

ExtractTable

Extracting Tables from Plain Text

Leonardo Hübscher

Prof. Felix Naumann
Lan Jiang

12.02.2021

Data preparation

Tasks



Collecting data



Cleaning data

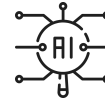


Normalizing data

Benefits



Analysis



Machine learning



Data warehouses

Downside



Tedious tasks



Time consuming



Error prone

Tables are great, right?

- Structured
- Store information in high densities
- Often used for data sharing
- Interpretable by humans and machines
- Tables look similar, yet they are so different...

```
2/28/2008, 1:00:00, NaN  
2/28/2008, 2:00:00, NaN  
2/28/2008, 3:00:00, NaN
```

X	Y	Z
0.00	0.00	0.00

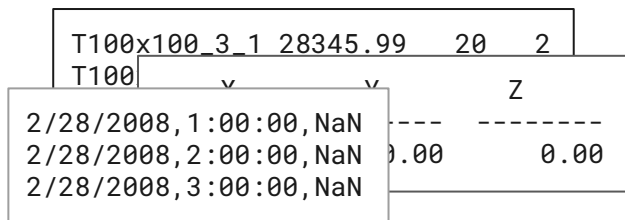


T100x100_3_1	28345.99	20	2
T100x100_3_2	29580.17	23	3
T100x100_3_3	27062.23	20	2

Vision

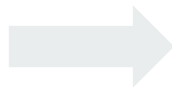
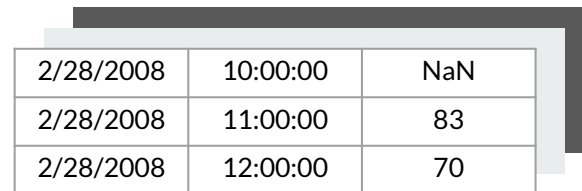
- Data preparation fuses information from different sources
- Valuable information is often stored in tables
- Tables are represented differently
- Manual pre-processing is time consuming

We want to decrease the time data scientists spend with data wrangling by extracting tables from plain text files automatically.



T100x100_3_1	28345.99	20	2
T100	v	v	Z
2/28/2008,1:00:00,NaN	---	---	---
2/28/2008,2:00:00,NaN	0.00	0.00	0.00
2/28/2008,3:00:00,NaN			

All kinds of table formats

2/28/2008	10:00:00	NaN
2/28/2008	11:00:00	83
2/28/2008	12:00:00	70

Standardized table format

Outline

1. Introduction
2. Background
3. Mission
4. Demo
5. Algorithm
6. Evaluation
7. Outlook

Parsing Instructions

CSV

```
"name","age","comment"  
"Max Mustermann",20,"max says \"hello\""  
Peter Pan,30,
```

delimited by character

ASCII

name	age	comment
Max Mustermann	20	Max says "hello"
Peter Pan	30	

delimited by layout

Parsing Instructions

CSV

```
"name","age","comment"  
"Max Mustermann",20,"max says \"hello\""  
Peter Pan,30,
```

delimited by character

Dialect

Delimiter character

Quotation character

Escape character

required

optional

ASCII

name	age	comment
Max Mustermann	20	Max says "hello"
Peter Pan	30	.

delimited by layout

Indexes

19

30

Related Work

CleverCSV²
2019



Gertjan van den Burg
*Postdoctoral researcher at
The Alan Turing Institute*

Dialect detection of CSV files based
on row-patterns and cell data types

Both

- Work on file-level
- Expect single table per file
- Detect dialect of CSV tables

Hypoparsr³
2016



Till Döhmen
*Research Assistant at
Fraunhofer FIT*

Improve quality by deferring
decision-making for sub-problems
like encoding & dialect detection

[2] van den Burg, G. J. J. (2019), 'Wrangling messy CSV files by detecting row and type patterns'

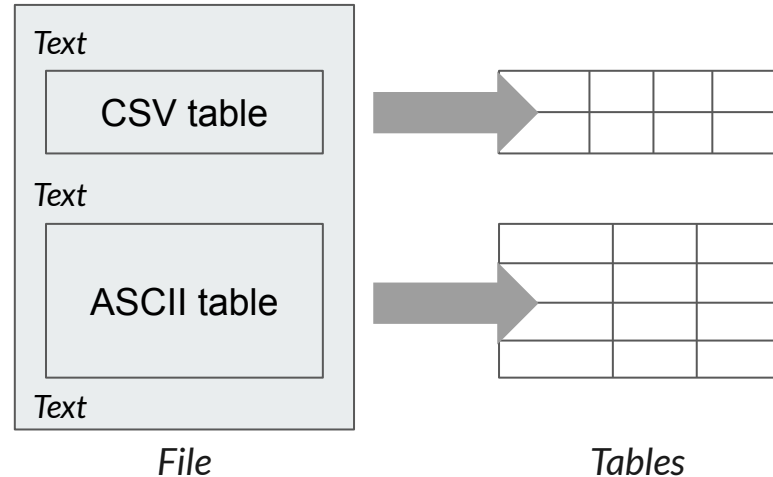
[3] Döhmen, T. (2016), 'Multi-Hypothesis Parsing of Tabular Data in Comma-Separated Values (CSV) Files'

Mission

We want to develop an algorithm that is able to extract tables from more *complex* files compared to existing solutions.

Key features

- ASCII table support
- Support for more dialects (e.g. multi-character delimiters)
- Detection of multiple tables within a file (of different dialects)
- Support files with surrounding text



Demo

Algorithm Workflow



Input Lines

Tasks

1. Detect valid dialects for CSV table candidates (line-based)
2. Detect valid split indexes for ASCII table candidates (table-based)

Output Lines incl. parsing instructions

Example

2/28/2008,11:00:00,83

Detected parsing instructions

delimiter: / quotation: none escape: none
delimiter: , quotation: none escape: none
delimiter: : quotation: none escape: none
(no valid layout parsing instruction)

Algorithm Workflow



Input Lines incl. parsing instructions

Tasks

1. Apply parsing instructions to build solution space

Output Solution space for each line, containing one interpretation per parsing instruction

Example

2/28/2008,11:00:00,83

Solution space

2	28	2008,11:00:00,83
2/28/2008	11:00:00	83
2/28/2008,11	00	00,83

Algorithm Workflow



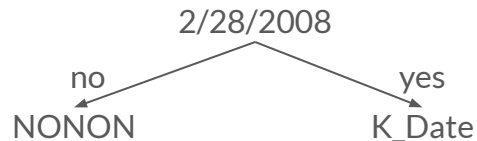
Input Solution space per line

Tasks

1. Describe table cell content by:
 - Recognizing known data types (RegEx-based)
 - Splitting into atomar components (String, *Number*, *Other*)

Output Solution space for each line, including features

Does a data type pattern match?



Number covers all kinds of numbers (int, signed number, floats, scientific)

Other matches everything that isn't a string or number component (similar to non-alphanumeric)

Algorithm Workflow



Input Solution space for each line, including features

Tasks

1. Group solution spaces into bins, use column count and parsing instruction as key
2. Split bins into compatibility blocks using some pattern-based consistency score

Output One table candidate per compatibility block

Algorithm Workflow

Task 1 Group solution spaces into bins, use column count and parsing instruction as key

date,time,humidity
2/28/2008,11:00:00,83
2/28/2008,12:00:00,70

Bin 1 3 columns, CSV, [delimiter=/]			Bin 2 3 columns, CSV, [delimiter=,]			Bin 3 3 columns, CSV, [delimiter=:]		
			date	time	humidity			
2	28	2008,11:00:00,83	2/28/2008	11:00:00	83	2/28/2008,11	00	00,83
2	28	2008,12:00:00,70	2/28/2008	12:00:00	70	2/28/2008,12	00	00,70

Algorithm Workflow

Task 2 Split bins into compatibility blocks using a column consistency score

Bin 1 3 columns, CSV, [delimiter=/]			Bin 2 3 columns, CSV, [delimiter=,]			Bin 3 3 columns, CSV, [delimiter=:]		
			date	time	humidity			
2	28	2008,11:00:00,83	2/28/2008	11:00:00	83	2/28/2008,11	00	00,83
2	28	2008,12:00:00,70	2/28/2008	12:00:00	70	2/28/2008,12	00	00,70

Bin 1 3 columns, CSV, [delimiter=/]			Bin 2 3 columns, CSV, [delimiter=,]			Bin 3 3 columns, CSV, [delimiter=:]		
			K_Text	K_Text	K_Text			
N	N	NONNONN	K_Date	K_Time	N	NONON	NN	NN
N	N	NONNONN	K_Date	K_Time	N	NONON	NN	NN

Algorithm Workflow

Task 2 Split bins into compatibility blocks using a column consistency score

Bin 1 3 columns, CSV, [delimiter=/]			Bin 2 3 columns, CSV, [delimiter=,]			Bin 3 3 columns, CSV, [delimiter=:]		
			date	time	humidity			
2	28	2008,11:00:00,83	2/28/2008	11:00:00	83	2/28/2008,11	00	00,83
2	28	2008,12:00:00,70	2/28/2008	12:00:00	70	2/28/2008,12	00	00,70

Bin 1 3 columns, CSV, [delimiter=/]			Bin 2 3 columns, CSV, [delimiter=,]			Bin 3 3 columns, CSV, [delimiter=:]		
			K_Text	K_Text	K_Text			
N	N	NONNONN	K_Date	K_Time	N	NONON	NN	NN
N	N	NONNONN	K_Date	K_Time	N	NONON	NN	NN

Algorithm Workflow

Bin 1 3 columns, CSV, [delimiter=/]			Bin 2 3 columns, CSV, [delimiter=,]			Bin 3 3 columns, CSV, [delimiter=:]		
			K_Text	K_Text	K_Text			
N	N	NONNONN	K_Date	K_Time	N	NONON	NN	NN
N	N	NONNONN	K_Date	K_Time	N	NONON	NN	NN

N	N	NONNONN
N	N	NONNONN

K_Text	K_Text	K_Text
K_Date	K_Time	N
K_Date	K_Time	N

K_Text	K_Text	K_Text
K_Date	K_Time	N
K_Date	K_Time	N

NONON	NN	NN
NONON	NN	NN

Algorithm Workflow



Input Table candidates

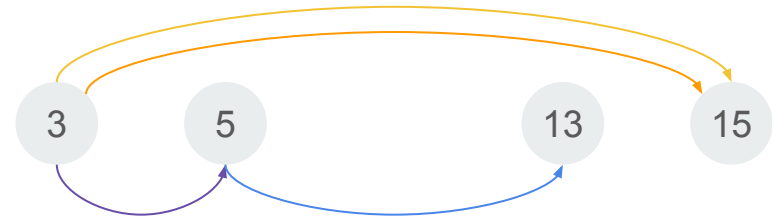
Tasks

- map table selection problem to multi-edge DAG
V = line indexes E = table candidate
distances are based on consistency and size
- use shortest path to find optimal solution

Output Selection of standardized tables

Table candidates

Table 1	start: 3	end: 15	distance: -135.2	✓
Table 2	start: 3	end: 5	distance: -9.0	✗
Table 3	start: 5	end: 13	distance: -56.7	✗
Table 4	start: 3	end: 15	distance: -101.4	✗



Insight: ASCII Table Detection



Task Find the correct split indexes for a given ASCII table

	Student	Supervision
1: Topic	Leonardo Hübscher	Felix Naumann, Lan Jiang
2: Extracting Plain Tables from Text	Jakob Köhler	Thorsten Papenbrock
3: Distributed Duplicate Detection on Streaming-Data	Lasse Kohlmeyer	Ralf Krestel, Tim Repke
4: Multi-Aspect Embeddings for Fiction Novels	Noel Danz	Ralf Krestel, Tim Repke
5: Generating Rap Lyrics with Flow and Rhythm		

Idea Detect split indexes by finding vertical lines of white spaces

Insight: ASCII Table Detection

1: Topic	Student	Supervision
2: Extracting Plain Tables from Text	Leonardo Hübscher	Felix Naumann, Lan Jiang
3: Distributed Duplicate Detection on Streaming-Data	Jakob Köhler	Thorsten Papenbrock
4: Multi-Aspect Embeddings for Fiction Novels	Lasse Kohlmeier	Ralf Krestel, Tim Repke
5: Generating Rap Lyrics with Flow and Rhythm	Noel Danz	Ralf Krestel, Tim Repke

transform to
bitmap
1 = white space

[illegible]

Insight: ASCII Table Detection

1: Topic	Student	Supervision
2: Extracting Plain Tables from Text	Leonardo Hübscher	Felix Naumann, Lan Jiang
3: Distributed Duplicate Detection on Streaming-Data	Jakob Köhler	Thorsten Papenbrock
4: Multi-Aspect Embeddings for Fiction Novels	Lasse Kohlmeier	Ralf Krestel, Tim Repke
5: Generating Rap Lyrics with Flow and Rhythm	Noel Danz	Ralf Krestel, Tim Repke

transform to
bitmap
1 = white space

[illegible]

Insight: ASCII Table Detection

[illegible]

Idea Group consecutive vertical lines into blocks

Block width (number of characters)
 height (number of rows)

Insight: ASCII Table Detection

```

1: 0000011111|11111|11111|1111|1111|11111111111111111111|0000001|11111111|1|0000000000111|1111|1111
2: 0000000000|100000|100000|10000|10000|111111111111111111|00000000|10000000|1|00000100000000|1000|10000
3: 00000000000100000000010000000010010000000000001110000100000011111110000000100000000011111
4: 000000000001000000000100010000000100000011111111000010000000001110000100000001000100001
5: 00000000010001000000100010000100010000001111111110001000011111110000100000001000100001

```

When is a table started? If a block of length equal to T_{\min_rows} (2) appears

Insight: ASCII Table Detection

[illegible]

What happens if a block is continued only partially? Split the block

What happens if a block is discontinued?

- If width = 1 and height < $T_{\text{essential}}$ or block is leading/ trailing
 - a. Remove block from table t
- Otherwise:
 - a. Copy table t and remove block from copy
 - b. Add table t to final table set

Insight: ASCII Table Detection

[illegible]

We reached the final line: all blocks are discontinued

What happens if a block is discontinued?

- If width = 1 and height < $T_{\text{essential}}$ or block is leading/ trailing
 - a. Remove block from table t
- Otherwise:
 - a. Copy table t and remove block from copy
 - b. Add table t to final table set

Insight: ASCII Table Detection

[illegible]

1: Topic	Student	Supervision
2: Extracting Plain Tables from Text	Leonardo Hübscher	Felix Naumann, Lan Jiang
3: Distributed Duplicate Detection on Streaming-Data	Jakob Köhler	Thorsten Papenbrock
4: Multi-Aspect Embeddings for Fiction Novels	Lasse Kohlmeyer	Ralf Krestel, Tim Repke
5: Generating Rap Lyrics with Flow and Rhythm	Noel Danz	Ralf Krestel, Tim Repke

Evaluation

Dataset 1,000 annotated files from Mendeley, UKData, GitHub

Annotations

- Parsing instructions (line-level)
- Table ranges
- Row types

Complexity levels of files

- Simple single: 1 Table, w/o surrounding text
- Complex single: 1 Table, with surrounding text
- Complex multi: >1 Table

Annotate Line 1 (1/1628) APPLY TO REMAINING

Body,Body Name,Date,Transaction Number,Invoice Number,Amount,Supplier Name,Supplier ID,VAT Registration Number

Character Delimiter Quote Escape Header

Delimiter Type Delimiter Quote Escape Header

Preview

Body	Body Name	Date	Transaction Number	Invoice Number	Amount	Supplier Name	Supplier ID	VAT Registration Number	Exi Art

< PREVIOUS ROW NEXT ROW >

Screenshot of line annotation

Evaluation

Line classification
(table/ non-table)

File complexity	Accuracy	Precision	Recall	F1	Balanced accuracy
simple	0.999	1.000	0.999	1.000	1.000
complex-single	0.995	0.996	0.998	0.997	0.979
complex-multi	0.919	0.917	0.997	0.955	0.697

Evaluation

Line classification (table/ non-table)

File complexity	Accuracy	Precision	Recall	F1	Balanced accuracy
Simple single	0.999	1.000	0.999	1.000	1.000
Complex single	0.995	0.996	0.998	0.997	0.979
Complex multi	0.919	0.917	0.997	0.955	0.697

Parsing correctness (line-level)

Table type	Correct (cleaned)*
CSV	88.5 %
ASCII	71.8 %

What's Next?

Next Steps

- Continue with evaluation (comparison to existing solutions)
- (Work on improvements as indicated by evaluation)
- If we have some spare time left: demo website
- Deadline: 22.04 → start writing in two weeks

Future Work

- Add support for spanning rows/ columns

Overview



Key features

- ASCII table support
- Support for more dialects (e.g. multi-character delimiters)
- Detection of multiple tables within a file (of different dialects)
- Support files with surrounding text

