Advanced Data Management (CSCI 680/490)

Data Transformation

Dr. David Koop





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





2

Reading & Writing Data in Pandas

Format	Data Description
text	<u>CSV</u>
text	Fixed-Width Text File
text	JSON
text	<u>HTML</u>
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 680/490, Spring 2022

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]









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 $(', ', '', '', ' \setminus t', ' \setminus s+')$
 - 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





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











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









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, supports compression
 - Interfaces in C, Java, MATLAB, etc.
 - Use pd.HDFStore to access, shortcuts: read hdf/to hdf,
- Excel: need to specify sheet when a spreadsheet has multiple sheets
 - pd.ExcelFile Or pd.read excel
- Parquet: big data format, can use compression







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 680/490, Spring 2022

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]









Filtering and Cleaning Data

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





Replacing Values

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

In [<mark>61]: d</mark> ata		In [(52]
Out [6:	1]:	Out[52]
0	1.0	Ο	
1	-999.0	1	
2	2.0	2	
3	-999.0	3	
4 - :	1000.0	4 -	- 10
5	3.0	5	
dtype	: 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





Classifying Data Quality Problems



D. Koop, CSCI 680/490, Spring 2022





Northern Illinois University 12

HoloClean

Input

	Dataset to be cleaned					
	DBAName	Address	City	State	Zip	
t1	John Veliotis Sr.	3465 S Morgan ST	Chicago	IL	60608	
t2	John Veliotis Sr.	3465 S Morgan ST	Chicago	IL	60609	
t3	John Veliotis Sr.	3465 S Morgan ST	Chicago	IL	60609	
t4	Johnnyo's	3465 S Morgan ST	Cicago	IL	60608	

Denial Constraints

- c1: DBAName \rightarrow Zip
- c2: Zip \rightarrow City, State
- c3: City, State, Address \rightarrow Zip

Matching Dependencies

m1: $Zip = Ext_Zip \rightarrow City = Ext_City$ m2: $Zip = Ext_Zip \rightarrow State = Ext_State$ m3: City = Ext_City \land State = Ext_State \land \land Address = Ext_Address \rightarrow Zip = Ext_Zip

External Information

Ext_Address	Ext_City	Ext_State	Ext_Zip
3465 S Morgan ST	Chicago	IL	60608
1208 N Wells ST	Chicago	IL	60 <mark>61</mark> 0
259 E Erie ST	Chicago	IL	60611
2806 W Cermak Rd	Chicago	IL	60623

D. Koop, CSCI 680/490, Spring 2022



Output

	Proposed Cleaned Dataset						
	DBAName	Address	City	State	Zip		
t1	John Veliotis Sr.	3465 S Morgan ST	Chicago	IL	60608		
t2	John Veliotis Sr.	3465 S Morgan ST	Chicago	IL	60608		
t3	John Veliotis Sr.	3465 S Morgan ST	Chicago	IL	60608		
t4	John Veliotis Sr.	3465 S Morgan ST	Chicago	IL	60608		

Marginal Distribution of Cell Assignments

Cell	Possible Values	Probability
10 71-	60608	0.84
tz.zip	60609	0.16
14 01	Chicago	0.95
t4.City	Cicago	0.05
	John Veliotis Sr.	0.99
t4.DBAName	Johnnyo's	0.01







13

Mask Policy

- Masks not required in this class
- Respect all
- Office hours or other interactions: you may ask me to wear a mask





<u>Assignment 3</u>

- Same dataset as A1 and A2...
- ...but dealing with the full raw data now
- Want to clean and transform using Trifacta and pandas
 - Medium
 - Date cleanup
 - Tags expansion
 - [CSCI 680] Artist Data





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',
```







Regular Expression Methods

Argument	Description
findall	Return all non-overlapping matching patte
finditer	Like findall, but returns an iterator
match	Match pattern at start of string and optional matches, returns a match object, and other
search	Scan string for match to pattern; returning the string as opposed to only at the beginn
split	Break string into pieces at each occurrence
sub, subn	Replace all (sub) or first n occurrences (sult $1, 2, \ldots$ to refer to match group e

D. Koop, CSCI 680/490, Spring 2022

erns in a string as a list

- ally segment pattern components into groups; if the pattern rwise None
- a match object if so; unlike match, the match can be anywhere in ning
- of pattern
- ubn) of pattern in string with replacement expression; use symbols elements in the replacement string

[W. McKinney, Python for Data Analysis]







Pandas String Methods with Regexs

In [172]: pattern Out [172]: '([A-Z0-9. %+-]+)@([A-Z0-9.-]+)\\.([A-Z]{2,4})' In [173]: data.str.findall(pattern, flags=re.IGNORECASE) Out[173]: [(dave, google, com)] Dave [(rob, gmail, com)] Rob [(steve, gmail, com)] Steve Wes NaN dtype: object In [175]: matches Out[175]: Dave True Rob True Steve True Wes NaN dtype: object

D. Koop, CSCI 680/490, Spring 2022

- In [174]: matches = data.str.match(pattern, flags=re.IGNORECASE)

[W. McKinney, Python for Data Analysis]









Foofah: Transforming Data By Example

Z. Jin, M. R. Anderson, M. Cafarella, and H. V. Jagadish













• What is the paper's contribution?









- What is the paper's contribution?
- What questions do you have about what is going on?









- What is the paper's contribution?
- What questions do you have about what is going on?
- What does the technique do well/have issues with?









- What is the paper's contribution?
- What questions do you have about what is going on?
- What does the technique do well/have issues with?
- How does its approach compare with Trifacta?









Starting Point: Raw Data

less split	# spl	it1	#	split2	#	split3	#	split4
	2004		2004		2004		2003	
STATE	Participation	Rate 2004	Mean SAT I	Verbal	Mean SAT	I Math	Partici	pation Rate
New York	87		497		510		82	
Connecticut	85		515		515		84	
Massachusetts	85		518		523		82	
New Jersey	<mark>83</mark>		501		514		85	
New Hampshire	<mark>80</mark>		522		521		75	
D.C.	77		489		476		77	
Maine	76		505		501		70	
Pennsylvania	74		501		502		73	
Delaware	73		500		499		73	
Georgia	73		494		493		66	
split	# fo	ld	Abc	fold1	#	value		
New York	2004		Participat	ion Rate 2004	87			
New York	2004		Mean SAT I	Verbal	497			
New York	2004		Mean SAT I	Math	510			
New York	2003		Participat	ion Rate 2003	82			
New York	2003		Mean SAT I	Verbal	496			
New York	2003		Mean SAT I	Math	510			
Connecticut	2004		Participat	ion Rate 2004	85			
Connecticut	2004		Mean SAT I	Verbal	515			
Connecticut	2004		Mean SAT I	Math	515			
Connecticut	2003		Participat	ion Rate 2003	84			
Connecticut	2003		Mean SAT I	Verbal	512			
C	2002		MARY CAT T	Math				









(foal

- Focus on data transformation
- Data transformation tools suffer usability issues:
 - High Skill: familiarity with operations and the effect or their order
 - High Effort: user effort increases as the program becomes longer
- Repetitive and tedious
- Goal: minimize a user's effort and reduce the required background knowledge for data transformation tasks









Getting Lost in Transformations

Bureau of I.A.	
Regional Director	Numbers
Niles C.	Tel: (800)645-8397
	Fax: (907)586-7252
Jean H.	Tel: (918)781-4600
	Fax: (918)781-4604
Frank K.	Tel: $(615)564-6500$
	Fax: (615)564-6701



Intermediate Table

D. Koop, CSCI 680/490, Spring 2022

	Tel	Fax
Niles C.	(800)645-8397	
		(615)564-6701
Jean H.	(918)781-4600	
Frank K.	(615)564-6500	

Problem Table

		Tel	Fax
•	Niles C.	(800)645-8397	(907)586-7252
	Jean H.	(918)781-4600	(918)781-4604
l+	Frank K.	(615)564-6500	(615)564-6701

Desired Solution







Foofah Design: Programming by Example









Input, Output, and Transformations

Fax:(918)781-4604



Raw Data:

- A grid of values, i.e., spreadsheets "Somewhat" structured - must have some regular structure or is automatically generated.



User Input:

 Sample from raw data Transformed view of the sample

Tel:(800)645-839



Program to synthesize: A loop-free Potter's Wheel [2] program

Transformations Targeted: 1. Layout transformation



D. Koop, CSCI 680/490, Spring 2022















26

Transformations

Operator	Descriptio
Drop	Deletes a o
Move	Relocates
	other in the
Сору	Duplicates
	column to
Merge	Concatena
	merged co
Split	Separates
Eald	at the occ
FOID	Collapses
Unfold	"Inflatton"
Uniola	data value
Fill	Fill empty
Divide	Divide is ı
	columns b
Delete	Delete rov
	predicate
Extract	Extract fir
	sion each o
Transpose	Transpose
$Wrap\ (\mathrm{added})$	Concatena

on
column in the table
a column from one position to an-
he table
es a column and append the copied
the end of the table
ates two columns and append the
olumn to the end of the table
a column into two or more halves
currences of the delimiter
all columns after a specific column
column in the output table
n" tables and move information from
es to column names
y cells with the value from above
used to divide one column into two
based on some predicate
ws or columns that match a given
rst match of a given regular expres-
cell of a designated column
e the rows and columns of the table
ate multiple rows conditionally









Proposed Solution

- No loops
- Assumes relational tables
- ... and perfect data?

D. Koop, CSCI 680/490, Spring 2022

• Use a small, manually transformed portion of the data to infer a program (in Potter's Wheel syntax) based on the specified data transformation operations









Foofah Solution















Need a Heuristic Function to Prune



D. Koop, CSCI 680/490, Spring 2022

Most transformations are composed of cell-based operations

Mike Anders	on l	Jniversity	y of Michigan	PhD Student
				6
Mike	e Ander	rson l	Jniversity of M	lichigan

Remove a cell

Transform a cell



NIU







Use Add/Remove/Modify + Move

Table Edit Distance (TED) Definition: The cost of transforming Table T_1 to Table T_2 using the cell-level operators Add/Remove/Move/Transform cell.

$$\mathrm{TED}(T_1,T_2) = ($$

 $\min_{(p_1,...,p_k)\in P(T_1,T_2)}\sum_{i=0}^{\cdot} cost(p_i)$ • $P(T_1, T_2)$: Set of all "paths" transforming T_1 to T_2 using cell-level operators













Table Edit Distance Batch



8 Transform operations

D. Koop, CSCI 680/490, Spring 2022

2 "batched" Transform operations [Z. Jin et al., 2017]









Geometric Patterns Used to Batch

Pattern	Formulation $(X \text{ is a table edit})$
Horizontal to Horizontal Horizontal to Vertical Vertical to Horizontal Vertical to Vertical One to Horizontal One to Vertical Remove Horizontal Remove Vertical	$ \{X((x_i, y_i), (x_j, y_j)), X((x_i, y_i) \\ \{X((x_i, y_i), (x_j, y_j)), X((x_i, y_i) \\ \{X((x_i, y_i), (x_j, y_j)), X((x_i + 1) \\ \{X((x_i, y_i), (x_j, y_j)), X((x_i, y_i) \\ \{X((x_i, y_i), (x_j, y_j)), X((x_i, y_i) \\ \{X((x_i, y_i)), (x_j, y_j)), X((x_i, y_i) \\ \{X((x_i, y_i)), X((x_i, y_i + 1)), \dots \\ \{X((x_i, y_i)), X((x_i + 1, y_i)), \dots \\ \{X((x_i, y_i)), X((x_i + 1, y_$

operator)	Related Operators
$ + 1), (x_{j}, y_{j} + 1)), \dots \} $ $ + 1), (x_{j} + 1, y_{j})), \dots \} $ $ -, y_{i}), (x_{j}, y_{j} + 1)), \dots \} $ $ -, (x_{j}, y_{j} + 1)), \dots \} $ $ +, (x_{j} + 1, y_{j})), \dots \} $ $ +, \{x_{j} + 1, y_{j})), \dots \} $	Delete(Possibly) Fold, Transpose Unfold,Transpose Move, Copy, Merge, Split, Extract, Drop Fold(Possibly), Fill(Possibly) Fold, Fill Delete Drop, Unfold













Other Pruning Rules

- Global:
 - Missing Alphanumerics: check that character maintained
 - No effect: meaningless operation
 - Introducing Novel Symbols: check that no new characters added
- Property-specific:
 - Generating Empty Columns
 - Null in Column

















Search Strategies and Pruning Rules



D. Koop, CSCI 680/490, Spring 2022



[Z. Jin et al., 2017]







User Study Results

WRANGLER

Test	Complex	$\geq 4 \text{ Ops}$	Time	Mouse	Key	Time vs Wrangler	Mouse	Key
PW1	No	No	104.2	17.8	11.6	49.4 > 52.6%	20.8	22.6
PW3 (modified)	No	No	96.4	28.8	26.6	$38.6 \ 60.0\%$	14.2	23.6
ProgFromEx13	Yes	No	263.6	59.0	16.2	$145.8 \searrow 44.7\%$	43.6	78.4
PW5	Yes	No	242.0	52.0	15.2	58.8 > 75.7%	31.4	32.4
ProgFromEx17	No	Yes	72.4	18.8	11.6	48.6 > 32.9%	18.2	15.2
PW7	No	Yes	141.0	41.8	12.2	44.4 > 68.5%	19.6	35.8
Proactive1	Yes	Yes	324.2	60.0	13.8	$104.2 \ \mathbf{67.9\%}$	41.4	57.0
Wrangler3	Yes	Yes	590.6	133.2	29.6	137.0 > 76.8%	58.6	99.8

D. Koop, CSCI 680/490, Spring 2022

Foofah









Comparisons with other tools

Success rates on pure layout transformation benchmark tasks



Foofah FlashRelate ProgFromEx Wrangler

Sizes of input-output examples required for benchmark tests

D. Koop, CSCI 680/490, Spring 2022

Success rates on benchmark tasks requiring syntactic transformations



Foofah FlashRelate ProgFromEx Wrangler















C	
Customer Name	
John K. Doe Jr.	Doe
Mr. Doe, John	Doe
Jane A. Smith	Smi
MS. Jane Smith	Smi
Smith, Jane	Smi
Dr Anthony R Von Fange III	Vor
Peter Tyson	Tys
Dan E. Williams	Wil
James Davis Sr.	Dav
James J. Davis	Dav
Mr. Donald Edward Miller	Mil













C	D
Address	Output
4297 148th Avenue NE L105, Bellevue, WA 98007	Bellevue, WA, 98007
2720 N Mesa St, El Paso, 79902, USA	El Paso, TX, 79902
3524 W Shore Rd APT 1002, Warwick,02886	Warwick, RI, 02886
4740 N 132nd St, Omaha, 68164	Omaha, NE, 68164
10508 Prairie Ln, Oklahoma City	Oklahoma City, OK, 73162
525 1st St, Marysville, WA 95901	Marysville, CA, 95901
211 W Ridge Dr, Waukon,52172	Waukon, IA, 52172
1008 Whitlock Ave NW, Marietta, 30064	Marietta, GA, 30064
602 Highland Ave, Shinnston, 26431	Shinnston, WV, 26431
840 W Star St, Greenville, 27834	Greenville, NC, 27834

=		,				
				Show	/ Instru	ction
	Get	Transfor	matio	ns		
	Sear	ch:				
	Le .			>	-	ş
Bui	itins Address	Parser Parse	=WitbBir	ngMa	ns	
(Sy	stem.String)			.g	22	
				>	-	4
Bu	itins.Address	Parser				
Pa	rseWithBingM	lapsAndPyt	honUsA	ddre	ssLibra	ry
(Sy	stem.String)					
				>	-	ş
Hu	manizer.ToTit	tleCase Trai	nsform(Syster	m.Strin	ig)











С	D
Transaction Date	output
Wed, 12 Jan 2011	2011-01-12-Wednesday
Thu, 15 Sep 2011	2011-09-15-Thursday
Mon, 17 Sep 2012	
2010-Nov-30 11:10:41	
2011-Jan-11 02:27:21	
2011-Jan-12	
2010-Dec-24	
9/22/2011	
7/11/2012	
2/12/2012	

С	D	Transform Data by Example
Transaction Date	output	=
Wed, 12 Jan 2011	2011-01-12-Wednesday	Show Instructions Get Transformations
Thu, 15 Sep 2011	2011-09-15-Thursday	> + <i>f</i>
Mon, 17 Sep 2012	2012-09-17-Monday	System.DateTime Parse(System.String)
2010-Nov-30 11:10:41	2010-11-30-Tuesday	System.Convert ToDateTime(System.String)
2011-Jan-11 02:27:21	2011-01-11-Tuesday	> ¥ \$
2011-Jan-12	2011-01-12-Wednesday	DateFormat.Program Parse(System.String)
2010-Dec-24	2010-12-24-Friday	
9/22/2011	2011-09-22-Thursday	
7/11/2012	2012-07-11-Wednesday	
2/12/2012	2012-02-12-Sunday	© Microsoft Privacy Terms Feedback





TDE: Synthesized Function

Input Examples		Return Object Dump					Member method result dump					Desired Outpu	
		Year	Month	Day	Day-of-week	Day-of-Year	 ToLongDateString()	ToTimeStr()	ToUTC()	ToBinary()			
Wed, 12 Jan 2011	f	2011	01	12	Wednesday	12	 Wednesday, January 12, 2011	12:00:00 AM			Synthesize	2011-01-12 (We	
Thu, 15 Sep 2011		2011	09	15	Thursday	258	 Thursday, September 15, 2011	12:00:00 AM			Synthesize	2011-09-15 (Th	
Mon, 17 Sep 2012		2012	09	17	Monday	261	 Monday, September 17, 2012	12:00:00 AM					
2010-Nov-30 21:10:41		2010	11	30	Tuesday	334	 Tuesday, November 30, 2010	09:10:41 PM					
2011-Jan-11 02:27:21			2011	01	11	Tuesday	11	 Tuesday, January 11, 2011	02:27:21 AM				
2011-Jan-12		2011	01	12	Wednesday	12	 Wednesday, January 12, 2012	12:00:00 AM					











- Row-to-row translation only
- Search System, GitHub, and StackOverflow for functions
- Given dataset with examples
 - Use L1 from library
 - Compose synthesized programs (L2)
 - Rank best transformations







TDE Benchmarks

System	Total cases (239)	FF-GR-Trifacta (46)	Head cases (44)	StackOverflow (49)	BingQL-Unit (50)	BingQL-Other (5
TDE	72%~(173)	91% (42)	82%~(36)	63%~(31)	96%~(48)	32%~(16)
TDE-NF	53% (128)	87% (40)	41% (18)	35%~(17)	96% (48)	10% (5)
FlashFill	23% (56)	57% (26)	34%~(15)	31%~(15)	0% (0)	0% (0)
Foofah	3% (7)	9% (4)	2% (1)	4% (2)	0% (0)	0% (0)
DataXFormer-UB	38% (90)	7% (3)	36%~(16)	35%~(17)	62% (31)	46%~(23)
System-A	13% (30)	52% (24)	2% (1)	10%~(5)	0% (0)	0% (0)
OpenRefine-Menu ⁸	4% (9)	13% (6)	2% (1)	4% (2)	0% (0)	0% (0)

- TDE and FlashFill focused on row-to-row transformations
- Foofah considers a wider range of transformations (table reformatting)













TDE Benchmarks

System	Total cases (239)	FF-GR-Trifacta (46)	Head cases (44)	StackOverflow (49)	BingQL-Unit (50)	BingQL-Other (5
TDE	$\boxed{\begin{array}{c c} 72\% \ (173) \end{array}}$	91% (42)	82%~(36)	63%~(31)	96% (48)	32%~(16)
TDE-NF	53% (128)	87% (40)	41% (18)	35%~(17)	96% (48)	10% (5)
FlashFill	23% (56)	57%~(26)	34%~(15)	31%~(15)	0% (0)	0% (0)
Foofah	3% (7)	9% (4)	2% (1)	4% (2)	0% (0)	0% (0)
DataXFormer-UB	38% (90)	7% (3)	36% (16)	35%~(17)	62% (31)	46%~(23)
System-A	13% (30)	52% (24)	2% (1)	10%~(5)	0% (0)	0% (0)
OpenRefine-Menu ⁸	4% (9)	13% (6)	2% (1)	4% (2)	0% (0)	0% (0)

- TDE and FlashFill focused on row-to-row transformations
- Foofah considers a wider range of transformations (table reformatting)













Trifacta's Transform by Example



