | | Out[30]: | |
| 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: | int64 | dtype: float | 64 | Utah | NaN |
| | | | | dtype: float | 64 |

- also have .add, .subtract, ... that allow fill value argument
- obj3.add(obj4, fill value=0)









Mutating Dataframes

- assign allows new columns to be created, returns "new" dataframe
 - df2 = df.assign(Total=df.Points1 + df.Points2)
- More reusable:
- If you have columns that are not proper identifiers, can use **kwargs
 - df.Points2})

D. Koop, CSCI 640/490, Spring 2023

- df2 = df.assign(Total=lambda df: df.Points1 + df.Points2) - df2 = df.assign(**{"Total Points": lambda df: df.Points1 +





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 640/490, Spring 2023



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







String Methods

- Can manipulate columns of strings - Use the .str modifier
- Most string and regex operations are available
- Examples:
 - df.first name.str.startswith("Jo")
 - df.name.str.split(' ')







DuckDB

- SQL syntax with extras
 - read csv auto
 - similar string, datetime, array functions to pandas







lbis

- More Pythonic interface to database or dataframe systems
- Uses DuckDB as default backend, can be configured to use others
- Syntax aligns better with SQL, potentially clearer
 - select
 - filter
 - mutate
 - group_by
 - order_by
 - unnest

D. Koop, CSCI 640/490, Spring 2023

e or dataframe systems can be configured to use others entially clearer





<u>Assignment 2</u>

- Assignment 1 Questions with pandas, DuckDB, and Ibis
- CS 640 students do all, CS 490 do pandas & DuckDB (lbis is EC)
- Can work by framework or by query
- Most questions can be answered with a single statement... but that statement can take a while to write
 - Read documentation
 - Check hints





Test 1

- Monday, Feb. 27
- In-class, 9:30-10:45am
- Format:
 - Multiple Choice
 - Free Response
- Information will be posted online





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]: | | | _ |

top

freq

| | one | two | |
|-------|----------|-----------|--|
| count | 3.000000 | 2.000000 | |
| mean | 3.083333 | -2.900000 | |
| std | 3.493685 | 2.262742 | |
| min | 0.750000 | -4.500000 | |
| 25% | 1.075000 | -3.700000 | |
| 50% | 1.400000 | -2.900000 | |
| 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
```





Statistics

| Method | Description |
|----------------|------------------------------|
| count | Number of non-NA values |
| describe | Compute set of summary sta |
| min, max | Compute minimum and max |
| argmin, argmax | Compute index locations (in |
| idxmin, idxmax | Compute index values at wh |
| quantile | Compute sample quantile ra |
| sum | Sum of values |
| mean | Mean of values |
| median | Arithmetic median (50% qua |
| mad | Mean absolute deviation fro |
| var | Sample variance of values |
| std | Sample standard deviation of |
| skew | Sample skewness (3rd mom |
| kurt | Sample kurtosis (4th momer |
| cumsum | Cumulative sum of values |
| cummin, cummax | Cumulative minimum or ma |
| cumprod | Cumulative product of value |
| diff | Compute 1st arithmetic diffe |
| pct_change | Compute percent changes |
| | |

D. Koop, CSCI 640/490, Spring 2023

- atistics for Series or each DataFrame column
- iximum values
- ntegers) at which minimum or maximum value obtained, respectively
- hich minimum or maximum value obtained, respectively
- anging from 0 to 1
- uantile) of values
- om mean value
- of values
- nent) of values
- ent) of values
- aximum of values, respectively
- es
- ference (useful for time series)

[W. McKinney, Python for Data Analysis]





14

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})





Handling Missing Data

| Argument | Description |
|----------|--|
| dropna | Filter axis labels based on whether values for missing data to tolerate. |
| fillna | Fill in missing data with some value or using |
| isnull | Return like-type object containing boolean |
| notnull | Negation of isnull. |

D. Koop, CSCI 640/490, Spring 2023

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

g an interpolation method such as 'ffill' or 'bfill'. values indicating which values are missing / NA.

[W. McKinney, Python for Data Analysis]











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 640/490, Spring 2023

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 640/490, Spring 2023

median 37 * mean 58.52632 * variance 9252.041 *





Northern Illinois University









This takes a lot of time!



D. Koop, CSCI 640/490, Spring 2023

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 640/490, Spring 2023

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











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











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 640/490, Spring 2023

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 640/490, Spring 2023

• 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





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"







Your Critique/Questions







Example Critique

- operations and demonstrates them on demand
- other path to take. In addition, a user has to have some idea of the correction. Perhaps a more example-based strategy could help.

• Summary: Wrangler tackles data wrangling tasks by combining a language for specifying operations with an interface allowing users to specify the types of changes they are interested; the system can then generate suggested

 Critique: The suggestions may lead to states that a user cannot recover from easily. Suppose a suggestion looks like it works well, but a user later realizes was incorrect. They can backtrack, but it's often unclear where to and which constructs of the language in order to edit parameters. Without a good idea of the impact of the parameters, the work may become as tedious as manual







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 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 | $\lambda(R, i_1, i_2, \dots i_k)$ | = | $\{(a_1,\ldots,a_{i_1-1})\}$ |
| | | | (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 640/490, Spring 2023

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

[V. Raman and J. Hellerstein, 2001]











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 640/490, Spring 2023



[V. Raman and J. Hellerstein, 2001]











Potter's Wheel: Inferring Structure from Examples

| Example Values Split By User (is user specified split position) | Inferred Structure | Comments | |
|---|--|--|--|
| 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 640/490, Spring 2023

[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





Interface

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

н н

Table

Promote row 0 to header

Delete rows 0,1

Fill row 0 by copying

values from the left

Clea

Columns Rows

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

D. Koop, CSCI 640/490, Spring 2023

 2 New
 3 Con
 4 Mas
 5 New 6 New 7 D.C 8 Mai 9 Pen 10 De **11** Geo . New New New New Con Con Con

| splitsplit1split2split302004200420041STATEParticipation Rate 2004Mean SAT I VerbalMean SAT I Math | 2003 |
|---|------------------|
| 02004200420041STATEParticipation Rate 2004Mean SAT I VerbalMean SAT I Math | 2003 |
| I STATE Participation Rate 2004 Mean SAT I Verbal Mean SAT I Math | |
| | Participation Ro |
| 2 New York 87 497 510 | 82 |
| 3 Connecticut 85 515 515 | 84 |
| 4 Massachusetts 85 518 523 | 82 |
| 5 New Jersey 83 501 514 | 85 |
| 6 New Hampshire 80 522 521 | 75 |
| 7 D.C. 77 489 476 | 77 |
| 8 Maine 76 505 501 | 70 |
| 9 Pennsylvania 74 501 502 | 73 |
| l O Delaware 73 500 499 | 73 |
| 11 Georgia 73 494 493 | 66 |
| | |
| split # fold At fold1 # value | |
| New York 2004 Participation Rate 2004 87 | |
| New York 2004 Mean SAT I Verbal 497 | |
| New York 2004 Mean SAT I Math 510 | |
| New York 2003 Participation Rate 2003 82 | |
| New York 2003 Mean SAT I Verbal 496 | |
| New York 2003 Mean SAT I Math 510 | |
| Connecticut 2004 2004 Participation Rate 2004 85 | |
| Connecticut 2004 Mean SAT I Verbal 515 | |
| Connecticut 2004 Mean SAT I Math 515 | |
| Connecticut 2003 Participation Rate 2003 84 | |
| Connecticut 2003 Mean SAT I Verbal 512 | |
| | |









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) | $ \{(), (`Alabat) \\ \{(`, '), (), ()\} \\ \{(`, '), (`Alabat) \\ \{(`, '), (`, '), (`, '), (`Alabat) \\ \{(`, '), (`, '$ | ma'),()} bama'),()} { },()} {),()} {(Alabama'),()} | <pre>{(),(word),()} {(word, ``),(),()} {(word, ``),('Alabama'),()} {(word, ``),(word),()} {(lowercase, ``),(),()} {(lowercase, ``),('Alabama'),() {(lowercase, ``),(word),()}</pre> |
| (\mathbf{n}) | [(lowaroas | a (') ('Alabama | $()) $ $(10 \ z \ (10 \ mo))$ |

(e) {(*lowercase*, ''), ('Alabama'), ()} \rightarrow /[a-z]+ (Alabama)/







v 43

Data Wrangler Demo

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

| Transform Script Impor | Transform Script Import Export | |
|--|--------------------------------|-----|
| Split data repeatedly on newline into rows | | |
| Split split repeatedly on ',' | | |
| Promote row 0 to header | | |
| Text Columns Rows Table | Clear | ļ |
| | | |
| | | |
| Delete row 7 | | |
| | | |
| Delete empty rows | | 10 |
| Fill row 7 by conving values from above | | 1 |
| Thirtow 7 by copying values non above | | 1. |
| | | 1: |
| | and the second | 14 |
| D. Koop, CSCI 640/490, Spring 2023 | | L : |



| | 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 |
| 1 | 2008 | 2928.3 |





44