# DATA CLEANING & DATA MANIPULATION

PETRA ISENBERG

VISUAL ANALYTICS

# WHAT IS "DIRTY DATA"?

BEFORE WE CAN TALK ABOUT CLEANING, WE NEED TO KNOW ABOUT TYPES OF ERROR AND WHERE THEY COME FROM

# SOURCES OF ERROR

DATA ENTRY ERRORS

MEASUREMENT ERRORS

DISTILLATION ERRORS

DATA INTEGRATION ERRORS

[HELLERSTEIN 2008]

# DATA ENTRY ERROR

## LOTS OF DATA IS ENTERED BY HAND

TYPOGRAPHIC ERRORS

MISUNDERSTANDING DATA OR CONVENTIONS

"SPURIOUS INTEGRITY"

# "SPURIOUS INTEGRITY"

ENTERING BAD DATA IN RESPONSE TO (OFTEN WELL-INTENTIONED) INTERFACE CONSTRAINTS
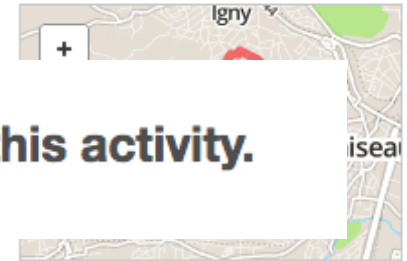
# "SPURIOUS INTEGRITY"

Step 1: Activity/Equipment Type → Step 2: Add a Map → **Step 3: Additional Details**

## Add An Activity

**Date of Activity:**

< September 2014 >

| Su | M | | | | | |
|----|---|---|---|---|---|---|
| 7 | | | | | | |
| 14 | | | | | | |
| 21 | 22 | 23 | 24 | 25 | 26 | 27 |
| 28 | **29** | 30 | | | | |

**Average Heart Rate (optional):**

[    ] bpm

**Duration:**

[ 00 ] : [ 00 ] : [ 00 ]

5.62    mi    [          ]

**Training Plan:**

None

**Oops! You forgot to enter a duration for this activity.**

### Activity Details

| | |
|---|---|
| Activity Type: | Running |
| Equipment Type: | None |
| Route: | None |
| Distance: | 5.62 mi. |
| Duration: | -:-:- |

# MEASUREMENT ERRORS

SENSOR ISSUES

MALFUNCTIONS

PLACEMENT

INTERFERENCE

MISCALIBRATION

# DISTILLATION ERRORS

SOME DATA MAY BE LOST OR COMPRESSED
BEFORE IT ENTERS
THE DATABASE

0.345413 ➡ 0.35

National Price Index ➡ NPI

1985, $2, Apples
1985, $2, Oranges ➡ 1985, $2, "Apples,Oranges,Cucumbers"
1985, $2, Cucumbers

# DATA INTEGRATION ERRORS

DATA OFTEN COMES FROM MULTIPLE SOURCES

SCHEMAS CHANGE OVER TIME

DATA IS OFTEN COERCED FROM
ONE TYPE TO ANOTHER

CAN LEAD TO DATA LOSS,
DUPLICATION, AND OTHER

# WHY IS THIS IMPORTANT?

# MOST OF THE TIME IN THE DATA ANALYSIS PROCESS IS ACTUALLY SPENT HERE!

*"I spend more than half my time integrating, cleansing, and transforming data without doing any actual analysis. Most of the time I'm lucky if I get to do any 'analysis' at all."*

**[Kandel 2012]**

# SOME DATA QUALITY ISSUES

**MISSING DATA** — MISSED MEASUREMENTS, REDACTED ITEMS, INCOMPLETE FORMS, ETC.

**ERRONEOUS VALUES** — MISSPELLINGS, OUTLIERS, "SPURIOUS INTEGRITY", ETC.

**ENTITY RESOLUTION** — DIFFERENT VALUES, ABBREVS., 2+ ENTRIES FOR THE SAME THING?

**TYPE CONVERSION** — E.G., ZIP CODE OR PLACE NAME TO LAT-LON

**DATA INTEGRATION** — MISMATCHES AND INCONSISTENCIES WHEN COMBINING DATA

# DETECTING ERRORS

LOOK FOR OUTLIERS / ANOMALIES

EXAMINE DATA TYPES

SCHEMA CHECKING

VALIDATE WITH OTHER DATA

OTHER HEURISTICS

HISTORICALLY – MORE FOCUS ON AUTOMATED APPROACHES

# DETECTION METHODS

+ CAN IDENTIFY POTENTIAL ANOMALIES

- HARD TO KNOW IF THEY'RE REALLY ANOMALOUS OR HOW TO CORRECT THEM

| Type | Issue | Detection Method(s) |
|---|---|---|
| **Missing** | Missing record | Outlier Detection \| Residuals then Moving Average w/ Hampel X84 |
| | | Frequency Outlier Detection \| Hampel X84 |
| | Missing value | Find NULL/empty values |
| **Inconsistent** | Measurement units | Clustering \| Euclidean Distance |
| | | Outlier Detection \| z-score, Hampel X84 |
| | Misspelling | Clustering \| Levenshtein Distance |
| | Ordering | Clustering \| Atomic Strings |
| | Representation | Clustering \| Structure Extraction |
| | Special characters | Clustering \| Structure Extraction |
| **Incorrect** | Erroneous entry | Outlier Detection \| z-score, Hampel X84 |
| | Extraneous data | Type Verification Function |
| | Misfielded | Type Verification Function |
| | Wrong physical data type | Type Verification Function |
| **Extreme** | Numeric outliers | Outlier Detection \| z-score, Hampel X84, Mahalanobis distance |
| | Time-series outliers | Outlier Detection \| Residuals vs. Moving Average then Hampel X84 |
| **Schema** | Primary key violation | Frequency Outlier Detection \| Unique Value Ratio |

# MISSING AND IMPOSSIBLE VALUES

1. LOOK AT EMPTY/MISSING VALUES
2. LOOK AT IMPOSSIBLE VALUES

Gender = 3

Heart Rate = 0

Unlikely Dates (e.g. "01/01/0001")

**JUST <u>SORTING</u> THE DATA CAN HELP HIGHLIGHT ISSUES LIKE THESE**

# OUTLIER DETECTION

1. EXAMINE DISTRIBUTIONS
2. MODEL DATA AND LOOK FOR RESIDUALS
3. PARTITION DATA

FOR ONE DATA DIMENSION OR MULTIPLE DIMENSIONS

# EXAMINE DISTRIBUTIONS

# DETECTING DUPLICATES

## Title
Ben-Hur
Ben Hur
BEN-HUR
Ben-Hur (1959 film)

## Name
Anand Vaskar
Anand Vaskkar
A. Vaskar
Vaskar, Anand

THESE *MICHT* ALL BE THE SAME

# SOME USEFUL DISTANCE METRICS

## LEVENSHTEIN ("STRING-EDIT") DISTANCE

How many edits do I need to change one value into another?

Ben-Hur
Ben Hur

**DISTANCE = 1**

Anand Vaskar
Anand Vaskkar

**DISTANCE = 1**

# SOME USEFUL DISTANCE METRICS

## LEVENSHTEIN ("STRING-EDIT") DISTANCE

How many edits do I need to change one value into another?

Ben-Hur
Ben-Hur (1959 film)

DISTANCE = 12

Anand Vaskar
Vaskar, Anand

DISTANCE = 12

# SOME USEFUL DISTANCE METRICS

## SOUNDEX / METAPHONE

How similar do they sound?

Ben-Hur            Anand Vaskar
Ben-Hurr           Anand Vaskkar
Been Her           Ahnund Vachkar

# SOME USEFUL DISTANCE METRICS

Strip away unimportant details.

(e.g., remove punctuation, capitals, and sort)

Anand Vaskar ➡ anand vaskar

Vaskar, Anand ➡ anand vaskar

# AND MANY MORE

## STRING/KEY COMPARISONS

## DISTANCE METRICS FOR NUMERIC DATA

e.g., HAMPEL X84 (UNIVARIATE), MAHALANOBIS (MULTIVARIATE)

**"Quantitative Data Cleaning for Large Databases"**
Hellerstein (2008)

# DECIDING HOW TO FIX PROBLEMS

YOU CAN DO ALMOST ALL OF
THIS IN **SQL** … BUT IT'S A LOT OF WORK

# DECIDING HOW TO FIX PROBLEMS

WHICH DUPLICATE TO KEEP?

OUTLIERS: KEEP, REMOVE, OR REPAIR?

BADLY-STORED DATES, ADDRESSES, OR KEYS MAY NEED TO BE PARSED MANUALLY

# DECIDING HOW TO FIX PROBLEMS

FUZZY MATCHING SYSTEMS

MACHINE LEARNING TO DETECT/RESOLVE ERRORS

USUALLY REQUIRES HUMAN JUDGMENT
(ESPECIALLY FOR NEW DATA)

# INTERACTIVE PROFILING



PROFILER [KANDEL ET AL. 2012]

# "PROFILING" DATA

<u>UNDERSTANDING</u> WHAT ASSUMPTIONS YOU CAN MAKE ABOUT DATA


<u>INTERACTIVELY</u> IDENTIFYING
DATA QUALITY ISSUES

# AN EXAMPLE

| Title | Release Date | MPAA Rating | Distributor | Rotten Tomatoes Rating | IMDB Rating |
|---|---|---|---|---|---|
| The Land Girls | Jun 12, 1998 | R | Gramercy | | 6.1 |
| First Love, Last Rites | Aug 7, 1998 | R | Strand | | 6.9 |
| I Married a Strange Person | Aug 28, 1998 | | Lionsgate | | 6.8 |
| Slam | Oct 9, 1998 | R | Trimark | 62 | 3.4 |
| Mississippi Mermaid | Jan 15, 1999 | | MGM | | |
| Following | Apr 4, 1999 | R | Zeitgeist | | 7.7 |
| Foolish | Apr 9, 1999 | R | Artisan | | 3.8 |
| Pirates | Jul 1, 1986 | R | | 25 | 5.8 |
| Duel in the Sun | Dec 31, 2046 | | | 86 | 7 |
| Tom Jones | Oct 7, 1963 | | | 81 | 7 |
| Oliver! | Dec 11, 1968 | | Sony Pictures | 84 | 7.5 |
| To Kill A Mockingbird | Dec 25, 1962 | | Universal | 97 | 8.4 |
| Tora, Tora, Tora | Sep 23, 1970 | | | | |
| Hollywood Shuffle | Mar 1, 1987 | | | 87 | 6.8 |
| Over the Hill to the Poorhouse | Sep 17, 2020 | | | | |
| Wilson | Aug 1, 2044 | | | | 7 |
| Darling Lili | Jan 1, 1970 | | | | 6.1 |
| The Ten Commandments | Oct 5, 1956 | | | 90 | 2.5 |
| 12 Angry Men | Apr 13, 1957 | | United Artists | | 8.9 |
| Twelve Monkeys | Dec 27, 1995 | R | Universal | | 8.1 |
| 1776 | Nov 9, 1972 | PG | Sony/ Columbia | 57 | 7 |

| Title | Release Date | MPAA Rating | Distributor | Rotten Tomatoes Rating | IMDB Rating |
|---|---|---|---|---|---|
| The Land Girls | Jun 12, 1998 | R | Gramercy | | 6.1 |
| First Love, Last Rites | Aug 7, 1998 | R | Strand | | 6.9 |
| I Married a Strange Person | Aug 28, 1998 | | Lionsgate | | 6.8 |
| Slam | Oct 9, 1998 | R | Trimark | 62 | 3.4 |
| Mississippi Mermaid | Jan 15, 1999 | | MGM | | |
| Following | Apr 4, 1999 | R | Zeitgeist | | 7.7 |
| Foolish | Apr 9, 1999 | R | Artisan | | 3.8 |
| Pirates | Jul 1, 1986 | R | | 25 | 5.8 |
| Duel in the Sun | Dec 31, 2046 | | | 86 | 7 |
| Tom Jones | Oct 7, 1963 | | | 81 | 7 |
| Oliver! | Dec 11, 1968 | | Sony Pictures | 84 | 7.5 |
| To Kill A Mockingbird | Dec 25, 1962 | | Universal | 97 | 8.4 |
| Tora, Tora, Tora | Sep 23, 1970 | | | | |
| Hollywood Shuffle | Mar 1, 1987 | | | 87 | 6.8 |
| Over the Hill to the Poorhouse | Sep 17, 2020 | | | | |
| Wilson | Aug 1, 2044 | | | | 7 |
| Darling Lili | Jan 1, 1970 | | | | 6.1 |
| The Ten Commandments | Oct 5, 1956 | | | 90 | 2.5 |
| 12 Angry Men | Apr 13, 1957 | | United Artists | | 8.9 |
| Twelve Monkeys | Dec 27, 1995 | R | Universal | | 8.1 |
| 1776 | Nov 9, 1972 | PG | Sony/Columbia | 57 | 7 |

| | | | |
|---|---|---|---|
| Arnolds Park | Oct 19, 2007 | PG-13 | The Movie Partners |
| Sweet Sweetback's Baad Asssss Song | Jan 1, 1971 | | |
| And Then Came Love | Jun 1, 2007 | Not Rated | Fox Meadow |
| Around the World in 80 Days | Oct 17, 1956 | PG | United Artists |
| Barbarella | Oct 10, 1968 | | Paramount Pictures |
| Barry Lyndon | 1975 | | Warner Bros. |
| Barbarians, The | March, 1987 | | |
| Babe | Aug 4, 1995 | G | Universal |
| Boynton Beach Club | Mar 24, 2006 | R | Wingate Distribution |
| Baby's Day Out | Jul 1, 1994 | PG | 20th Century |

| Bad Boys | Apr 7, 1995 | 6.6 | 53929 |
|---|---|---|---|
| Body Double | Oct 26, 1984 | 6.4 | 9738 |
| The Beast from 20,000 Fathoms | Jun 13, 1953 | | |
| Beastmaster 2: Through the Portal of Time | Aug 30, 1991 | 3.3 | 1327 |
| The Beastmaster | Aug 20, 1982 | 5.7 | 5734 |
| Ben-Hur | Dec 30, 2025 | 8.2 | 58510 |
| Ben-Hur | Nov 18, 1959 | 8.2 | 58510 |
| Benji | Nov 15, 1974 | 5.8 | 1801 |
| Before Sunrise | Jan 27, 1995 | 8 | 39705 |

# PROFILING IN OPEN REFINE

# INTERACTIVE DATA CLEANING



**Trifacta Wrangler**
https://www.trifacta.com/



**Wrangler** (Stanford HCI Group)
http://vis.stanford.edu/wrangler/



**OpenRefine** (formerly Google Refine)
http://openrefine.org/

# DATA CLEANING IN GOOGLE REFINE



Google Refine Intro Video

# REFERENCES

"Quantitative Data Cleaning for Large Databases"

Hellerstein (2008)

# CSVKIT

Search docs

Tutorial

Reference

Tips and Troubleshooting

Contributing to csvkit

Release process

License

Changelog

🔗 Edit on GitHub

# csvkit 1.0.3

## About

| build passing | FIXME Migrate to GitLab | coverage 87% | pypi v1.0.3 | license MIT |

| python 2.7 | 3.3 | 3.4 | 3.5 | 3.6 |

csvkit is a suite of command-line tools for converting to and working with CSV, the king of tabular file formats.

It is inspired by pdftk, gdal and the original csvcut tool by Joe Germuska and Aaron Bycoffe.

If you need to do more complex data analysis than csvkit can handle, use agate.

Important links:

- Repository: https://github.com/wireservice/csvkit