

Einsendeaufgabe 8 – Tool Supported Data Cleaning

1. Clean the dsm-beuth-edl-demodata-dirty.csv mini csv from the first exercise with Trifacta Wrangler (if you have no cloud access use the download version).

Not the red bars that give you instant feedback (on big datasets) where errors could be!

Create a recipe to clean the data as good as you can (it must not be a general script). Try to upload only one file (e.g. with screenshots and the end result). (10 points)

Recipe

settype col: age lockDataType: true type: Integer

filter type: custom rowType: single row: 0 >= age action: Delete

filter type: custom rowType: single row: 'old' == age action: Delete

filter type: custom rowType: single row: ('' == id) && ('' == full_name) action: Delete

set col: id value: IF('' == id, 22, \$col)

filter type: custom rowType: single row: IN(email, ['']) action: Delete

set col: gender value: IF(IN(gender, ['']), 'Divers', \$col)

filter type: custom rowType: single row: ISMISMATCHED(gender, ['Gender']) action: Delete

The screenshot displays the Trifacta Wrangler interface. On the left, a workflow diagram shows a 'Dataset' (dsm-beuth-edl-demodata-dirty.csv) connected to a 'Recipe' (dsm-beuth-edl-demodata-dirty) and then to an 'Output' (dsm-beuth-edl-demodata-dirty). The main area on the right shows the 'Details' panel for the recipe. It includes a 'Recipe' tab with a 'Steps Preview' section listing seven steps: 1. Lock age type to integer, 2. Delete rows where 0 >= age, 3. Delete rows where 'old' == age, 4. Delete rows where ('' == id) && ('' == full_name), 5. Set id to IF('' == id, 22, \$col), 6. Delete rows where IN(email, ['']), and 7. Set gender to IF(IN(gender, ['']), 'Divers', \$col). The bottom right corner shows the recipe's status: 7 steps, updated today at 12:05 PM, and created today at 11:35 AM.

| | id | email | gender | age | first_name | last_name | full_name |
|----|----|------------------------------|--------|-----|------------|-------------|--------------------|
| 1 | 1 | arfonigant@honda.gov | Female | 68 | Marisel | Famigien | Marisel Famigien |
| 2 | 2 | kpossek1@uoz.com | Male | 12 | Kenyon | Possek | Kenyon Possek |
| 3 | 3 | lmanifou1d2@pks.org | Male | 26 | Lalo | Manifou1d | Lalo Manifou1d |
| 4 | 4 | ncarow2@phoca.cz | Male | 4 | Nickola | Carow | Nickola Carow |
| 5 | 5 | ndubbin4@wikipedia.org | Male | 17 | Norman | Dubbin | Norman Dubbin |
| 6 | 6 | rcastell1@68.com | Male | 25 | Franz | Castello | Franz Castello |
| 7 | 7 | jtarney7@ft.com | Male | 77 | Jorge | Tarney | Jorge Tarney |
| 8 | 8 | eblakebrough@sohu.com | Female | 45 | Eunice | Blakebrough | Eunice Blakebrough |
| 9 | 9 | pdomotor@github.io | Male | 6 | Palu | Domotor | Palu Domotor |
| 10 | 10 | llansdowne@theguardian.com | Female | 16 | Luz | Lansdowne | Luz Lansdowne |
| 11 | 11 | skubiec@cmu.edu | Female | 91 | Modestia | Kubie | Modestia Kubie |
| 12 | 12 | obovis@bwheden.co.uk | Female | 22 | Stacey | Bovis | Stacey Bovis |
| 13 | 13 | ewacee@marriott.com | Female | 16 | Eden | Wace | Eden Wace |
| 14 | 14 | twaccer@marriott.com | Female | 16 | Eden | Wace | Eden Wace |
| 15 | 15 | taherburn@facebook.com | Male | 2 | Tobias | Sherburn | Tobias Sherburn |
| 16 | 16 | maddicot@acquirethisname.com | Male | 65 | Nathew | Addicott | Nathew Addicott |
| 17 | 17 | mshaw1@dmz.org | Male | 72 | Maurita | Shaw1 | Maurita Shaw1 |

Recipe steps:

1. Lock age type to Integer
2. Delete rows where 0 >= age
3. Delete rows where 'old' == age
4. Delete rows where (' == id) && (' == full_name)
5. Set id to IF(' == id, 22, \$col)
6. Delete rows where IN(email, [])
7. Set gender to IF(IN(gender, []), 'Divers', \$col)
8. Delete rows where ISMISMATCHED(gender, ['Gender'])

2. Load the Grid_Disruption_00_14_standardized - Grid_Disruption_00_14_standardized.csv Dataset from Kaggle: 15 YEARS OF POWER OUTAGES. Where are errors here? How would you clean this file? (5 Points)

- Die Spalte „Demand Loss (MW)“ hat sehr unterschiedliche Werte. Für nicht vorhandene Werte wird sowohl „Unknown“ als auch „N/A“ verwendet. Wenn der Wert bekannt ist wird dieser in unterschiedlicher Notation verwendet und ist daher nicht direkt vergleichbar. Die Daten müssten dazu zunächst skaliert werden
- Time of Restoration ich würde die Zeitangabe in ein 24 Stunden Format überführen um eine besser Vergleichbarkeit auf den ersten Blick zu haben
- Event Description, hier wird zum Teil der genauere Grund mit angegeben, zum Teil wird dieser nicht erwähnt. Ich würde für die Unterscheidung in Wetter bedingt sowie Menschen bedingt noch eine weitere Zeile einfügen
- Geografic Area, hier wird zum Teil keine genaue Angabe, zum Teil der Bundesstaat sowie zum Teil der Bundesstaat mit genauer Region angegeben. Ich würde die Daten dahingehend vereinheitlichen, dass in der Spalte nur keine Angabe oder Bundesstaat angegeben wird. Für die Angabe einer genauen Region würde ich eine weitere Spalte hinzufügen.