



DANDD INSTANCE GENERATION IN TEXTILE PLOT ENVIRONMENT

Natsuhiko KUMASAKA
Faculty of Science and Technology
Keio University

BACKGROUND

- Variety of data sources
 - Flat files (CSV, TSV, fixed format,...)
 - RDBMS (PostgreSQL, Access, etc...)
 - Spread sheet (Excel,...)
 - Webpage (HTML, CGI,...)
- DandD project
 - DandD instance
 - XML document describing data and its attributes
 - DandD client-server system (Yokouchi and Shibata 2004)
 - Integrate different data sources
 - Employ DandD Instance as an intermediate entity
 - DandDEditor (Yokouchi 2005)
 - Create / Edit DandD instances
 - Restricted in well organised data sources

TEXTILE PLOT ENVIRONMENT

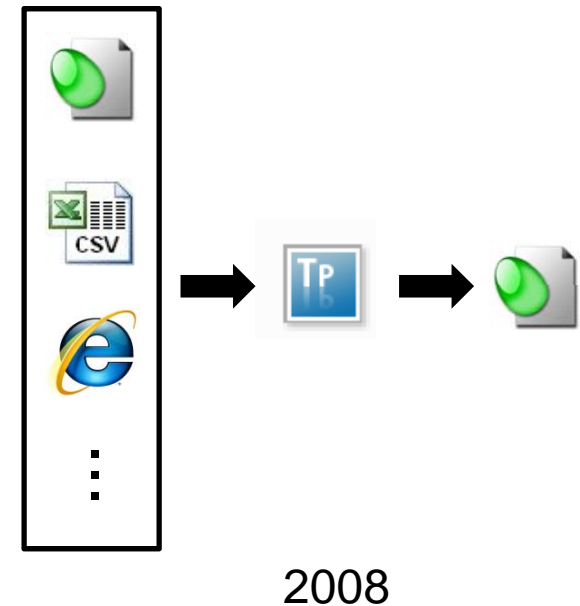
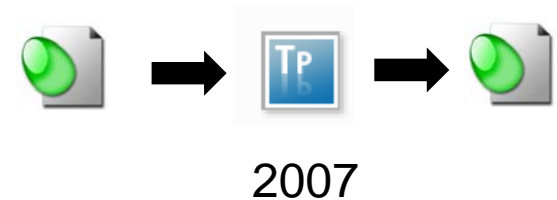
(KUMASAKA AND SHIBATA 2007)

○ Previous features

- Open DandD instance
- Modify textile plots
 - Remove/replace columns, cases and observations
 - Add/Remove/Modify attributes
 - Re-organise data structure
- Save DandD instance
 - Audit description of data analysis

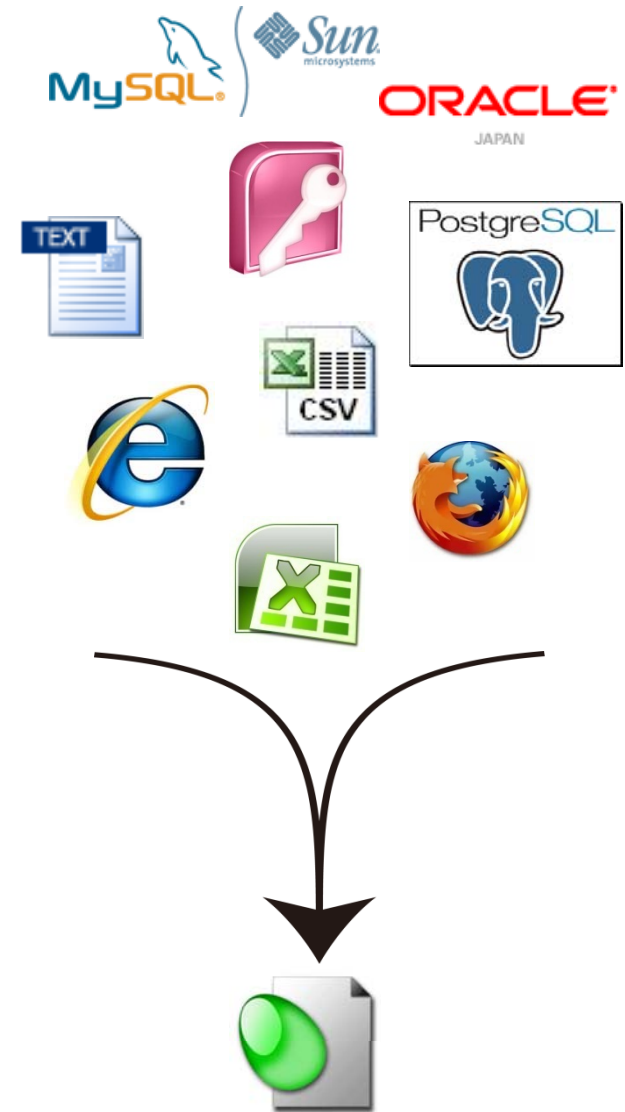
○ New feature in 2008

- Open various data sources
 - Extract tables embedded in data files or databases
 - Separate attributes from tables
 - Divide tables into a set of datavectors
 - Check consistency in datavectors

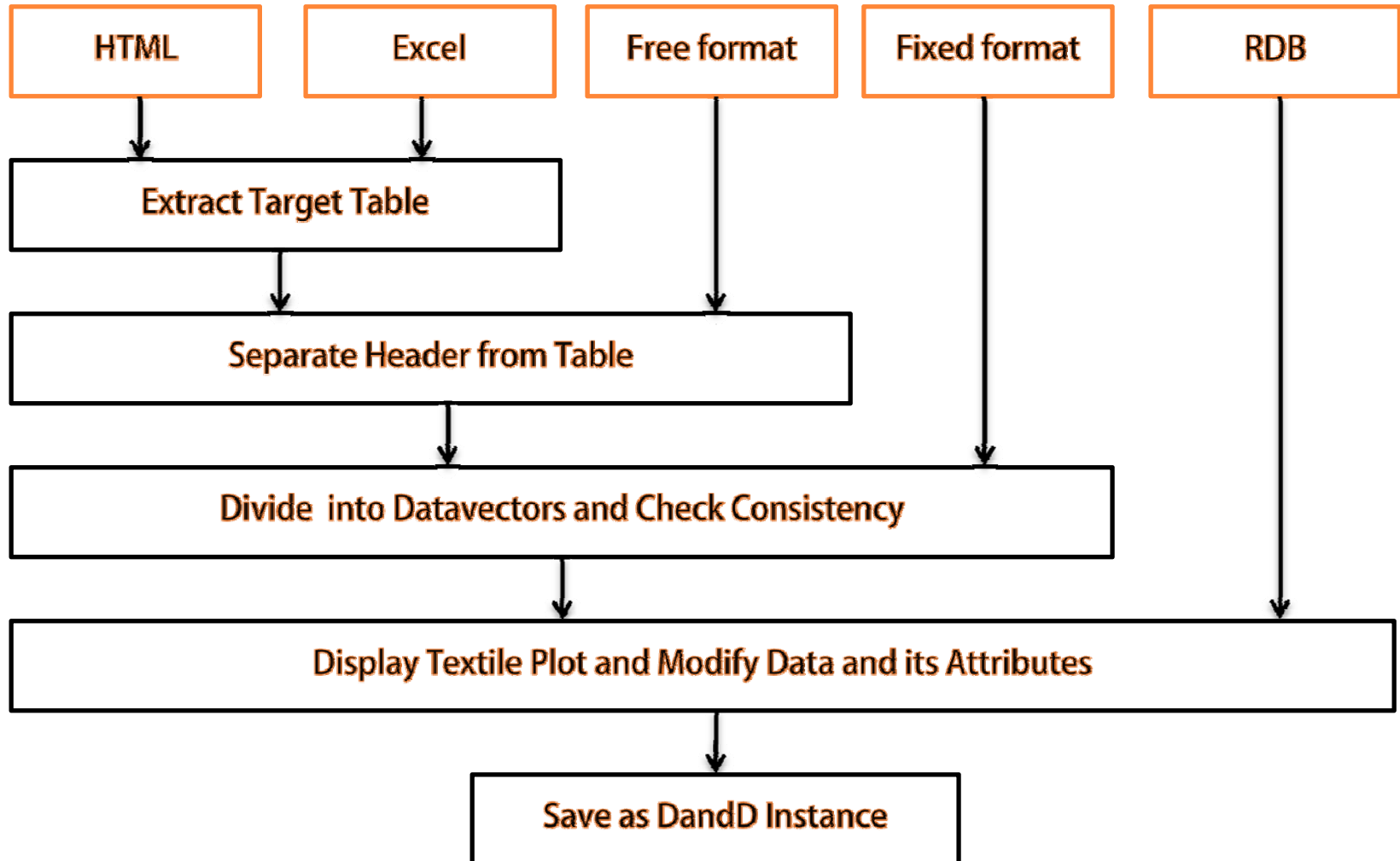


TARGET DATA SOURCES

- Relational Databases
 - RDBMS
 - PostgreSQL, Access, MySQL
 - Oracle (not tested yet)
- Data File
 - Flat file
 - Fixed format
 - Free format (CSV, TSV)
 - Non-flat file
 - Excel
 - HTML



STRATEGY



HTML

- Extract target <TABLE>s
 1. Create DOM Tree
 - DOM Level 1 (W3C 2008)
 - CyberNeko HTML Parser 1.9.6.1 (Clark 2008)
 2. Find <TABLE>s
 3. Count the number of <TR> and <TD>
 4. Select <TABLE>s which have a number of cells
- Separate header from table body
 1. Search <TBODY> and <THEAD>
 2. Find disjoint sets of <TR>s
 3. Measure dissimilarity between disjoint <TR>s
 4. Using clustering algorithm
- Merge <TABLE>s according to header information
- Divide into datavectors and check consistency

EXTRACT TARGET <TABLE>S

<TR>: 1, <TD>: 1
Num of cells: 1

<TR>: 1, <TD>: 2
Num of cells: 2

The screenshot shows the Yahoo! Real Estate website interface. At the top, there's a navigation bar with the Yahoo! logo and various links. Below that, there's a search bar and a main content area. The main content area features a table of search results for apartments. The table has 9 columns and 21 rows of data. The columns represent different attributes of the properties, such as location, price, and features. The rows list individual apartment listings. The table is highlighted with a red border, indicating it is the target for extraction. Below the table, there's a footer with additional information and a copyright notice.

<TR>: 21, <TD>: 9
Num of cells: 189

<TR>: 1, <TD>: 1
Num of cells: 1

SEPARATE HEADER FROM BODY: FIND DISJOINT SETS OF <TR>S

成田線(千葉県)の賃貸住宅 (賃貸マンション、賃貸アパート、賃貸一戸建て) - Yahoo!不動産 - Mozilla Firefox

http://rent.realestate.yahoo.co.jp/bin/rsearch?md=

821件中241~260件を表示しています。 前へ 9 10 11 12 13 14 15 16 17 18 次へ

一覧表示 間取り表示 選択した物件をまとめて 詳細表示 お問い合わせ

画像	交通 住所	バス 徒歩	賃料 管理費等	礼金(保証金) 敷金(敷引)	間取り 専有面積	築年月 (築年数)	詳細	選択
	成田線/安食 印旛郡栄町安食	- 7分	5.00万円 2,000円	なし 1か月	2DK 44.72m ²	'95/02 (築14年)	▶ 詳細を見る	<input type="checkbox"/>
	成田線/安食 印旛郡栄町安食	- 7分	5.00万円 2,000円	なし 1か月	2DK 44.72m ²	'95/02 (築14年)	▶ 詳細を見る	<input type="checkbox"/>
	総武本線/佐倉 佐倉市大崎台2丁目	- 12分	5.00万円 2,000円	なし 2か月	4K 42.97m ²	'97/03 (築11年)	▶ 詳細を見る	<input type="checkbox"/>
	成田線/湖北 我孫子市古戸	- 19分	4.90万円 3,000円	なし 2か月	2DK 40.57m ²	'93/02 (築16年)	▶ 詳細を見る	<input type="checkbox"/>
	成田線/湖北 我孫子市古戸	- 19分	4.90万円 3,000円	なし 2か月	2DK 40.57m ²	'93/02 (築16年)	▶ 詳細を見る	<input type="checkbox"/>
	成田線/湖北 我孫子市古戸	- 19分	4.90万円 3,000円	なし 2か月	2DK 40.57m ²	'93/02 (築16年)	▶ 詳細を見る	<input type="checkbox"/>

http://rent.realestate.yahoo.co.jp/

FIND DISJOINT SETS OF <TR>S: NESTED BODY

【アットホーム】借りる（賃貸）：賃貸アパート・賃貸マンション・賃貸住宅：検索結果（一覧で表示）【東京都（JR京葉線 新浦安）】 | 住まいを探す

ファイル(E) 編集(E) 表示(V) 履歴(S) ブックマーク(B) ツール(I) ヘルプ(H)

http://www.athome.co.jp/apps/sch/list/kr_01_rail_12_13_2164_2164080_default 賃貸

Firefox を使ってみよう 最新ニュース

COMPST... 再現性 - ... すこしJava 1.1.1. W... 2005/02... The Nev... 気象庁 | ... 【ア...

一覧で表示 間取図で表示 地図で表示 情報の見方

選択をすると、まとめてお問合せまたは検討中リストへ追加ができます。 すべての物件を選択する

交通 ▲ 所在地 ▼	駅徒歩 ▲ (バス停徒歩)	賃料 ▲ (管理費等)	敷金/保証金 礼金	間取り ▼ 面積 ▼	物件種目 ▼ 築年月 ▼	画像 ▼	選択
・JR京葉線/新浦安 浦安市入船4丁目	10分	4.5万円 (3,000円)	2ヶ月 / なし なし	1K 16.00m ²	貸アパート 1985年10月		<input type="checkbox"/>
・ビタットハウス新浦安店 スタートビタットハウス(株) (JR京葉線/新浦安 徒歩4分)						TEL:047-354-9011	
・JR京葉線/新浦安 浦安市海楽2丁目	12分	4.6万円 (2,000円)	1ヶ月 / なし 1ヶ月	ワンルーム 16.29m ²	貸マンション 1991年10月		<input type="checkbox"/>
・ユニオン不動産(株) (東京外口東西線/浦安 徒歩1分) 浦安で住まいをお探しの際は是非御来店下さい						TEL:047-351-2300	
・JR京葉線/新浦安 浦安市海楽1丁目	14分	4.8万円 (1,000円)	2ヶ月 / なし なし	1K 19.61m ²	貸アパート 1985年11月		<input type="checkbox"/>
・住宅情報館(株)ファーストジョイ (東京外口東西線/浦安 徒歩1分) ワンルームからファミリー物件まで物件豊富！						TEL:047-355-1511	
・JR京葉線/新浦安 浦安市海楽2丁目	15分	4.8万円 (-)	2ヶ月 / なし 1ヶ月	ワンルーム 17.00m ²	貸アパート 1987年8月		<input type="checkbox"/>
・(株)月和地所 浦安駅前支店 (東京外口東西線/浦安 徒歩1分) 何でも揃う「住まいのコンビニ」としてご利用下さい！						TEL:0120-380-794	

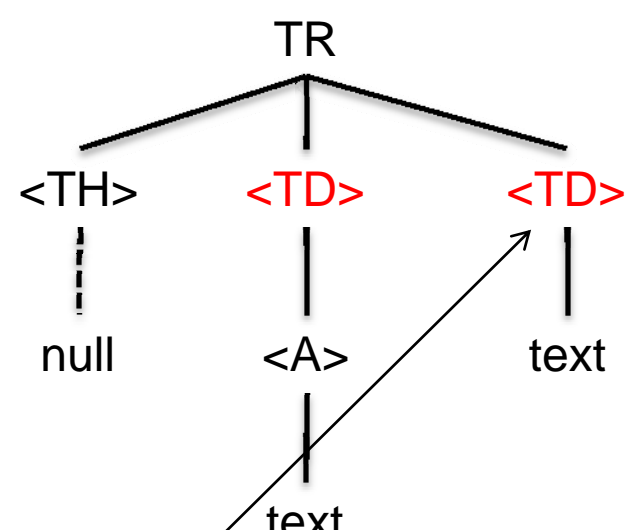
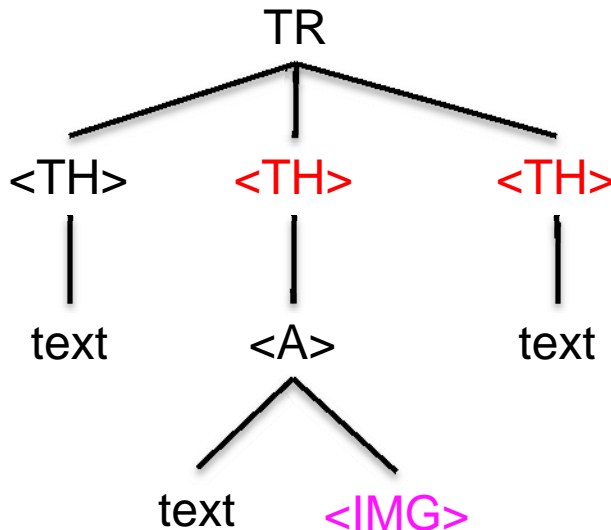
http://www.athome.co.jp/

FIND DISJOINT SETS OF <TR>S: NESTED HEADER

毎日の全国データ一覧表(日別値詳細版:2008年03月07日 17時現在)

地点	気圧		気温				蒸気圧	湿度			風向・風速				日照時間	全天日射量	平均雲量	降水量			降雪の深さ合計	最深積雪	天気概況	
	現地	海面	平均	最高	最低	平均	平均	最小	平均	最大		最大瞬間		日合計				日最大		1時間			10分間	1800-1800
	平均	平均		値	起時					値	起時	値	起時		風速	風速	風向	起時	風速		風向	起時		
札幌		1015.6	15:52	3.9	13:32	-4.7	05:44		40	12:29	6.2	北北西	14:10	9.0	北北西	14:16	5.0	13.2	0.0	—	3	76		
稚内		1013.1	17:00	3.2	12:13	-3.4	06:09		66	12:30	7.0	西南西	13:00	11.4	西	12:33	9.8	15.4	0.0	—	3	35		
北見枝幸		1013.3	16:26	3.9	11:53	-2.8	06:00		44	13:23	6.3	西	14:30	10.6	西	16:06	5.1		0.0	—	2	65		
旭川		1015.1	16:07	2.1	14:59	-6.7	05:02		56	14:54	5.9	西南西	15:20	9.2	西南西	15:20	3.5	12.1	0.5		4	61		
留萌		1015.2	17:00	2.5	14:28	-3.6	06:20		60	15:06	6.6	西南西	16:20	9.8	西	14:48	4.5	11.6	5.0		10	42		
羽幌		1014.7	17:00	2.9	14:23	-3.8	04:11		58	15:01	5.9	西	15:10	8.9	西	15:27	5.7		1.0		2	70		
岩見沢		1015.8	16:07	3.0	14:41	-5.6	03:25		57	14:36	3.3	西南西	13:50	5.9	北西	15:00	6.7		0.0	—	—	88		
小樽		1015.6	16:12	3.7	12:47	-2.9	03:32		39	12:50	4.2	西南西	17:00	7.7	西南西	16:32	6.8		0.0	—	1	96		
寿都		1016.5	15:56	3.5	12:32	-2.9	06:21		47	12:49	3.7	北西	13:30	9.8	北西	13:36	5.3	11.2	1.0		2	37		
倶知安		1016.1	16:07	2.8	12:18	-7.5	05:27		46	12:46	6.2	西北北	15:00	9.8	西北北	14:44	7.4		3.5		8	171		
網走		1013.7	17:00	2.2	13:35	-4.1	05:30		41	14:49	4.2	北西	13:00	7.6	西北北	13:07	1.4	12.0	0.0	—	6	62		
紋別		1013.2	15:49	3.3	12:34	-3.1	03:34		43	12:20	6.2	西	16:30	9.7	西	16:25	3.7		0.0	—	—	41		
雄武		1013.3	16:59	3.8	14:14	-2.8	03:06		42	13:55	6.4	西南西	14:50	10.2	西	14:48	3.9		0.0	—	—	41		
根室		1013.8	15:53	1.5	13:14	-4.4	05:58		59	13:15	4.5	北北西	14:40	6.6	北北西	14:32	3.9	13.6	0.0	—	—	—		
釧路		1014.3	14:52	1.4	12:59	-4.6	05:31		67	13:21	6.7	南西	13:30	8.5	南西	12:49	8.4		0.0	—	—	—		
帯広		1013.4	15:39	4.8	14:01	-6.7	06:20		37	13:29	6.8	西北北	15:10	11.3	西北北	14:59	10.2	17.3	0.5		2	16		
広尾		1014.1	15:49	1.7	16:33	-6.5	06:01		59	16:33	3.4	西	06:00	5.0	東南東	11:36	10.6		0.0	—	—	47		
室蘭		1016.2	16:35	3.3	14:31	-2.1	05:36		60	16:11	5.5	西	13:40	7.6	西北北	14:12	9.9	17.1	0.0	—	3	9		
苫小牧		1016.0	16:06	2.4	16:07	-7.3	06:23		60	16:17	4.1	南東	11:20	5.2	北	13:11	7.3		0.0	—	1	8		
蒲河		1015.8	16:35	2.5	14:54	-5.2	05:14		61	09:12	8.6	西	13:30	12.1	西北北	15:18	10.4	17.2	0.0	—	—	—		
函館		1016.6	16:27	4.2	12:44	-4.6	04:44		38	12:21	6.6	北西	14:20	10.0	西北北	14:09	7.0	14.2	0.5		3	18		
江差		1017.1	16:06	5.0	12:27	-0.5	07:50		43	12:46	6.6	西北北	04:00	12.2	北西	00:52	6.9		1.5		1	1		
青森		1017.4	14:54	4.8	13:12	-2.4	05:58		54	15:01	6.4	西南西	10:30	10.7	西北北	14:58	4.7	13.2	0.5		2	43		
八戸		1017.1	14:43	4.2	10:21	-2.3	06:00		48	10:24	5.6	南西	10:00	10.5	南西	10:59	4.7		0.0	—	—	1		
深浦		1018.4	15:32	4.8	12:35	-1.0	08:53		48	12:00	7.8	北西	08:50	14.6	北西	08:43	0.8		3.0		5	22		
むつ		1017.0	16:19	2.6	14:45	-5.6	01:48		66	14:43	4.2	南南西	13:10	7.9	西	08:06	3.0		3.5		8	11		
秋田		1018.9	16:38	3.1	11:43	-2.4	02:13		75	10:54	7.5	西北北	11:50	14.0	北北西	11:50	3.2	10.8	8.5		11	5		
盛岡		1017.4	15:00	5.3	12:29	-5.0	06:22		37	12:23	6.4	西	13:10	10.6	南西	12:50	6.6	14.6	0.5		—	—		
大船渡		1016.8	14:52	7.0	13:41	-2.3	05:42		32	13:49	5.6	北北西	15:30	11.4	西北北	15:19	7.1		0.0	—	—	—		
宮古		1017.0	13:53	6.6	11:07	-3.8	06:08		46	13:13	5.3	東南東	13:40	7.6	東南東	13:32	1.8		0.5		2	3		
仙台		1017.6	14:29	8.0	11:32	-0.3	05:27		32	11:37	9.2	西北北	14:10	13.9	西北北	15:40	7.5	15.8	0.0	—	—	—		
石巻		1017.3	15:21	8.1	11:22	-0.1	05:43		31	14:40	10.1	西北北	15:00	15.0	西北北	15:55	9.1		0.0	—	—	—		

DISSIMILARITY MEASURE BETWEEN DISJOINT <TR>S

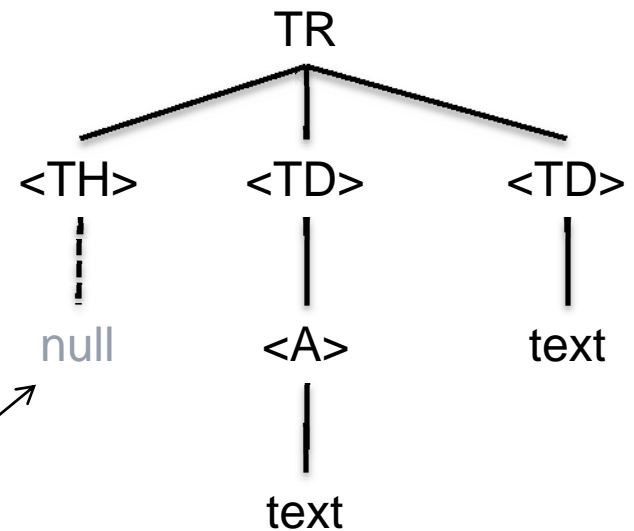
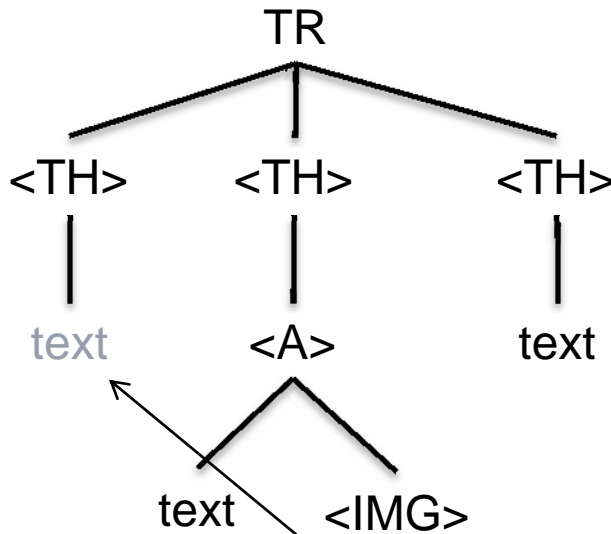


↑
Additional Node

↙ ↘
Different Node

Count the number of differences

DISSIMILARITY MEASURE BETWEEN DISJOINT <TR>S



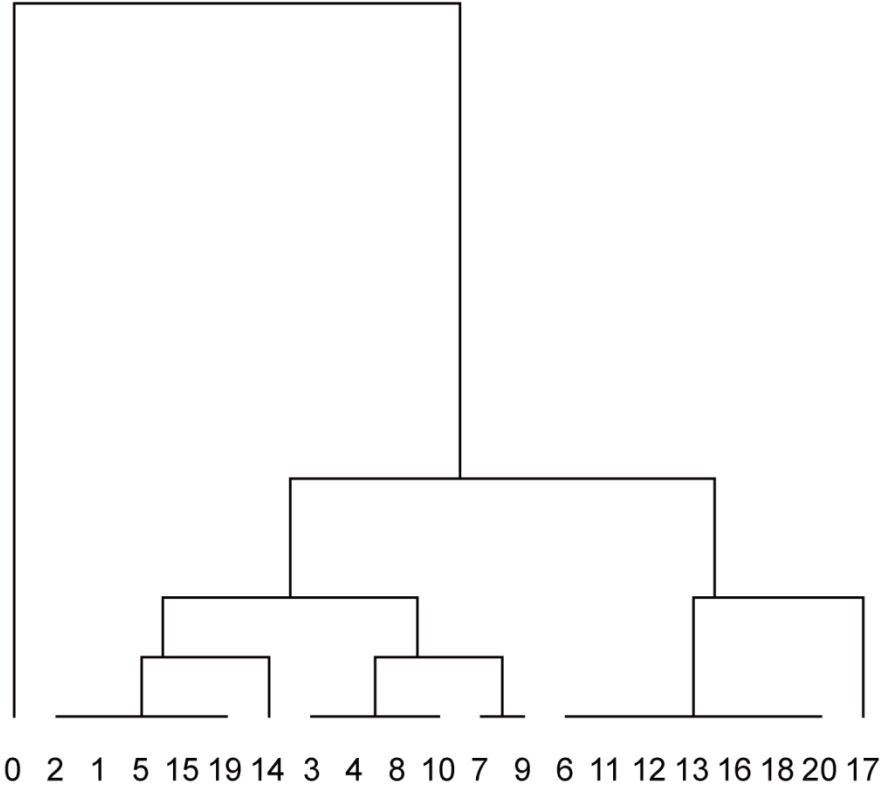
Text Node and Null Node

CLUSTERING ALGORITHM

駅名	所在地	面積 ㎡	築年数 年	価格 万円	築年数 管理費 円	駅別 徒歩	画像の有 無	物件の 詳細
下北沢駅	渋谷区九山町	9.72㎡	41年	3.2万円	2,000円	2ヶ月 2ヶ月	マンション 鉄骨造	写真あり
下北沢駅	世田谷区代田6丁目	9.1㎡	44年	3.6万円	2,000円	1ヶ月 1ヶ月	アパート 木造	写真あり
下北沢駅	世田谷区北沢4丁目	11.38㎡	28年	4.5万円	3,000円	2ヶ月 2ヶ月	マンション 鉄骨造	写真あり
下北沢駅	世田谷区北沢4丁目	10.84㎡	28年	4.5万円	2,000円	2ヶ月 2ヶ月	アパート 木造	写真あり
下北沢駅	世田谷区羽根木1丁目	12.5㎡	20年	5.2万円	7,000円	1ヶ月 1ヶ月	マンション RC	写真あり
下北沢駅	世田谷区羽根木1丁目	12.25㎡	28年	5.2万円	4,500円	1ヶ月 1ヶ月	マンション RC	写真あり
下北沢駅	世田谷区代田5丁目	11㎡	21年	5.2万円	2,000円	2ヶ月 2ヶ月	アパート 木造	写真あり
下北沢駅	世田谷区代田5丁目	11㎡	22年	5.2万円	2,000円	2ヶ月 2ヶ月	アパート 木造	写真あり
下北沢駅	世田谷区代田5丁目	11㎡	22年	5.2万円	2,000円	2ヶ月 2ヶ月	アパート 木造	写真あり
下北沢駅	世田谷区代田5丁目	11㎡	20年	5.2万円	2,000円	2ヶ月 2ヶ月	アパート 木造	写真あり
下北沢駅	世田谷区羽根木1丁目	14.5㎡	15年	5.27万円	2,000円	2ヶ月 0.5ヶ月	アパート 木造	写真あり
下北沢駅	世田谷区羽根木1丁目	17.23㎡	15年	5.4万円	8,000円	1ヶ月 1ヶ月	マンション RC	写真あり
下北沢駅	世田谷区代田6丁目	17㎡	31年	5.5万円	2,000円	2ヶ月 2ヶ月	アパート 木造	写真あり
下北沢駅	世田谷区羽根木1丁目	12.28㎡	20年	5.5万円	8,100円	1ヶ月 1ヶ月	マンション RC	写真あり
下北沢駅	世田谷区代田5丁目	15㎡	28年	5.6万円	2,000円	2ヶ月 2ヶ月	アパート 木造	写真あり
下北沢駅	世田谷区代田5丁目	15㎡	28年	5.6万円	2,000円	2ヶ月 2ヶ月	アパート 木造	写真あり
下北沢駅	世田谷区代田5丁目	13㎡	20年	5.6万円	2,000円	2ヶ月 2ヶ月	アパート 木造	写真あり
下北沢駅	世田谷区代田5丁目	13㎡	21年	5.6万円	2,000円	2ヶ月 2ヶ月	アパート 木造	写真あり
下北沢駅	世田谷区代田5丁目	11㎡	22年	5.6万円	2,000円	2ヶ月 2ヶ月	アパート 木造	写真あり

Header

Body



Header

Body

<http://chintai.homes.co.jp/>

MERGE <TABLE>S ACCORDING TO HEADER INFORMATION

Same headers

Different bodies

The screenshot shows a web browser window displaying a large table of data. The table is divided into several sections, each with a red box highlighting its header row. The headers are identical across all sections, but the data rows (bodies) are different. The browser window title is "気象庁 | 毎日の全国データ一覧表 日別値" and the URL is "http://www.data.jma.go.jp/obd/stats/data/mdr/synopday/data1s.html". The table contains columns for date, time, and various meteorological data points. The red boxes highlight the header rows of the first, second, third, and fourth sections of the table.

CHECK CONSISTENCY IN DATAVECTOR

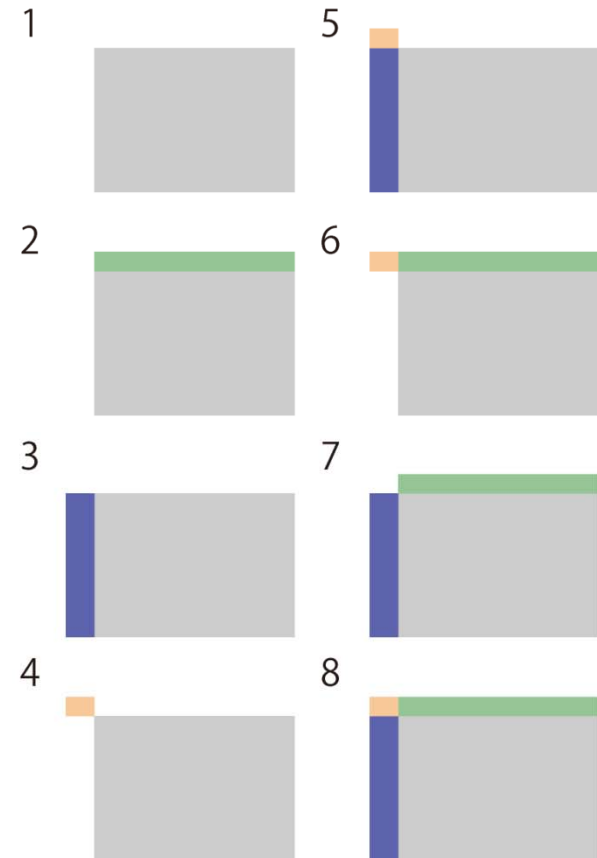
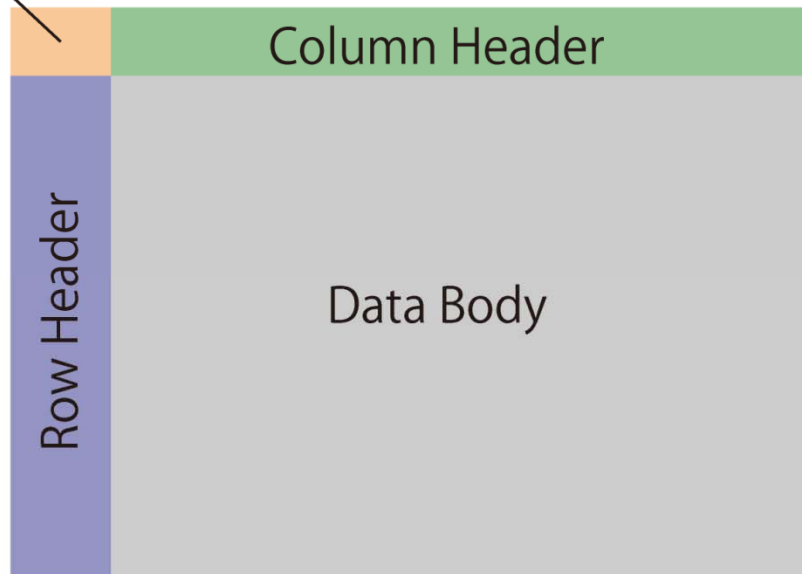
- Find Missing / Invalid
 - “NA”, “-”, “?”, “”, “.”, space, null
- Estimate Data Type
 - Numeric (1.34, 12,000,...)
 - Check data range
 - Double -> Measurement
 - Integer -> Ordinal Number
 - Non-negative Integer -> Count
 - {0, 1} -> Logical
 - Semi-numeric (2000yen, 5min,...)
 1. Find pattern (e.g. “2,000 yen”, “3,000 yen”)
 2. Separate value and unit (e.g. “2,000yen” -> “2,000”, “yen”)
 3. Go to numeric
 - Character (Hiyoshi, Tsunashima,...)
 - Check repeated values -> Category; otherwise -> ID

EXCEL

- Extract tables from a sheet
 1. Create Sheet Object
 - JExcelAPI (JExcel API group, 2008)
 2. Find disjoint sets of cells
 3. Determine rectangle hull for each set
- Separate header and body
 1. Check table type
 2. Cell type clustering
- Merge tables
- Divide into datavectors and check consistency

TARGET TABLE TYPE

Table Header



CHECK TABLE TYPE

EUR	1-Dec-99	2-Dec-99	3-Dec-99	6-Dec-99	7-Dec-99	8-Dec-99
1w	3.07375	3.07563	3.07000	3.07375	3.07188	3.06000
1m	3.51063	3.51963	3.52125	3.53125	3.54500	3.54000
2m	3.45000	3.45000	3.45000	3.46000	3.46000	3.45875
3m	3.45000	3.45000	3.45000	3.46000	3.46000	3.45875
4m	3.45000	3.45000	3.45575	3.46000	3.46000	3.46000
5m	3.46000	3.46000	3.46463	3.46813	3.47000	3.47000
6m	3.46938	3.47313	3.48000	3.48113	3.48000	3.48000
7m	3.50750	3.51000	3.51563	3.51625	3.51688	3.51875
8m	3.55125	3.55588	3.56000	3.56000	3.56000	3.56000
9m	3.60250	3.60063	3.60875	3.60713	3.60375	3.60000
10m	3.65000	3.65000	3.66000	3.65938	3.65000	3.65000
11m	3.70375	3.70875	3.71938	3.71400	3.71000	3.70875
12m	3.76375	3.76588	3.78000	3.77825	3.77000	3.77000

CHECK TABLE TYPE

Table header

EUR	1-Dec-99	2-Dec-99	3-Dec-99	4-Dec-99	7-Dec-99	8-Dec-99
1w	3.07375	3.07563	3.07000	3.07375	3.07188	3.06000
1m	3.51063	3.51963	3.52125	3.53125	3.54500	3.54000
2m	3.45000	3.45000	3.45000	3.46000	3.46000	3.45875
3m	3.45000	3.45000	3.45000	3.46000	3.46000	3.45875
4m	3.45000	3.45000	3.45575	3.46000	3.46000	3.46000
5m	3.46000	3.46000	3.46463	3.46813	3.47000	3.47000
6m	3.46938	3.47313	3.48000	3.48113	3.48000	3.48000
7m	3.50750	3.51000	3.51563	3.51625	3.51688	3.51875
8m	3.55125	3.55588	3.56000	3.56000	3.56000	3.56000
9m	3.60250	3.60063	3.60875	3.60713	3.60375	3.60000
10m	3.65000	3.65000	3.66000	3.65938	3.65000	3.65000
11m	3.70375	3.70875	3.71938	3.71400	3.71000	3.70875
12m	3.76375	3.76588	3.78000	3.77825	3.77000	3.77000

FREE FORMAT

- Separate header from body
 - Clustering technique
 - Distance between value types
 - Numeric (1.45, 25000, etc...)
 - Semi-Numeric (4,000yen, 5min, etc...)
 - Character (Hiyoshi, Shibuya, etc...)
 - Minimum distance algorithm
- Divide into datavectors and check consistency

RDB AND FIXED FORMAT

- Header and body are separated
- Value type of each column is given by
 - Metadata (VARCHAR, INTEGER, DOUBLE,...)
 - Format text (F3.1, S5,...)

FUTURE WORKS

- HTML documents generated by CGI and JavaScript
- Excel with Common header
- Free format file with a lot of missing values
- Large RDB

REFERENCES

Clark, A. (2008) CyberNeko THML Parser Home Page,
<http://nekohtml.sourceforge.net/>.

JExcelApi Group (2008) JExcel API Homepage:
<http://jexcelapi.sourceforge.net/>

Kumasaka, N. and Shibata, R. (2007) Textile Plot Environment, 統計数理特集号「統計データの可視化」, **55** pp. 47-68.

Yokouchi, D. and Shibata, R. (2004) DandD Client-Server System,
Proceedings in COMPSTAT 2004.

W3C (2008) DOM Homepage: <http://www.w3c.org/DOM/>