



POLITECNICO  
MILANO 1863

SYSTEMS AND METHODS FOR BIG AND UNSTRUCTURED DATA

# Wrangling

Marco Brambilla

[marco.brambilla@polimi.it](mailto:marco.brambilla@polimi.it)

 @marcobrambi



Marco Brambilla. **POLITECNICO**  
MILANO 1863

# Conventional Definition of Data Quality

## **Accuracy**

The data was recorded correctly.

## **Completeness**

All relevant data was recorded.

## **Uniqueness**

Entities are recorded once.

## **Timeliness**

The data is kept up to date (and time consistency is granted).

## **Consistency**

The data agrees with itself.



# Problems ...

## Unmeasurable

Accuracy and completeness are extremely difficult, perhaps impossible to measure.

## Context independent

No accounting for what is important. E.g., if you are computing aggregates, you can tolerate a lot of inaccuracy.

## Incomplete

What about interpretability, accessibility, metadata, analysis, etc.

## Vague

The conventional definitions provide no guidance towards practical improvements of the data.



Isn't data science sexy?

Wrangler®  
The Western Original

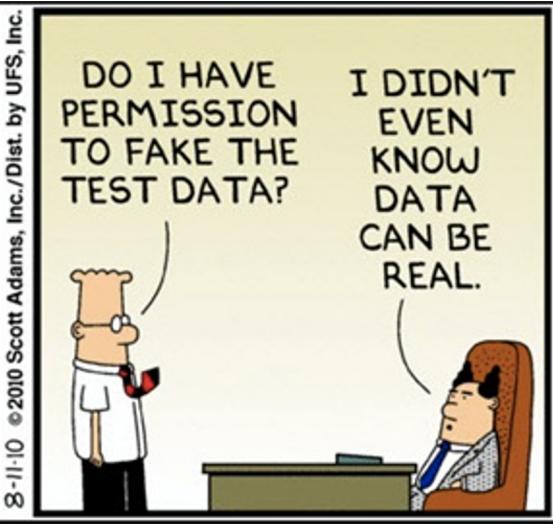




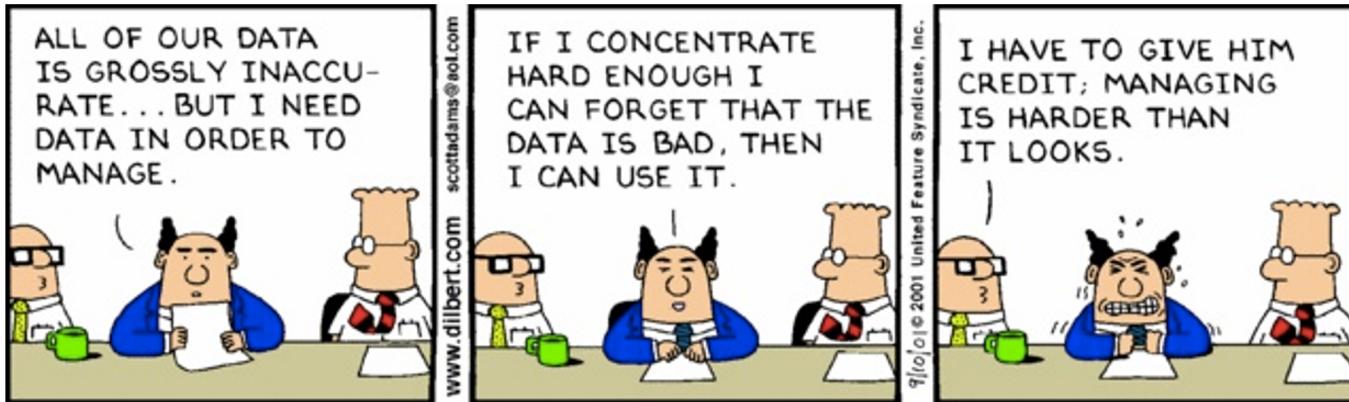
# 1. When Data Is Wrong



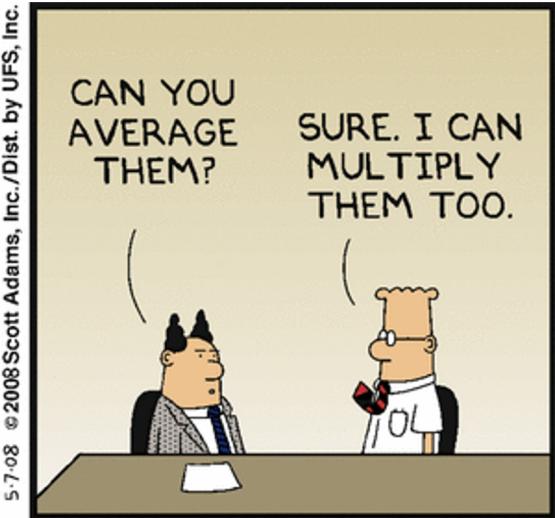
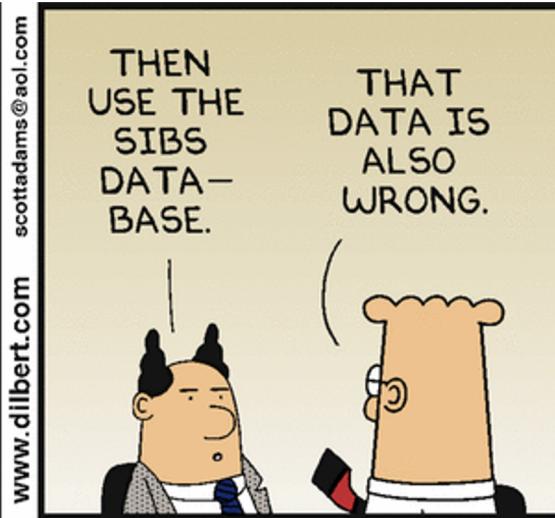
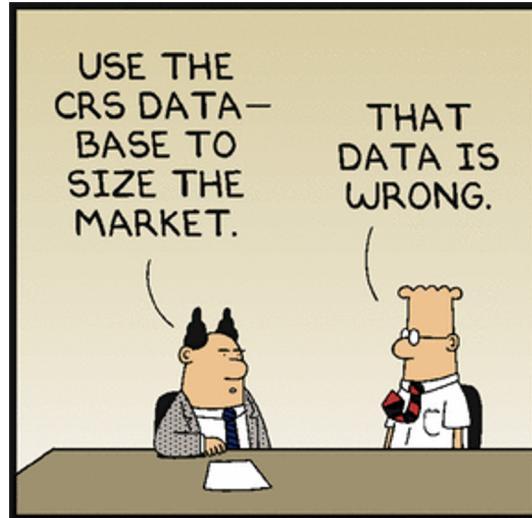
# The skeptic approach



# The pragmatic approach



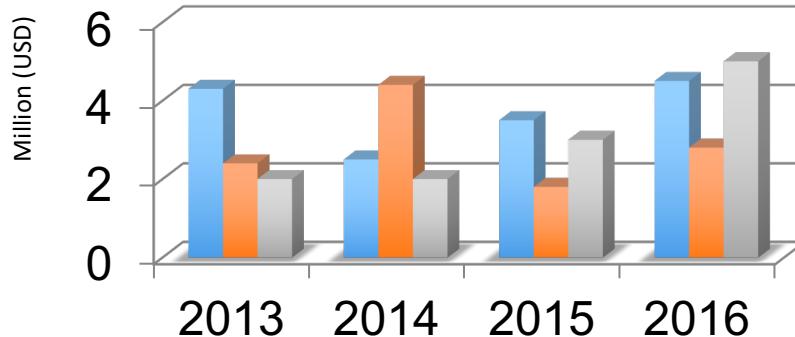
# The (pseudo) practitioner approach



[www.dilbert.com](http://www.dilbert.com)

5-7-08 © 2008 Scott Adams, Inc./Dist. by UFS, Inc.

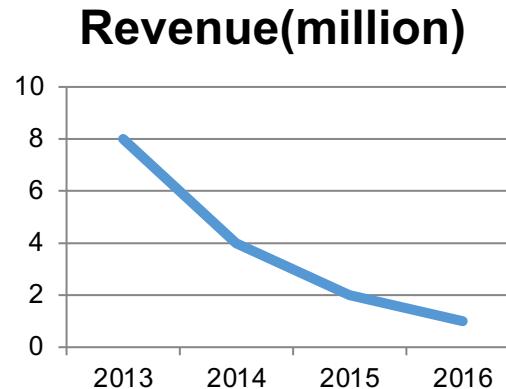
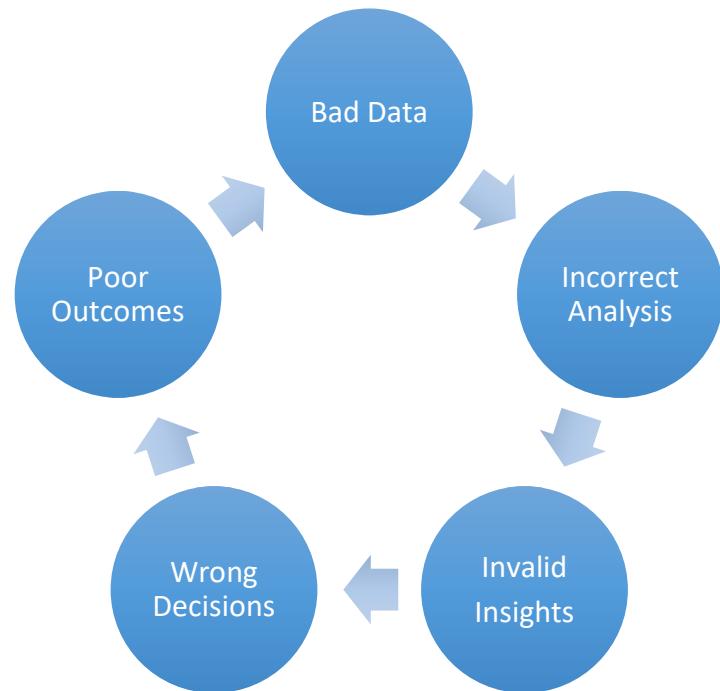
# Goal: Better Faster Cheaper!



\$\$\$



# The Vicious Cycle of Bad Data



*Data Quality is an issue...*



# Data Quality Issue

## Gartner Report

By 2017, 33% of the largest global companies will experience an information crisis due to their inability to adequately value, govern and trust their enterprise information.

*If you torture the data long enough,  
it will confess to anything*

*– Darrell Huff*





## 2. Making a Wrong Right



# Data Wrangling is ...



The process of transforming “raw” data into data that can be analyzed to generate valid actionable insights

Data scientists spend more time  
on preparing data than on analyzing it



# Data Wrangling a.k.a.

Data Preprocessing

Data Preparation

Data Cleansing

Data Scrubbing

Data Munging

Data Transformation

Data Fold, Spindle, Mutilate...

(good old) ETL



# ETL (Extract-Transform-Load)

Information architecture pattern for  
datawarehouse solutions

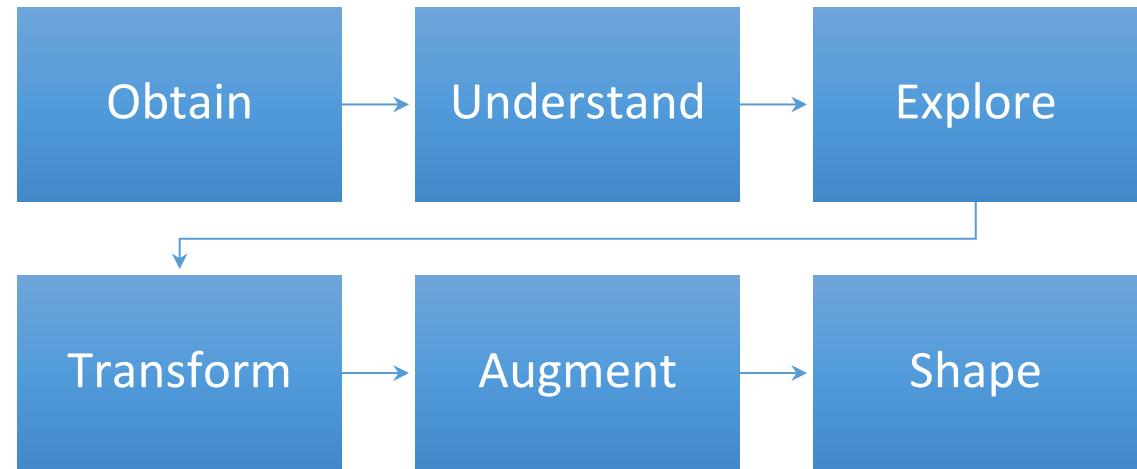
Batch oriented data integration and  
delivery



# Data Wrangling Steps

Iterative process

- Understand
- Explore
- Transform
- Augment
- Visualize



# What is Data Cleansing?

Data cleansing or data scrubbing is the act of detecting and correcting (or removing) corrupt or inaccurate records from a data set.

The term refers to identifying incomplete, incorrect, inaccurate, partial or irrelevant parts of the data and then replacing, modifying, filling in or deleting this dirty data.



# Why is Data “Dirty” ?

Dummy Values,  
Absence of Data,  
Multipurpose Fields,  
Cryptic Data,  
Contradicting Data,  
Shared Field Usage,

Inappropriate Use of Fields,  
Violation of Business Rules,  
Reused Primary Keys,  
Non-Unique Identifiers,  
Data Integration Problems



# Data Cleansing in Practice

- Parsing
- Correcting
- Standardizing
- Matching
- Consolidating



# Parsing

Parsing locates and identifies individual data elements in the source files and then isolates these data elements in the target files.



# Parsing

*Input Data from Source File*  
Beth Christine Parker, SLS MGR  
Regional Port Authority  
Federal Building  
12800 Lake Calumet  
Hedgewisch, IL



## *Parsed Data in Target File*

<b>First Name:</b>	Beth
<b>Middle Name:</b>	Christine
<b>Last Name:</b>	Parker
<b>Title:</b>	SLS MGR
<b>Firm:</b>	Regional Port Authority
<b>Location:</b>	Federal Building
<b>Number:</b>	12800
<b>Street:</b>	Lake Calumet
<b>City:</b>	Hedgewisch
<b>State:</b>	IL



# Correcting

Corrects parsed individual data components using sophisticated data algorithms and secondary data sources.



# Correcting

## Parsed Data

**First Name:** Beth  
**Middle Name:** Christine  
**Last Name:** Parker  
**Title:** SLS MGR  
**Firm:** Regional Port Authority  
**Location:** Federal Building  
**Number:** 12800  
**Street:** Lake Calumet  
**City:** Hedgewisch  
**State:** IL



## Corrected Data

**First Name:** Beth  
**Middle Name:** Christine  
**Last Name:** Parker  
**Title:** SLS MGR  
**Firm:** Regional Port Authority  
**Location:** Federal Building  
**Number:** 12800  
**Street:** South Butler Drive  
**City:** Chicago  
**State:** IL  
**Zip:** 60633  
**Zip+Four:** 2398



# Standardizing

Standardizing applies conversion routines to transform data into its preferred (and consistent) format using both standard and custom business rules, as well as coherent measurement units, ...



# Standardizing

## Corrected Data

First Name: Beth  
Middle Name: Christine  
Last Name: Parker  
Title: SLS MGR  
Firm: Regional Port Authority  
Location: Federal Building  
Number: 12800  
Street: South Butler Drive  
City: Chicago  
State: IL  
Zip: 60633  
Zip+Four: 2398

## Corrected Data

Pre-name: Ms.  
First Name: Beth  
**1st Name Match Standards:** Elizabeth, Bethany, Bethel  
Middle Name: Christine  
Last Name: Parker  
Title: Sales Mgr.  
Firm: Regional Port Authority  
Location: Federal Building  
Number: 12800  
Street: S. Butler Dr.  
City: Chicago  
State: IL  
Zip: 60633  
Zip+Four: 2398



# Matching

Searching and matching records within and across the parsed, corrected and standardized data based on predefined business rules to eliminate duplications.



# Match Patterns

Business Name	Street	Branch Type	Customer #/Tax ID	City	Vendor	Pattern	Pattern I.D.
Exact	Exact	Exact	Exact	Exact	Code Exact	AAAAAAA	P110
Exact	VClose	Exact	VClose	Exact	Blanks	ABAAA-	P115
Exact	VClose	Exact	Blanks	Exact	Exact	ABA-AA	P120
Exact	VClose	Close	Close	Exact	Exact	ABCCAA	S300
VClose	VClose	Exact	Close	Exact	Exact	BBACAA	S310



# Matching

## Corrected Data (Data Source #1)

Pre-name: Ms.  
First Name: Beth  
**1st Name Match**  
Standards: Elizabeth, Bethany, Bethel  
Middle Name: Christine  
Last Name: Parker  
Title: **Sales Mgr.**  
Firm: Regional Port Authority  
Location: Federal Building  
Number: 12800  
Street: S. Butler Dr.  
City: Chicago  
State: IL  
Zip: 60633  
Zip+Four: 2398

## Corrected Data (Data Source #2)

Pre-name: Ms.  
First Name: **Elizabeth**  
**1st Name Match**  
Standards: Beth, Bethany, Bethel  
Middle Name: Christine  
Last Name: **Parker-Lewis**  
Title:  
Firm: Regional Port Authority  
Location: Federal Building  
Number: 12800  
Street: **S. Butler Dr., Suite 2**  
City: Chicago  
State: IL  
Zip: 60633  
Zip+Four: 2398  
Phone: **708-555-1234**  
Fax: **708-555-5678**

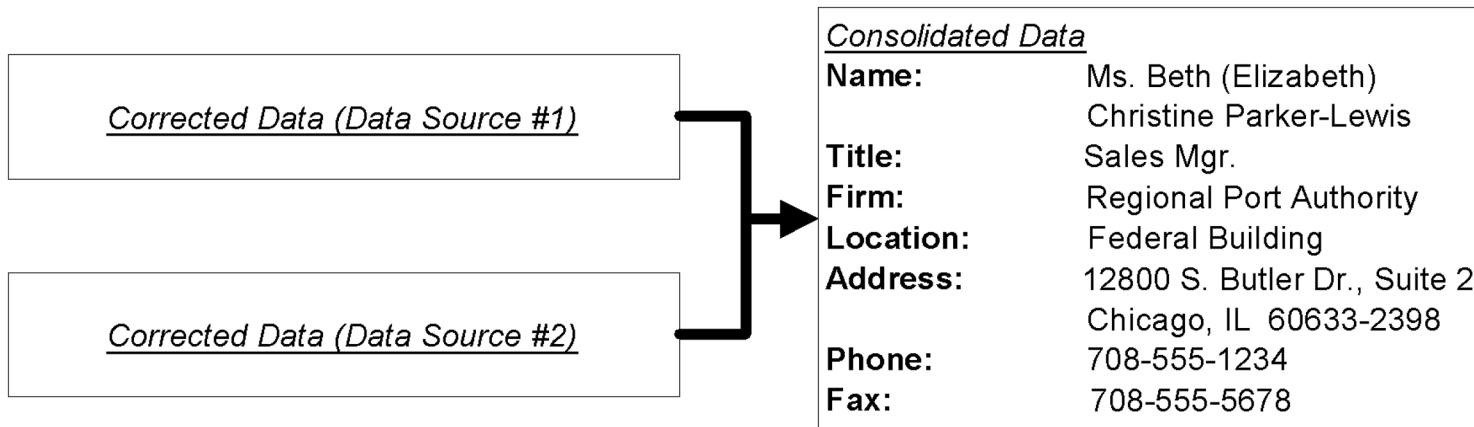


# Consolidating

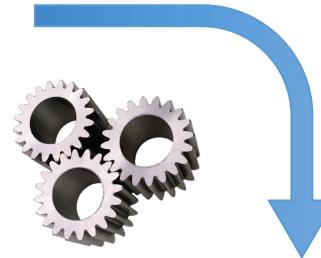
Analyzing and identifying relationships between matched records and consolidating/merging them into ONE representation.



# Consolidating



# Understanding Data: PDF



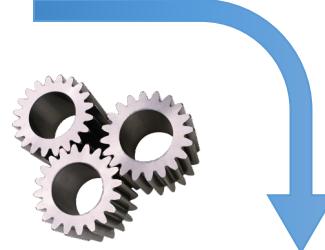
Company	Address	Bill To	Customer Address	Invoice	Date	Subtotal	Discount	Total	Description	Quantity	Unit	Amount
Musculinity @Last, Inc	123 Top Secret lane   BestLands, TE   84950	Justin BeefUpNow	1234 Healthy Street   Irving, TX, 75038   USA	1	5/6/2016	139.78	6.99	132.79	Naturally Extracted Magic Muscle Builder Powder	4	29.95	119.8
Musculinity @Last, Inc	123 Top Secret lane   BestLands, TE   84950	Justin BeefUpNow	1234 Healthy Street   Irving, TX, 75038   USA	1	5/6/2016	139.78	6.99	132.79	Pre Workout Super Shake	2	9.99	19.98



# Understanding Data: Free Text



Python + Textract +Tesseract



```
>>> text = textract.process('/Users/akuntamu/Downloads/Drawing3.png',method='tesseract-ocr')
>>> text
'    \n    \nhow is it going?\nThis is pretty\ninteresting\n\n'
```



# Understanding Data: more

“Looks like my V8 Chevy is running low on fuel.  
Didn’t I fill up just the day before?”

UNSTRUCTURED



STRUCTURED

Owner	Vehicle	Type	Fuel Level	Engine	Last Fill
AK	Chevy	Gas	5%	V8	05/04/16



# Understanding Data: more

Decode the following secret message:

DALDFWSFOEWRBOSDCALAXORDJFKMCO



DAL DFW SFO EWR BOS  
DCA LAX ORD JFK MCO



# Data Munging

Potentially lossy transformations applied to a piece of data or a file

Vague data transformation steps that are not yet completely clear

E.g., removing punctuation or html tags, data parsing, filtering, and transformation.



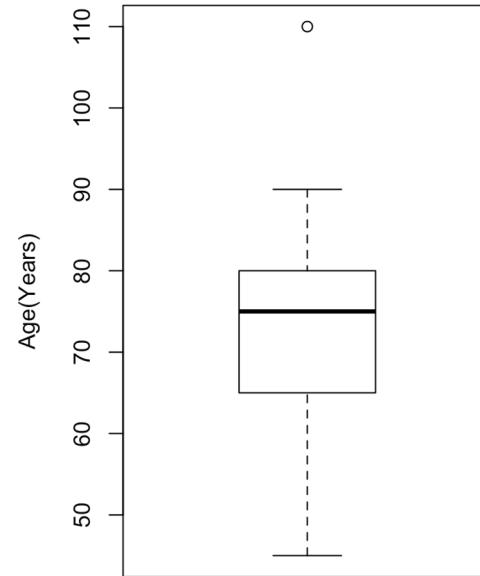
# Semantics

???
75
80
65
55
67
78
88
90
45
58
69
80
110

?

# Semantics and Outliers

Age(Years)
75
80
65
55
67
78
88
90
45
58
69
80
110



The value stands in the abnormal

# Missing Data: Detection

## Overtly missing data

- Match data specifications against data - are all the attributes present?
- Scan individual records - are there gaps?
- Rough checks : number of files, file sizes, number of records, number of duplicates
- Compare estimates (averages, frequencies, medians) with “expected” values and bounds; check at various levels of granularity since aggregates can be misleading



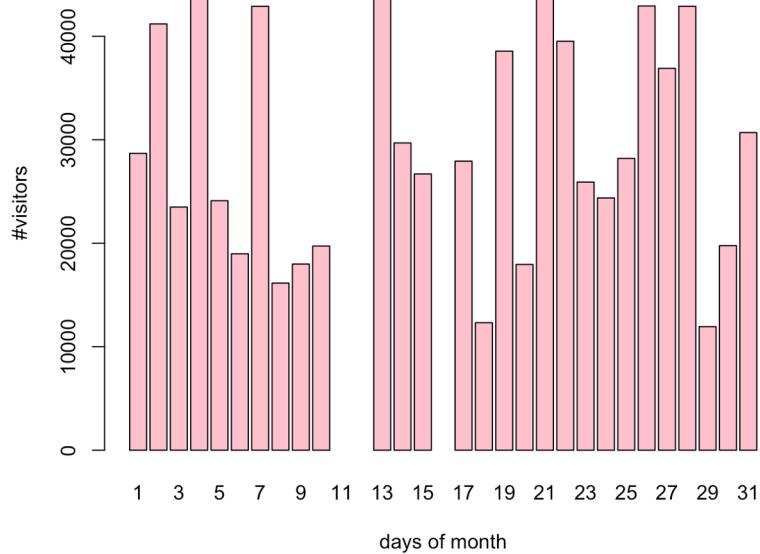
# Missing data: Detection (cont.)

## Hidden damage to data

- Values are truncated or censored - check for spikes and dips in distributions and histograms
- Missing values and defaults are indistinguishable - too many missing values? metadata or domain expertise can help
- Errors of omission e.g. all calls from a particular area are missing - check if data are missing randomly or are localized in some way



# Missing Values: Random



System failures  
Complete miss



# Missing Values: Wrong Ingestion

CSV to table / excel

Merged fields

Missing fields

Antony Mc James IV, 123 Untidy Cir, Suite# 234, Cumbbersome, TX, 76849  
Miller Johnson, 1102 Messy Data St, Middletow, OH  
Betty Flyier 6483 Phew Lane, Apt A4, FixMe, CA, 91103



Missing due to invalid data and ingestion

Customer	Address Line 1	Address Line 2	City	State	Zip
Antony Mc James IV	123 Untidy Cir	Suite# 234	Cumbbersome	TX	76849
Miller Johnson	1102 Messy Data St	Middletow	OH		
Betty Flyier 6483 Phew Lane	Apt A4	FixMe	CA	91103	

# Missing Values: Inapplicability

Date of your most recent **IMMUNIZATIONS**:

Hepatitis A \_\_\_\_\_ Hepatitis B \_\_\_\_\_ Influenza (flu shot) \_\_\_\_\_ MMR \_\_\_\_\_ Tetanus (Td) \_\_\_\_\_  
Pneumovax (pneumonia) \_\_\_\_\_ Meningitis \_\_\_\_\_ Varicella (chicken pox) shot or Illness \_\_\_\_\_  
HPV \_\_\_\_\_ Zostavax (Shingles) \_\_\_\_\_

**HEALTH MAINTENANCE SCREENING TESTS:**

Lipid (cholesterol) & Sugar \_\_\_\_\_ Date \_\_\_\_\_ Abnormal?  Yes  No  
Sigmoidoscopy \_\_\_\_\_ or Colonoscopy \_\_\_\_\_ Date \_\_\_\_\_ Abnormal?  Yes  No

**Women:** Mammogram Date \_\_\_\_\_ Abnormal?  Yes  No Pap Smear Date \_\_\_\_\_ Abnormal?  Yes  No  
Have you ever had a **Dexa** Scan (for bone density)  Yes  No Abnormal? Date \_\_\_\_\_

**Women's Health History:** Number of: pregnancies \_\_\_\_\_ deliveries \_\_\_\_\_ abortions \_\_\_\_\_ miscarriages \_\_\_\_\_  
Age at start of periods: \_\_\_\_\_ Date of Last Period: \_\_\_\_\_ Age at end of periods: \_\_\_\_\_

**Name and phone number of your OB/Gynecologist:** \_\_\_\_\_

**Men's Health History:**  
Have you had a blood test for: PSA (prostate)  Yes  No If yes, what date \_\_\_\_\_ Was it normal? \_\_\_\_\_

**PERSONAL MEDICAL HISTORY:** Please indicate whether **YOU** have had any of the following medical problems (with dates).

Alcoholism \_\_\_\_\_  
Anxiety / Depression \_\_\_\_\_  
Cancer, specify type \_\_\_\_\_  
Heart disease \_\_\_\_\_  
Diabetes \_\_\_\_\_  
Stroke \_\_\_\_\_  
Urinary / Kidney Problems \_\_\_\_\_

Arthritis / Rheumatologic \_\_\_\_\_  
Bleeding or Clotting disorder or \_\_\_\_\_  
Blood Clot \_\_\_\_\_  
Thyroid problems \_\_\_\_\_  
High Blood Pressure \_\_\_\_\_  
High Cholesterol \_\_\_\_\_  
Other \_\_\_\_\_

**Do you see any other doctors?** If yes, please list their name, office phone number and specialty.

Partial data by nature

Remember to leave empty slots



# Imputing Values to Missing Data

In federated data, between 30%-70% of the data points will have at least one missing attribute - data wastage if we ignore all records with a missing value

Remaining data is seriously biased

Lack of confidence in results

Understanding pattern of missing data unearths data integrity issues



# Missing Value Imputation – 1

## Standalone imputation

Mean, median, other point estimates

Assume: Distribution of the missing values is the same as the non-missing values.

Does not take into account inter-relationships

Introduces bias

Convenient, easy to implement



# Missing Value Imputation - 2

Better imputation - use attribute relationships

Assume : all prior attributes are populated

That is, *monotonicity* in missing values.

X1	X2	X3	X4	X5
1.0	20	3.5	4	.
1.1	18	4.0	2	.
1.9	22	2.2	.	.
0.9	15	.	.	.

Two techniques

Regression (parametric),

Propensity score (nonparametric)



# Missing Value Imputation – 3

## Regression method

Use linear regression, sweep left-to-right

$$X_3 = a + b * X_2 + c * X_1;$$

$$X_4 = d + e * X_3 + f * X_2 + g * X_1, \text{ and so on}$$

$X_3$  in the second equation is estimated from the first equation if it is missing



# Missing Value Imputation - 3

## Propensity Scores (nonparametric)

Let  $Y_j=1$  if  $X_j$  is missing, 0 otherwise

Estimate  $P(Y_j = 1)$  based on  $X_1$  through  $X_{(j-1)}$  using logistic regression

Group by propensity score  $P(Y_j = 1)$

Within each group, estimate missing  $X_j$ s from known  $X_j$ s using approximate Bayesian bootstrap.

Repeat until all attributes are populated.



# Missing Value Imputation - 4

Arbitrary missing pattern

Markov Chain Monte Carlo (MCMC)

Assume data is multivariate Normal, with parameter  $\Theta$

(1) Simulate missing X, given  $\Theta$  estimated from observed X ; (2)  
Re-compute  $\Theta$  using filled in X

Repeat until stable.

Expensive: Used most often to induce monotonicity

Note that imputed values are useful in aggregates but can't be trusted individually



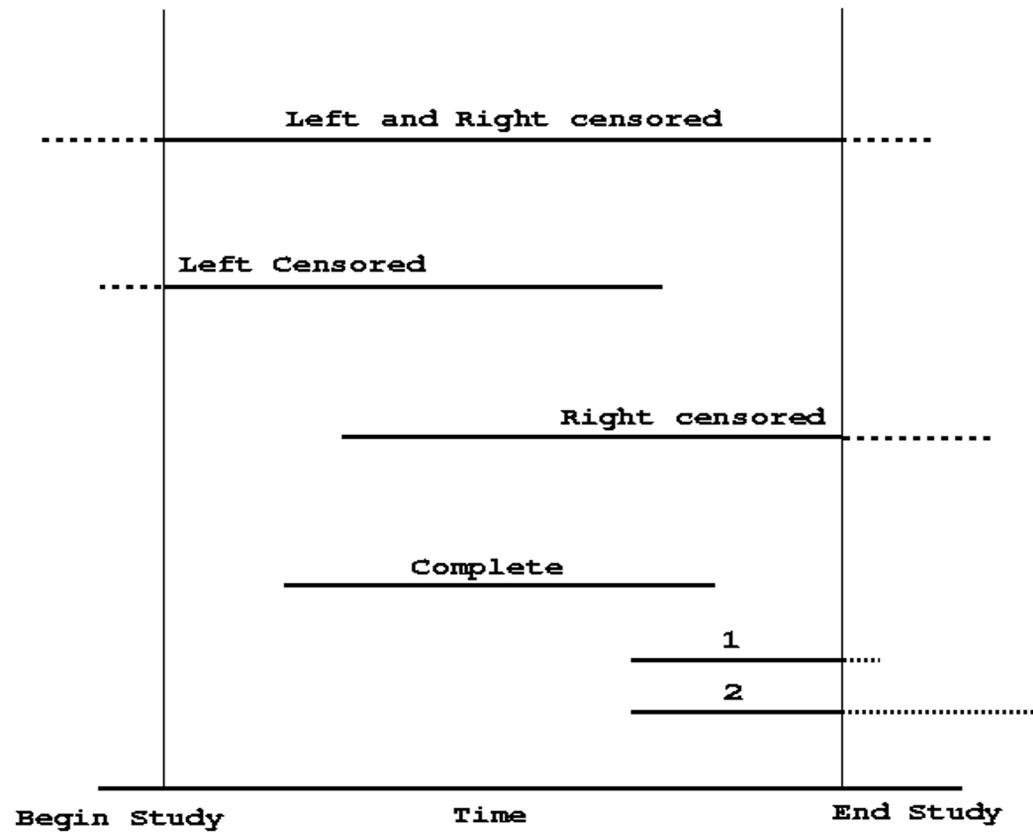
# Censoring and Truncation

Well studied in Biostatistics, relevant to time dependent data e.g. duration

*Censored* – Measurement is bounded but not precise e.g. Call duration  $> 20$  are recorded as 20

*Truncated* – Data point dropped if it exceeds or falls below a certain bound e.g. customers with less than 2 minutes of calling per month





## Censored time intervals

# Censoring/Truncation (cont.)

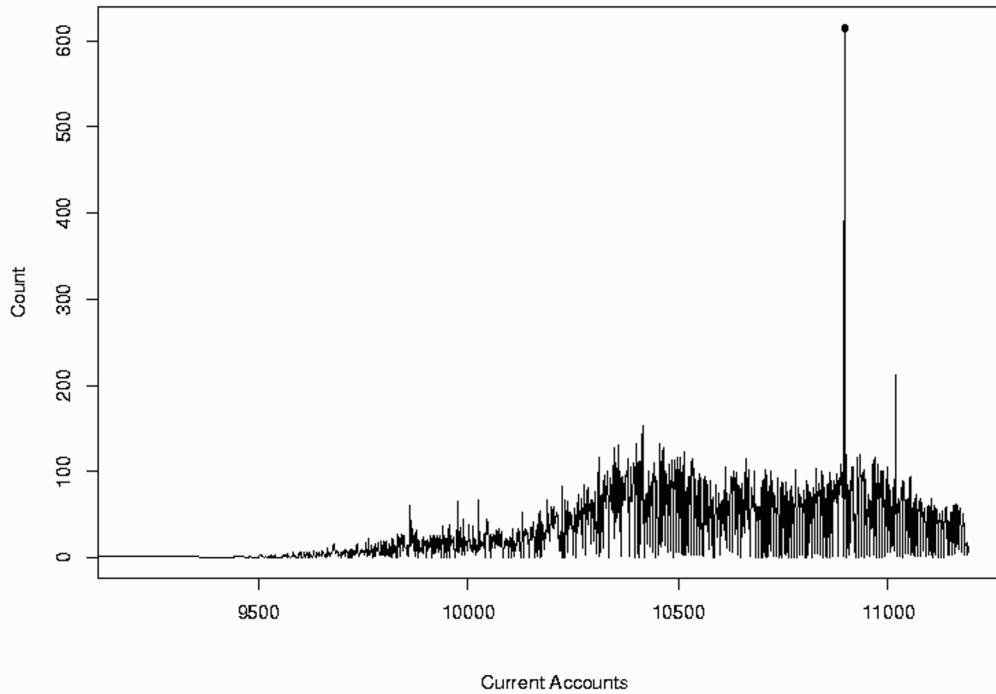
If censoring/truncation mechanism not known,  
analysis can be inaccurate and biased.

But if you know the mechanism, you can mitigate  
the bias from the analysis.

Metadata should record the existence as well as  
the nature of censoring/truncation



## Potential Data Problem



Spikes usually indicate censored time intervals  
caused by resetting of timestamps to defaults

Marco Brambilla.

# Suspicious Data

Consider the data points

3, 4, 7, 4, 8, 3, 9, 5, 7, 6, 92

“92” is suspicious – an *outlier*

Outliers are potentially legitimate

Often, they are data or model glitches

Or, they could be a data miner’s dream, e.g. highly profitable customers



# Outliers

Outlier – “departure from the expected”

Types of outliers – defining “expected”

Many approaches

Error bounds, tolerance limits – control charts

Model based – regression depth, analysis of residuals

Geometric

Distributional

Time Series outliers



# Control Charts

Quality control of production lots

Typically univariate: X-Bar, R, CUSUM

Distributional assumptions for charts not based on means e.g.  
R-charts

Main steps (based on statistical inference)

Define “expected” and “departure” e.g. Mean and standard error  
based on sampling distribution of sample mean (aggregate);

Compute aggregate each sample

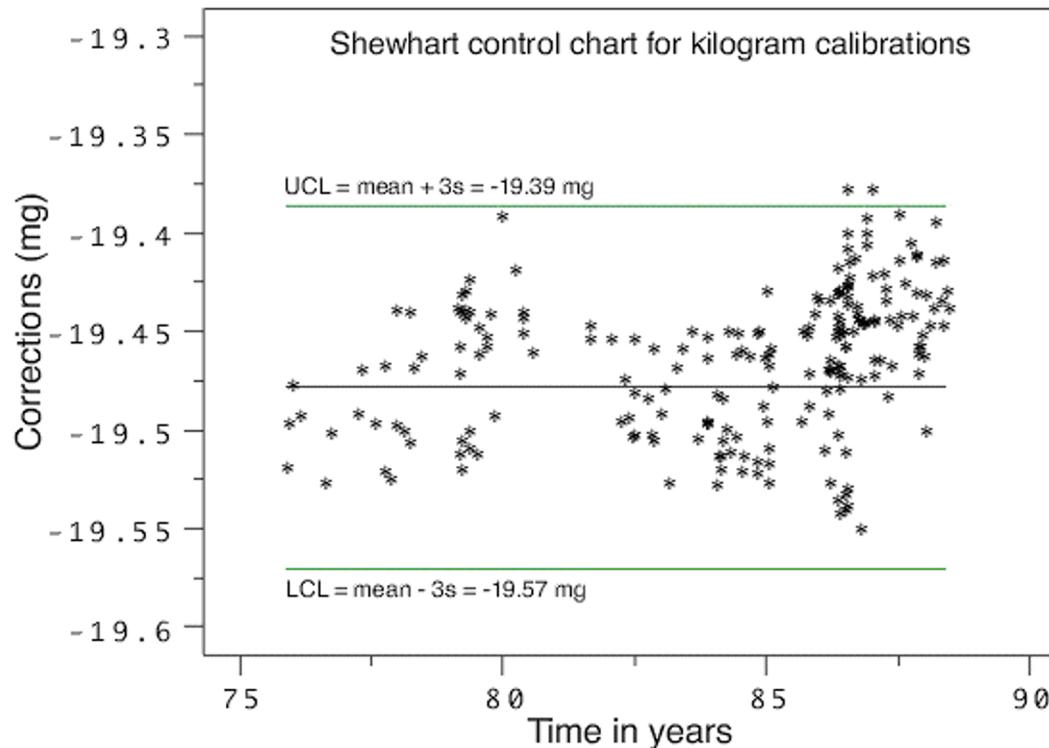
Plot aggregates vs expected and error bounds

“Out of Control” if aggregates fall outside bounds



# An Example

(<http://www.itl.nist.gov/div898/handbook/mpc/section3/mpc3521.htm>)



# Multivariate Control Charts - 1

Bivariate charts:

- based on bivariate Normal assumptions

- component-wise limits lead to Type I, II errors

Depth based control charts (nonparametric):

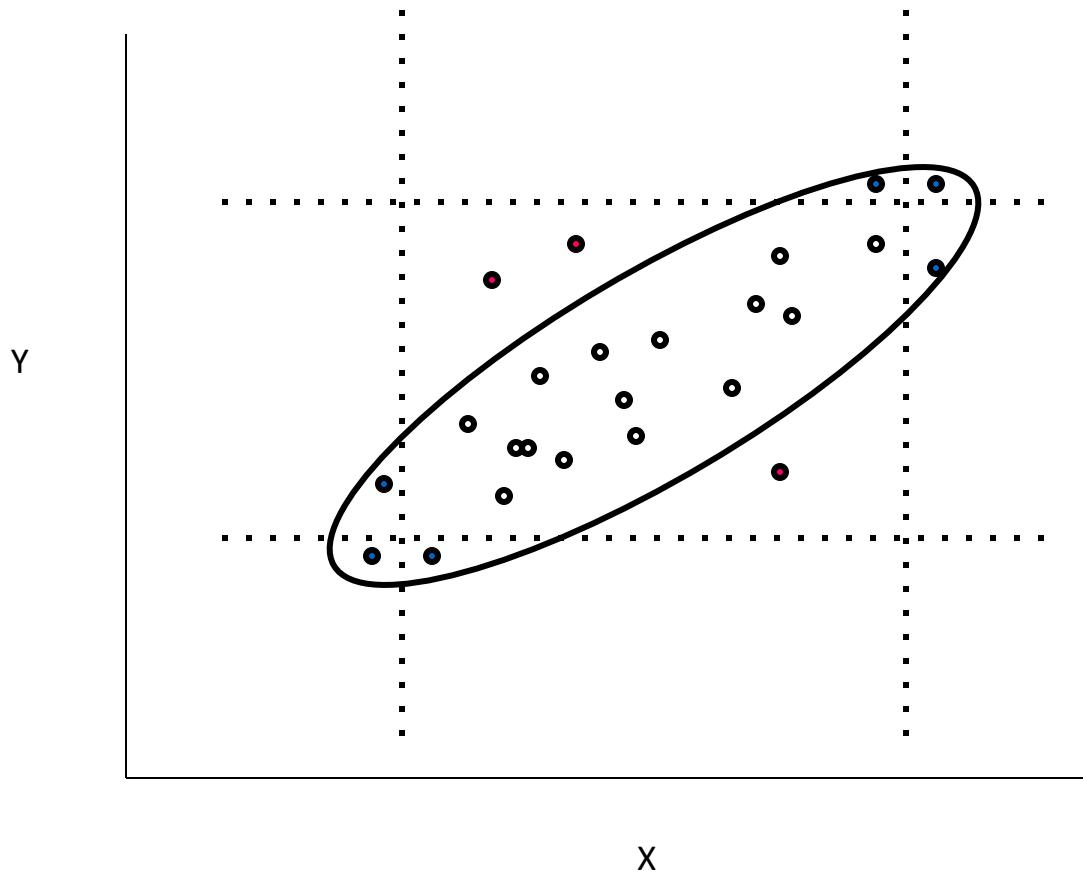
- map n-dimensional data to one dimension using depth e.g.  
Mahalanobis

- Build control charts for depth

- Compare against benchmark using depth e.g. Q-Q plots of depth  
of each data set



## Bivariate Control Chart



# Multivariate Control Charts - 2

Multiscale process control with wavelets:

- Detects abnormalities at multiple scales as large wavelet coefficients.

- Useful for data with heteroscedasticity

- Applied in chemical process control



# Model Fitting and Outliers

Models summarize general trends in data

more complex than simple aggregates

e.g. linear regression, logistic regression focus on attribute relationships

Data points that do not conform to well fitting models are  
*potential outliers*

Goodness of fit tests (DQ for analysis/mining)

check suitability of model to data

verify validity of assumptions

data rich enough to answer analysis/business question?



# Set Comparison and Outlier Detection

“Model” consists of partition based summaries

Perform nonparametric statistical tests for a rapid section-wise comparison of two or more massive data sets

If there exists a baseline “good” data set, this technique can detect potentially corrupt sections in the test data set



# Types of data

Categorical

Qualitative  
Subjective

Quantitative  
Discrete  
Continuous



# Types of data

Categorical

Qualitative  
Subjective

Quantitative  
Discrete  
Continuous

Color

Nice, Good, For birthday



6 balloons

Pressure 15 PSI,

139 m. over sea level



# Data Source Selection Criteria

Credibility

Completeness

Accurateness

Verifiability

Currency

Accessibility

Compliance

Cost

Legal issues

Security

Storage

Provenance



# Not all tables are created equal

The diagram illustrates two different ways of organizing data into tables:

- Row-major table (top):** This table has "School" as columns and "Year" as rows. The columns are labeled 2012, 2013, and 2014. The rows are labeled School (Good Samaritans and Percy Grammar). The data shows student counts for each school per year.
- Column-major table (bottom):** This table has "School" as rows and "Year" as columns. The rows are labeled School (Good Samaritans and Percy Grammar). The columns are labeled Year (2012, 2013, 2014). The data shows student counts for each year per school.

Annotations with arrows point to specific parts of the tables:

- A blue arrow labeled "Column" points to the "School" header in the first table.
- A blue arrow labeled "Row" points to the "School" header in the second table.
- A blue bracket labeled "year" spans the three "Year" columns in the first table.
- A blue arrow labeled "Variable" points to the "Year" header in the second table.
- A blue arrow labeled "Observation" points to the "School" header in the second table.

School	2012	2013	2014
Good Samaritans	2321	4550	1293
Percy Grammar	1540	1400	2949

School	Year	Student Count
Good Samaritans	2012	2321
Good Samaritans	2013	4550
Good Samaritans	2014	1293
Percy Grammar	2012	1540
Percy Grammar	2013	1400
Percy Grammar	2014	2949



# Not all tables are created equal

Year	Comedy-Q1	Thriller-Q1	Action-Q1	...
2014	2	1	0	
2015	0	3	2	

Find total comedy movies in all of 2014? -> Not easy in current form

Category	Quarter	Year	#Hits
Comedy	Q1	2014	2
Thriller	Q1	2014	1
Action	Q1	2014	0
Comedy	Q1	2015	0
Thriller	Q1	2015	3
Action	Q1	2015	2

Find % of hit comedy movies in a 2015?

Very easy to add a new column



# Not all tables are created equal

Category	Rating	Q1	Q2	Q3	...
Comedy	Excellent	1	0	1	
Comedy	Good	2	0	2	
Thriller	Excellent	0	1	1	
Thriller	Good	1	0	3	

Very messy data

Variables in both rows and columns

Category	Quarter	Excellent	Good
Comedy	Q1	1	2
Comedy	Q2	0	0
Comedy	Q3	1	2
Thriller	Q1	0	1
Thriller	Q2	1	0
Thriller	Q3	1	3

Each row is complete observation



# Not all tables are created equal

Invoice	Bill To	Sales %	Total(\$)	SKU#	Item	Qty	Unit Price (\$)
1	Jim Jones	8	8.03	A123	Hammer	1	3.55
1	Jim Jones	8	8.03	Q34	Screw Driver	2	2.05
2	Mike Z'Kale	8	97.20	W23	Hair Dryer	1	59.25
2	Mike Z'Kale	8	97.20	E452	Cologne	3	10.25

Normalize to avoid duplication

Invoice	Bill To	Sales %	Total(\$)	Invoice	SKU#	Item	Qty	Unit Price (\$)
1	Jim Jones	8	8.03	1	A123	Hammer	1	3.55
2	Mike Z'Kale	8	97.20	1	Q34	Screw Driver	2	2.05
				2	W23	Hair Dryer	1	59.25
				2	E452	Cologne	3	10.25



# Not all tables are created equal

Category	Quarter		Year	#Hits
	Category	Quarter	Year	#Hits
Comedy	Comedy	Q1	2014	2
Thriller	Thriller	Q1	2014	1
Action	Action	Q1	2014	0

Multiple Tables  
Divided by Time

Category	Quarter	Year	#Hits
Comedy	Q1	2014	2
Thriller	Q1	2014	1
Action	Q1	2014	0
Comedy	Q1	2014	2
Thriller	Q1	2014	1
Action	Q1	2014	0

Combine all tables  
accommodating  
varying formats



# The “Key” (matching) problem

Keys are crucial in DB

Many DBs --> Many keys

How to align?

Identification to certain  
degree of accuracy  
likely-identities

e.g., same user match



# The “Duplicates” problem

Related to Key problem

Identification to certain degree of accuracy  
likely-duplicates



e.g., duplicate posts



# The “Duplicates” problem

Related to Key problem

Identification to certain degree of accuracy  
likely-duplicates



# Lessons Learnt on Tables

(Multiple) variables in columns

Never values as columns!

Shape may depend on convenience of queries

Matching identities and duplications are crucial in data science!



# Lessons Learnt on Tables

Each observation is complete and atomic

Each variable belongs to (only!) one column!



# Schema-On-Write Vs Schema-On-Read

Traditional DBMSs enforced writing only data consistent with a pre-designed schema

Today data modeling at design time is a luxury

In schema on read, data is applied to a plan or schema as it is pulled out of a stored location, rather than as it goes in



# Popular Open Source Tools



IP[y] + python™ + pandas  
 $y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$

**Refine<sup>OPEN</sup>** DATA**SCIENCE**TOOLKIT



# Other Resources



BeautifulSoup

Tesseract

textract



OCRFeeder



# Commercial Vendors



# Triflacta Wrangler

Triflacta - Transformer

raw-top-50-pop

Full Datasource - 2.7KB | 5 Columns | 50 Rows | 3 Data Types | Grid | Filter in grid | Generate Results | ?

#	rank	ABC	artist	ABC	title	ABC	company	year
1	50	48 Categories	Björk	50 Categories	Vespertine	35 Categories	Elektra	2001
2	49	Libertines	Up The Bracket			Rough Trade		2002
3	48	Loretta Lynn	Van Lear Rose			Interscope		2004
4	47	Arctic Monkeys	Whatever People Say I Am, That's What I'm No	Domino			Columbia	2006
5	46	Once	Music From The Motion Picture			Self-released		2007
6	45	Radiohead	In Rainbows			Lost Highway		2003
7	44	The Jayhawks	Rainy Day Music			Secretly Canadian		2007
8	43	Jens Lekman	Night Falls Over Kortedala			Roc-A-Fella		2001
9	42	Jay-Z	The Blueprint			Capitol		2007
10	41	LCD Soundsystem	Sound of Silver			Interscope		2006
11	40	TV on the Radio	Return to Cookie Mountain			Merge		2007
12	39	Arcade Fire	Neon Bible			Douchemaster		2008
13	38	Gentleman Jesse	Introducing Gentleman Jesse					2004
14	37	Iron & Wine	Our Endless Numbered Days					
15	36	Pedro The Lion	Control		.....		Jade Tree	2002

TRANSFORM EDITOR

Enter transform expression | Add to Script

splitrows col: column1 on: ...  
replace col: column1 on: /...  
replace col: column1 on: /...  
delete row: matches([column1])  
split col: column1 on: {d...  
split col: column3 on: '{d...  
split col: column4 on: '[]'  
split col: column5 on: ']'  
extract col: column3 on: '...  
replace col: column3 on: '...  
merge col: column5,column6...  
drop col: column5,column6...  
set col: column7 value: tr...  
replace col: column7 on: '...  
replace col: column7 on: '...  
rename col: column2 to: 'r...  
rename col: column1 to: 'a...



# Google's Open Refine

# Hands on Data Wrangling

## Data Ingestion

- CSV
- PDF
- API/JSON
- HTML Web Scraping
- XLS, Access, ...!

## Data Exploration

- Visual inspection
- Graphing

## Data Shaping

- Tidying Data

## Data Cleansing

- Missing values
- Format
- Measurement Units
- Outliers
- Data Errors Per Domain
- Fat Fingered Data

## Data Augmenting

- Aggregate data sources
- Fuzzy/Exact match

# R Libraries for Data Wrangling



stringr

dplyr

tidyverse

readxl, xlsx

lubridate

gtools

plyr

rvest



# References – Web and Books

Web:

[www2.gbif.org/DataCleaning.pdf](http://www2.gbif.org/DataCleaning.pdf)

[www.webopedia.com/TERM/D/data\\_cleansing.html](http://www.webopedia.com/TERM/D/data_cleansing.html)

Books:

Data Mining by Ian H. Witten and Eibe Frank

Exploratory Data Mining and Data Quality by Dasu and Johnson (Wiley, 2004)



# References - Tools

- Stanford Wrangler  
<http://vis.stanford.edu/papers/wrangler>  
<http://vis.stanford.edu/wrangler>
- <http://openrefine.org/>
- <http://okfnlabs.org/>
- <http://schoolofdata.org/>



POLITECNICO  
MILANO 1863

SYSTEMS AND METHODS FOR BIG AND UNSTRUCTURED DATA

# Wrangling

Marco Brambilla

[marco.brambilla@polimi.it](mailto:marco.brambilla@polimi.it)

 @marcobrambi



Marco Brambilla. POLITECNICO  
MILANO 1863