# Advanced Data Management (CSCI 680/490)

## Data Wrangling

Dr. David Koop

# DataFrame Access and Manipulation

- `df.values` → `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`

# 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

- `s = Series(np.arange(4.), index=[4,3,2,1])`
- `s[3]`
- `s.loc[3]`
- `s.iloc[3]`
- `s2 = pd.Series(np.arange(4), index=['a','b','c','d'])`
- `s2[3]`

# Filtering

- Same as with numpy arrays but allows use of column-based criteria
  - `data[data < 5] = 0`
  - `data[data['three'] > 5]`
- `data < 5` → 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
Out[28]:
Ohio        35000
Oregon      16000
Texas       71000
Utah         5000
dtype: int64
```

```
In [29]: obj4
Out[29]:
California      NaN
Ohio         35000
Oregon       16000
Texas        71000
dtype: float64
```

```
In [30]: obj3 + obj4
Out[30]:
California       NaN
Ohio           70000
Oregon         32000
Texas         142000
Utah             NaN
dtype: float64
```

- also have `.add`, `.subtract`, … that allow `fill_value` argument

- `obj3.add(obj4, fill_value=0)`

# Arithmetic between DataFrames and Series

- Broadcasting: e.g. apply single row operation across all rows

- Example:

```
In [148]: frame          In [149]: series          In [150]: frame - series
Out[148]:                Out[149]:                 Out[150]:
        b   d   e         b    0                           b   d   e
Utah    0   1   2         d    1                   Utah    0   0   0
Ohio    3   4   5         e    2                   Ohio    3   3   3
Texas   6   7   8         Name: Utah, dtype: float64 Texas  6   6   6
Oregon  9  10  11                                   Oregon  9   9   9
```

- To broadcast over **columns**, use methods (`.add`, …)

```
In [154]: frame          In [155]: series3         In [156]: frame.sub(series3, axis=0)
Out[154]:                Out[155]:                 Out[156]:
        b   d   e         Utah      1                       b   d   e
Utah    0   1   2         Ohio      4               Utah   -1   0   1
Ohio    3   4   5         Texas     7               Ohio   -1   0   1
Texas   6   7   8         Oregon   10               Texas  -1   0   1
Oregon  9  10  11         Name: d, dtype: float64   Oregon -1   0   1
```

# Sorting by Index (sort_index)

- Sort by index (lexicographical):

```
In [168]: obj = Series(range(4), index=['d', 'a', 'b', 'c'])

In [169]: obj.sort_index()
Out[169]:
a    1
b    2
c    3
d    0
dtype: int64
```

- DataFrame sorting:

```
In [170]: frame = DataFrame(np.arange(8).reshape((2, 4)), index=['three', 'one'],
   .....:                   columns=['d', 'a', 'b', 'c'])
```

```
In [171]: frame.sort_index()          In [172]: frame.sort_index(axis=1)
Out[171]:                             Out[172]:
       d  a  b  c                            a  b  c  d
one    4  5  6  7                     three  1  2  3  0
three  0  1  2  3                     one    5  6  7  4
```

- axis controls sort rows (0) vs. sort columns (1)

# Sorting by Value (sort_values)

- `sort_values` method on series
  - `obj.sort_values()`
- Missing values (`NaN`) are at the end by default (`na_position` controls, can be 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

# Assignment 2

- Basically the same as Assignment 1, now with pandas and duckdb
- Can either do each task at the same time (one in pandas, one in duckdb), or all tasks in pandas then all tasks in duckdb

# Test 1

- Next Wednesday, Feb. 23

- In-class, 3:30-4:45pm in PM 153

- Format:

  - Multiple Choice

  - Free Response

- Information 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]: df.describe()
Out[204]:
            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
```

```
In [205]: obj = Series(['a', 'a', 'b', 'c'] * 4)

In [206]: obj.describe()
Out[206]:
count       16
unique       3
top          a
freq         8
dtype: object
```

# Statistics

| Method | Description |
| --- | --- |
| count | Number of non-NA values |
| describe | Compute set of summary statistics for Series or each DataFrame column |
| min, max | Compute minimum and maximum values |
| argmin, argmax | Compute index locations (integers) at which minimum or maximum value obtained, respectively |
| idxmin, idxmax | Compute index values at which minimum or maximum value obtained, respectively |
| quantile | Compute sample quantile ranging from 0 to 1 |
| sum | Sum of values |
| mean | Mean of values |
| median | Arithmetic median (50% quantile) of values |
| mad | Mean absolute deviation from mean value |
| var | Sample variance of values |
| std | Sample standard deviation of values |
| skew | Sample skewness (3rd moment) of values |
| kurt | Sample kurtosis (4th moment) of values |
| cumsum | Cumulative sum of values |
| cummin, cummax | Cumulative minimum or maximum of values, respectively |
| cumprod | Cumulative product of values |
| diff | Compute 1st arithmetic difference (useful for time series) |
| pct_change | Compute percent changes |

[W. McKinney, Python for Data Analysis]

Northern Illinois University

# 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 each label have missing data, with varying thresholds for how much missing data to tolerate. |
| `fillna` | Fill in missing data with some value or using an interpolation method such as `'ffill'` or `'bfill'`. |
| `isnull` | Return like-type object containing boolean values indicating which values are missing / NA. |
| `notnull` | Negation of `isnull`. |

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

# Dirty Data



[Flickr]

# Geolocation Errors

- Maxmind helps companies determine where users are located based on IP address

- "How a quiet Kansas home wound up with 600 million IP addresses and a world of trouble" [Washington Post, 2016]

# 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)



- median 37
- mean 58.52632
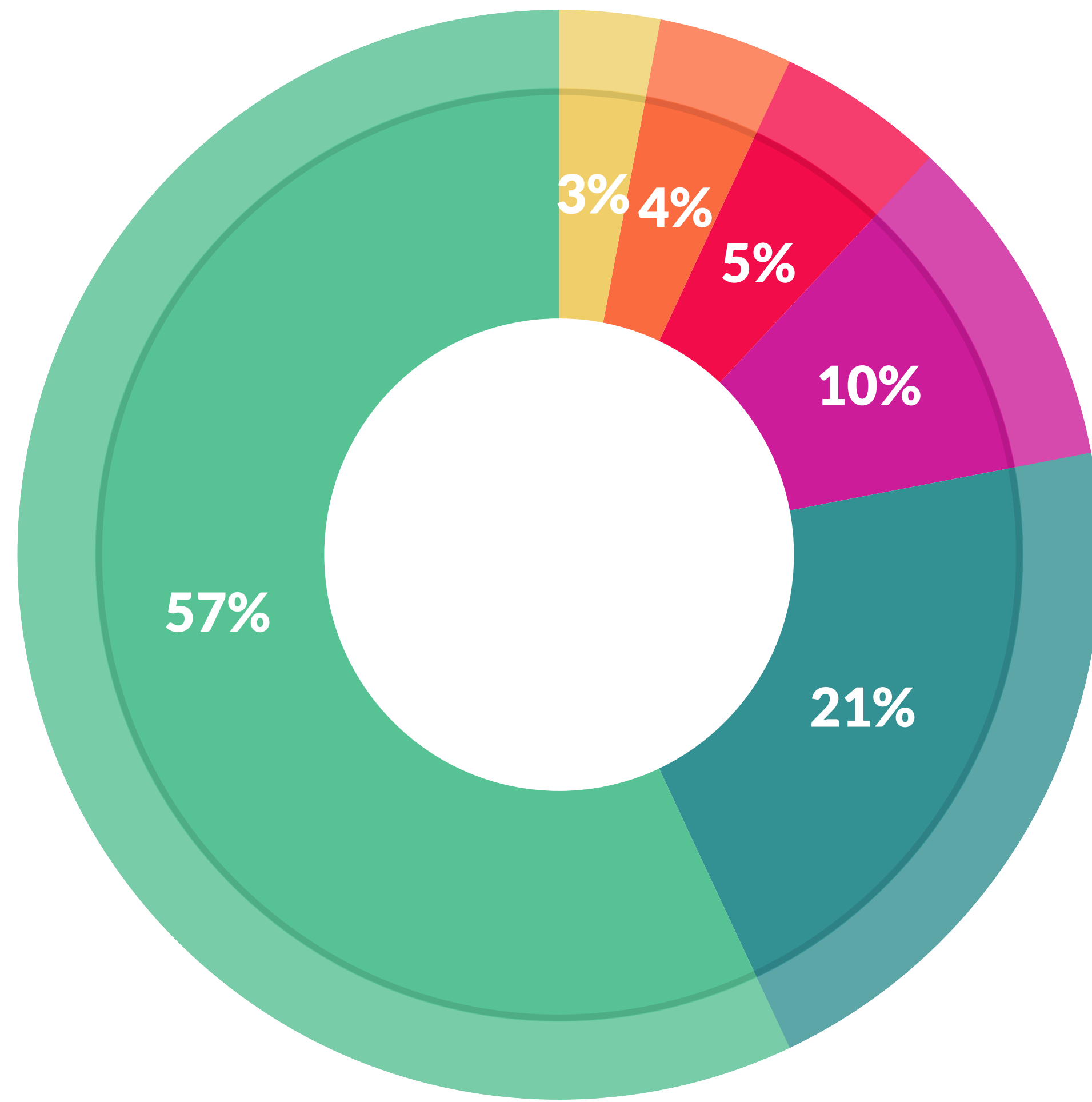- variance 9252.041

# This takes a lot of time!



**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



**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

[J. Canny et al.]

# 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

[J. Canny et al.]

# 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

[J. Canny et al.]
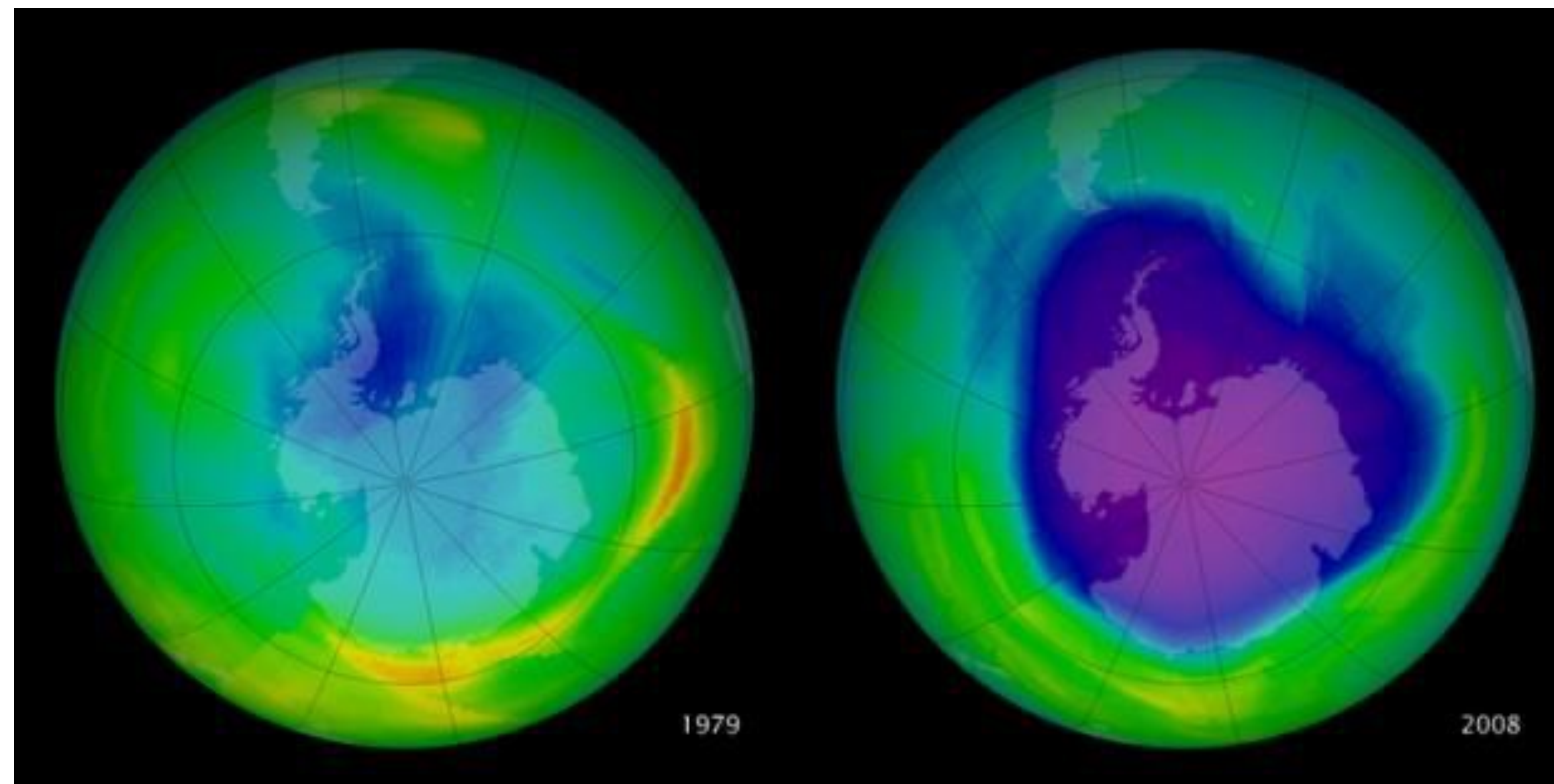
# Dirty Data: Data Scientist's View

- 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…

Northern Illinois University

# Be careful how you detect dirty data

- The appearance of a hole in the earth's ozone layer over Antarctica, first detected in 1976, was so unexpected that scientists didn't pay attention to what their instruments were telling them; they thought their instruments were malfunctioning.

  – National Center for Atmospheric Research

# 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`

- Truncated data: **`Janice Keihanaikukauakahihuliheekahaunaele`** becomes **`Janice Keihanaikukauakahihuliheek`** on Hawaii license

- Redundant records: may be exactly the same or have some overlap

- Formatting issues: `2017-11-07` vs. `07/11/2017` vs. `11/07/2017`

# Data Wrangling

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

- 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 operations and demonstrates them on demand

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

# 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

| Example Column Value (Example erroneous values) | # Structures Enumerated | Final Structure Chosen (Punc = Punctuation) |
|---|---|---|
| -60 | 5 | *Integer* |
| UNITED, DELTA, AMERICAN etc. | 5 | *IspellWord* |
| SFO, LAX etc. (JFK to OAK) | 12 | *AllCapsWord* |
| 1998/01/12 | 9 | *Int Punc(/) Int Punc(/) Int* |
| M, Tu, Thu etc. | 5 | *Capitalized Word* |
| 06:22 | 5 | *Int(len 2) Punc(:) Int(len 2)* |
| 12.8.15.147 (ferret03.webtop.com) | 9 | *Double Punc('.') Double* |
| "GET\b (\b) | 5 | *Punc(") IspellWord Punc(\)* |
| /postmodern/lecs/xia/sld013.htm | 4 | *ξ\** |
| HTTP | 3 | *AllCapsWord(HTTP)* |
| /1.0 | 6 | *Punc(/) Double(1.0)* |

[V. Raman and J. Hellerstein, 2001]

# Potter's Wheel: Transforms

| Transform | Definition | | |
|---|---|---|---|
| Format | $\phi(R, i, f)$ | $=$ | $\{(a_1, \ldots, a_{i-1}, a_{i+1}, \ldots, a_n, f(a_i)) \mid (a_1, \ldots, a_n) \in R\}$ |
| Add | $\alpha(R, x)$ | $=$ | $\{(a_1, \ldots, a_n, x) \mid (a_1, \ldots, a_n) \in R\}$ |
| Drop | $\pi(R, i)$ | $=$ | $\{(a_1, \ldots, a_{i-1}, a_{i+1}, \ldots, a_n) \mid (a_1, \ldots, a_n) \in R\}$ |
| Copy | $\kappa((a_1, \ldots, a_n), i)$ | $=$ | $\{(a_1, \ldots, a_n, a_i) \mid (a_1, \ldots, a_n) \in R\}$ |
| Merge | $\mu((a_1, \ldots, a_n), i, j, \text{glue})$ | $=$ | $\{(a_1, \ldots, a_{i-1}, a_{i+1}, \ldots, a_{j-1}, a_{j+1}, \ldots, a_n, a_i \oplus \text{glue} \oplus a_j) \mid (a_1, \ldots, a_n) \in R\}$ |
| Split | $\omega((a_1, \ldots, a_n), i, \text{splitter})$ | $=$ | $\{(a_1, \ldots, a_{i-1}, a_{i+1}, \ldots, a_n, \text{left}(a_i, \text{splitter}), \text{right}(a_i, \text{splitter})) \mid (a_1, \ldots, a_n) \in R\}$ |
| Divide | $\delta((a_1, \ldots, a_n), i, \text{pred})$ | $=$ | $\{(a_1, \ldots, a_{i-1}, a_{i+1}, \ldots, a_n, a_i, \text{null}) \mid (a_1, \ldots, a_n) \in R \wedge \text{pred}(a_i)\} \ \cup$ |
| | | | $\{(a_1, \ldots, a_{i-1}, a_{i+1}, \ldots, a_n, \ \text{null}, a_i) \mid (a_1, \ldots, a_n) \in R \wedge \neg\text{pred}(a_i)\}$ |
| Fold | $\lambda(R, i_1, i_2, \ldots i_k)$ | $=$ | $\{(a_1, \ldots, a_{i_1-1}, a_{i_1+1}, \ldots, a_{i_2-1}, a_{i_2+1}, \ldots, a_{i_k-1}, a_{i_k+1}, \ldots, a_n, a_{i_l}) \mid$ |
| | | | $(a_1, \ldots, a_n) \in R \wedge 1 \leq l \leq k\}$ |
| Select | $\sigma(R, \text{pred})$ | $=$ | $\{(a_1, \ldots, a_n) \mid (a_1, \ldots, a_n) \in R \wedge \text{pred}((a_1, \ldots, a_n))\}$ |

**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, $\text{left}(x, \text{splitter})$ is the left part of $x$ after splitting by splitter. pred is a function returning a boolean.

[V. Raman and J. Hellerstein, 2001]

# Potter's Wheel: Example

| | | Stewart,Bob |
|---|---|---|
| Anna | Davis | |
| | | Dole,Jerry |
| Joan | Marsh | |

Format
'(.*), (.*)' to '\2 \1'

→

| | | Bob Stewart |
|---|---|---|
| Anna | Davis | |
| | | Jerry Dole |
| Joan | Marsh | |

Split at ' '

| | | Bob | Stewart |
|---|---|---|---|
| Anna | Davis | | |
| | | Jerry | Dole |
| Joan | Marsh | | |

2 Merges ←

| Bob | Stewart |
|---|---|
| Anna | Davis |
| Jerry | Dole |
| Joan | Marsh |

[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<br>Blair, John \|, $73,238<br>Tony Smith \|, $1,00,533 | $(< \xi^* > < ','\ Money >)$ | Parsing is doable despite no good delimiter. A *regular expression* domain can infer a structure of $[0-9,]*$ for last component. |
| MAA \|to\| SIN<br>JFK \|to\| SFO<br>LAX \|–\| ORD<br>SEA \|/\| OAK | $(<len\ 3\ identifier> < \xi^* >$<br>$< len\ 3\ identifier> )$ | Parsing is possible despite multiple delimiters. |
| 321 Blake #7 \|, Berkeley \|, CA 94720<br>719 MLK Road \|, Fremont \|, CA 95743 | $(<number\ \xi^* > < ','\ word>$<br>$<','\ (2\ letter\ word)\ (5\ letter\ integer)>)$ | Parsing is easy because of consistent delimiter. |

[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→state)

- Reshape: Fold and Unfold (aka pivot)

- Positional: Fill and lag

- Sorting, aggregation, key generation, schema transforms

# Interface

- Automated Transformation Suggestions

- Editable Natural Langua



- Visual Transformation P

- Transformation History

[S. Kandel et al., 2011]

# 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

**(a)** Reported crime in `Alabama`

**(b)**
| | | |
|---|---|---|
| *before*: | {'in', ' '} | 'Alabama' → {'Alabama', *word*} |
| *selection*: | {'Alabama'} | 'in' → {'in', *word*, *lowercase*} |
| *after*: | ∅ | ' ' → {' '} |

**(c)**
| | |
|---|---|
| *before*: | {(' '), ('in', ' '), (*word*, ' '), (*lowercase*, ' ')} |
| *selection*: | {('Alabama'), (*word*)} |
| *after*: | ∅ |

**(d)**
| | |
|---|---|
| {(),('Alabama'),()} | ~~{(),(*word*),()}~~ |
| ~~{(' '),(),()}~~ | ~~{(*word*,' '),(),()}~~ |
| {(' '),('Alabama'),()} | {(*word*, ' '),('Alabama'),()} |
| ~~{(' '),(*word*),()}~~ | ~~{(*word*,' '),(*word*),()}~~ |
| {('in', ' '),(),()} | {(*lowercase*, ' '),(),()} |
| {('in', ' '),('Alabama'),()} | {(*lowercase*, ' '),('Alabama'),()} |
| {('in', ' '),(*word*),()} | ~~{(*lowercase*,' '),(*word*),()}~~ |

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

[S. Kandel et al., 2011]

# Data Wrangler Demo

- http://vis.stanford.edu/wrangler/app/



| 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 |

Delete **row 7**

Delete **empty rows**

Fill **row 7** by **copying** values from **above**

| | Year | Property_crime_rate |
|---|---|---|
| 0 | Reported crime in Alabama | |
| 1 | | |
| 2 | 2004 | 4029.3 |
| 3 | 2005 | 3900 |
| 4 | 2006 | 3937 |
| 5 | 2007 | 3974.9 |
| 6 | 2008 | 4081.9 |
| 7 | | |
| 8 | Reported crime in Alaska | |
| 9 | | |
| 10 | 2004 | 3370.9 |
| 11 | 2005 | 3615 |
| 12 | 2006 | 3582 |
| 13 | 2007 | 3373.9 |
| 14 | 2008 | 2928.3 |
| 15 | | |
| 16 | Reported crime in Arizona | |

# 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

- Found significant effect of tool and users found previews and suggestions helpful

- Complaint: No manual fallback, make implications of user choices more obvious for users

# Task Completion Times



User Study Task Completion Time (minutes)

● Wrangler  ● Excel

# Improvements in Prediction

Figure 12  Partially underlined qualified retrieval

| TYPE | ITEM | COLOR | SIZE |
|---|---|---|---|
|  | P. IKE | GREEN |  |

*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 q

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. This user displays both the TYPE table and the SALES table

[Heer et al. 2015]

# 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

[J. M. Hellerstein et al., 2018]

# 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

[J. M. Hellerstein et al., 2018]

# Trifacta Wrangler