# Discoveries on the forty datasets

#### 1- 186 - Wine Quality - Business

Dataset analysed: winequality-red.csv from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/wine-quality/winequality-red.csv">https://archive.ics.uci.edu/ml/machine-learning-databases/wine-quality/winequality-red.csv</a> (dataset\_file\_url column from Fortydatasets.xlsx) has header

Numerical format: All 1599 values are numerical in the range (4.6:15.9). volatile acidity (acidity): Numerical format: All 1599 values are numerical in the range (0.12:1.58). citric acid (acid): Numerical format: All 1599 values are numerical in the range (0.0:1.0). residual sugar (sugar): Numerical format: All 1599 values are numerical in the range (0.9:15.5). chlorides (chloride): Numerical format: All 1599 values are numerical in the range (0.012:0.611). free sulfur dioxide (sulfur): Numerical format: All 1599 values are numerical in the range (1.0:72.0). total sulfur dioxide (total): Numerical format: All 1599 values are numerical in the range (6.0:289.0). Numerical >=0 format: All 1599 values are numerical and greater or equal to 0 in the range (0.99007:1.00369). pH format: All 1599 values are numerical and valid in the range [0, 14]. Actual range of values: (2.74: 4.01) sulphates (sulphate): Numerical format: All 1599 values are numerical in the range (0.33:2.0). Numerical format: All 1599 values are numerical in the range (8.4:14.9). Numerical >=0 format: All 1599 values are numerical and greater or equal to 0 in the range (3:8). Last run on: 2024-01-26 12:45:06 Alerts 9 Reproduction Alerts Dataset has 220 (13.8%) duplicate rows citric acid is highly overall correlated with fixed acidity and 2 other fields density is highly overall correlated with fixed acidity fixed acidity is highly overall correlated with citric acid and 2 other fields High correlation free sulfur dioxide is highly overall correlated with total sulfur dioxide High correlation pH is highly overall correlated with citric acid and 1 other fields total sulfur dioxide is highly overall correlated with free sulfur dioxide volatile acidity is highly overall correlated with citric acid High correlation

YData did not find anything that we search for.

citric acid has 132 (8.3%) zeros

# 2- 222 - Bank Marketing - Business

Dataset analysed: bank-additional/bank-additional-full.csv from bank-additional.zip obtained at <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00222/bank-additional.zip">https://archive.ics.uci.edu/ml/machine-learning-databases/00222/bank-additional.zip</a> (dataset\_file\_url column from Fortydatasets.xlsx) <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00222/bank-additional.zip">https://archive.ics.uci.edu/ml/machine-learning-databases/00222/bank-additional.zip</a> (dataset\_file\_url column from Fortydatasets.xlsx)

```
age:
Age format: All 41188 values are numerical and valid in the range [0, 130].
Actual range of values: (17 : 98)
```

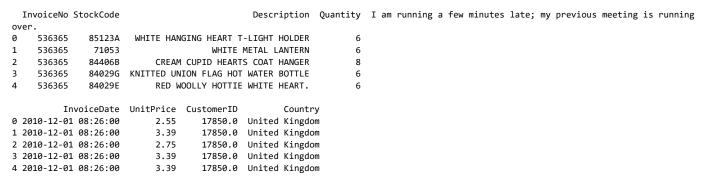
```
job:
All 41188 values are correctly categorical.
  Categorical format with 12 unique values:
     Category Frequency
                   10422
  blue-collar
                    9254
   technician
                    6743
    services
                    3969
  management
                    2924
     retired
                    1720
 entrepreneur
                    1456
                    1421
self-employed
                    1060
   housemaid
   unemployed
                    1014
      student
                     875
                     330
      unknown
marital:
All 41188 values are correctly categorical.
 Categorical format with 4 unique values:
Category Frequency
              24928
married
  single
              11568
divorced
               4612
unknown
                 80
education:
All 41188 values are correctly categorical.
  Categorical format with 8 unique values:
           Category Frequency
  university.degree
                         12168
        high.school
                          9515
           basic.9y
                          6045
professional.course
                          5243
                          4176
           basic.4v
           basic.6v
                          2292
            unknown
                          1731
         illiterate
default (categorical):
All 41188 values are correctly categorical.
  Categorical format with 3 unique values:
Category Frequency
     no
              32588
 unknown
               8597
    yes
                  3
housing (categorical):
All 41188 values are correctly categorical.
  Categorical format with 3 unique values:
Category Frequency
    yes
              21576
     no
              18622
 unknown
                990
loan (categorical):
All 41188 values are correctly categorical.
  Categorical format with 3 unique values:
Category Frequency
              33950
     no
               6248
     ves
 unknown
                990
contact (categorical):
All 41188 values are correctly categorical.
 Categorical format with 2 unique values:
 Category Frequency
 cellular
               26144
               15044
telephone
 Month format: All 41188 month values are valid.
Frequency Distribution:
    Month
           Frequency
     May
               13769
    July
                7174
   August
                6178
                5318
     Tune
 November
                4101
   April
                2632
  October 0
                 718
September
                 570
    March
                 546
```

```
December
                 182
day_of_week (day of week):
Weekday format: All 41188 weekday values are valid.
  Weekday format:
Frequency Distribution:
 Weekday Frequency
 Thursday
                8623
  Monday
                8514
Wednesday
                8134
                8090
  Tuesday
                7827
  Friday
duration:
 Numerical >=0 format: All 41188 values are correct (numeric and >=0) in the range (0:4918).
campaign (number):
 Numerical format: All 41188 values are numerical in the range (1:56).
  Numerical format: All 41188 values are numerical in the range (0:999).
previous (number):
  Numerical format: All 41188 values are numerical in the range (0:7).
poutcome (categorical):
All 41188 values are correctly categorical.
  Categorical format with 3 unique values:
   Category Frequency
nonexistent
                 35563
    failure
                  4252
    success
                  1373
emp.var.rate (rate):
  Numerical format: All 41188 values are numerical in the range (-3.4:1.4).
cons.price.idx (price):
 Numerical format: All 41188 values are numerical in the range (92.201:94.767).
cons.conf.idx (indicator):
  Numerical format: All 41188 values are numerical in the range (-50.8:-26.9).
euribor3m (indicator):
  Numerical format: All 41188 values are numerical in the range (0.634:5.045)
nr.employed (indicator):
  Numerical format: All 41188 values are numerical in the range (4963.6:5228.1).
y - has the client subscribed a term deposit? (binary):
 Binary format: All 41188 binary values are valid.
Frequency Distribution:
Value Frequency
  no
           36548
            4640
Last run on: 2024-01-25 00:02:33
Overview
              Alerts 5
                            Reproduction
  Alerts
   Dataset has 12 (< 0.1%) duplicate rows
   default is highly imbalanced (53.3%)
   loan is highly imbalanced (51.3%)
   poutcome is highly imbalanced (56.8%)
   previous has 35563 (86.3%) zeros
```

YData did not find anything that we care

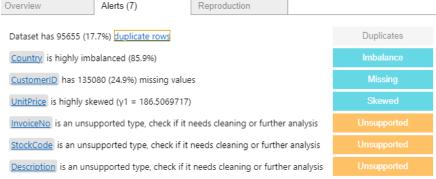
#### 3- 352 - Online Retail - Business

Dataset analysed: Online Retail.xlsx from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00352/Online%20Retail.xlsx">https://archive.ics.uci.edu/ml/machine-learning-databases/00352/Online%20Retail.xlsx</a> (dataset\_file\_url column from Fortydatasets.xlsx) <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00352/Online%20Retail.xlsx">https://archive.ics.uci.edu/ml/machine-learning-databases/00352/Online%20Retail.xlsx</a>



#### Alerts in Ydata-Profiling

1 Non-string value(s) at index(es): [(420391, 20713)]



It could not analyse the three last columns above, which are analysed below, with the Data Quality Issues found:

```
InvoiceNo (number):
  Numerical format: Error(s) found:
DQI #17 (Wrong Data Type - Consistency):
9291 Non-numeric value(s) at index(es): [(141, 'C536379'), (154, 'C536383'), (235, 'C536391'), (236, 'C536391'), (237, 'C536391'), (238, 'C536391'), (239, 'C536391'), (540141, 'C581468'), (540142, 'C581468'), (540176, 'C581470'), (540422, 'C581484'), (540448, 'C581490'), (540449, 'C581490'), (541541, 'C581499'), (541715, 'C581568'), (541716, 'C581569'), (541717, 'C581569')] (displaying only the first and last 10 items)
Range of values: (536365.0:581587.0).
StockCode (code):
    Error(s) found:
DQI #10 (Structural Conflicts - Consistency, Uniqueness):
Data seems not categorical or has too many categories (> 100). Sample values: ['85123A', 71053, '84406B', '84029G', '84029E', 22752, 21730, 22633, 22632, 84879]
 Categorical format with 4070 unique values:
Category Frequency
     85123A
       22423
     85099B
                           2159
       47566
                           1727
       20725
                           1639
       84879
                           1502
       22720
                           1477
       22197
                           1476
       21212
                           1385
       20727
                           1350
                             . . .
 DCGS0055
                                1
 DCGS0057
DCGS0066P
 DCGS0067
 DCGS0068
  DCGS0071
 DCGS0072
                                 1
 DCGS0073
                                 1
 DCGS0074
                                 1
              m
                                 1
Description:
 String format: Error(s) found:
DQI #1 (Missing Data - Completeness):
1454 Blank/Empty/Null/NaN value(s) at index(es): [(622, ''), (1970, ''), (1971, ''), (1972, ''), (1987, ''), (1988, ''), (2024, '(2025, ''), (2026, ''), (2406, ''), ('...', '...'), (524473, ''), (524475, ''), (529667, ''), (533711, ''), (533712, ''), (535322, (535326, ''), (535332, ''), (536981, ''), (538554, '')] (displaying only the first and last 10 items)
DQI #17 (Non-String Data Type - Consistency):
```

```
Quantity:
Numerical format: All 541909 values are numerical in the range (-80995:80995).

InvoiceDate (date and time):
Datetime format: All 541909 datetime values are valid in the YYYYMMDD format in the UnitPrice (price):
Numerical format: All 541909 values are numerical in the range (-11062.06:38970.0).

CustomerID (id):
ID column format: Error(s) found:
DQI #1 (Missing Data - Completeness):
135090 Blank (Empty(Mull (Malk value(s) at index(os)): [(622, 11) (1443, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11) (1444, 11)
```

CustomerID (id):

ID column format: Error(s) found:

DQI #1 (Missing Data - Completeness):

135080 Blank/Empty/Null/NaN value(s) at index(es): [(622, ''), (1443, ''), (1444, ''), (1445, ''), (1446, ''), (1447, ''), (1448, '')

(1449, ''), (1450, ''), (1451, ''), ('...', '...'), (541531, ''), (541532, ''), (541533, ''), (541534, ''), (541535, '')

(541537, ''), (541538, ''), (541539, ''), (541540, '')] (displaying only the first and last 10 items)

DQI #9 (Duplicates - Uniqueness):

541830 Duplicate value(s) at index(es): [[0, 17850], [1, 17850], [2, 17850], [3, 17850], [4, 17850], [5, 17850], [6, 17850], [7, 17850], [8, 17850], [9, 13047], ('...', '...'), [541899, 12680], [541900, 12680], [541901, 12680], [541902, 12680], [541903, 12680], [541904, 12680], [541905, 12680], [541906, 12680], [541907, 12680], [541908, 12680]] (displaying only the first and last 10 items)

DQI #19 (Uniqueness Violation - Uniqueness):

537536 Uniqueness violation(s) at index(es): [(1, 17850), (2, 17850), (3, 17850), (4, 17850), (5, 17850), (6, 17850), (7, 17850), (8, 17850), (10, 13047), (11, 13047), ('...', '...'), (541899, 12680), (541900, 12680), (541901, 12680), (541902, 12680), (541903, 12680), (541904, 12680), (541905, 12680), (541906, 12680), (541907, 12680), (541908, 12680)] (displaying only the first and last 10 items)

Alphanumeric range of values: (12346.0 : 18287.0)

CustomerID has 135080 (24.9%) missing values

Missing

This code also managed to find the 135080 Missing values, but it also found Duplicates and Uniqueness Violation errors in almost all values of this variable, while Ydata-Profiling found Duplicate rows in only 17,7% of the data:

Dataset has 95655 (17.7%) duplicate rows

Duplicates

```
Country:
```

Country format: All 541909 country values are valid.

Frequency Distribution (showing top and bottom 10 of 38 categories):

Country Frequency United Kingdom 495478 Germany 9495 8557 France EIRE 8196 Spain 2533 Netherlands 2371 2069 Belgium Switzerland Portugal 1519 Australia 1259 Malta 127 United Arab Emirates 68 European Community 61 RSA 58 Lebanon 45 Lithuania 35 32 Brazil Czech Republic 30 Bahrain Saudi Arabia

# 4- 602 - Dry Bean Dataset - Computer

Dataset analysed: DryBeanDataset/Dry\_Bean\_Dataset.xlsx from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00602/DryBeanDataset.zip">https://archive.ics.uci.edu/ml/machine-learning-databases/00602/DryBeanDataset.zip</a> (dataset\_file\_url column from Fortydatasets.xlsx) <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00602/DryBeanDataset.zip">https://archive.ics.uci.edu/ml/machine-learning-databases/00602/DryBeanDataset.zip</a> (dataset\_file\_url column from Fortydatasets.xlsx)

```
Area:
    Numerical >=0 format: All 13611 values are numerical and greater or equal to 0 in the range (20420:254616).

Perimeter:
    Numerical >=0 format: All 13611 values are numerical and greater or equal to 0 in the range (524.736:1985.37).

Major axis length (length):
    Numerical >=0 format: All 13611 values are numerical and greater or equal to 0 in the range (183.6011650038393:738.8601534818813).

Minor axis length (length):
    Numerical >=0 format: All 13611 values are numerical and greater or equal to 0 in the range (122.51265345074418:460.1984968278401).

Aspect ratio (ratio):
    Numerical format: All 13611 values are numerical in the range (1.0248675960667681:2.430306446836626).
```

```
Eccentricity:
 Numerical format: All 13611 values are numerical in the range (0.21895126335356507:0.9114229684680053).
 Numerical >=0 format: All 13611 values are numerical and greater or equal to 0 in the range (20684:263261).
Equivalent diameter (diameter):
 Numerical >=0 format: All 13611 values are numerical and greater or equal to 0 in the range (161.24376423134018:569.3743583287609).
Extent (area):
 Numerical >=0 format: All 13611 values are numerical and greater or equal to 0 in the range (0.55531471681117:0.8661946405648266).
Solidity (ratio):
 Numerical format: All 13611 values are numerical in the range (0.9192461570857022:0.9946774999456888).
Roundness (calculated):
 Numerical format: All 13611 values are numerical in the range (0.4896182562412148:0.9906853996160323).
 Numerical >=0 format: All 13611 values are numerical and greater or equal to 0 in the range (0.6405767589768725:0.9873029693778109).
ShapeFactor1 (factor):
 Numerical format: All 13611 values are numerical in the range (0.0027780126683855494:0.010451169324378654).
ShapeFactor2 (factor):
 Numerical format: All 13611 values are numerical in the range (0.0005641690180332927:0.0036649719644516834).
ShapeFactor3 (factor):
 Numerical format: All 13611 values are numerical in the range (0.41033858414131424:0.9747671533422431).
 Numerical format: All 13611 values are numerical in the range (0.9476874027098624:0.9997325300471389).
 All 13611 values are correctly categorical.
Categorical format with 7 unique values:
Category Frequency
DERMASON
               3546
   SIRA
               2636
   SEKER
               2027
  HOROZ
               1928
   CALI
               1630
BARBUNYA
               1322
  BOMBAY
                522
          Alerts 1
                     Reproduction
Overview
 Alerts
  Dataset has 68 (0.5%) duplicate rows
                                                              Duplicates
```

YData did not find anything that we care.

# 5-545 - Rice (Cammeo and Osmancik) - Computer

dataset file url was empty at Fortydatasets.xlsx

Dataset analysed: @DATA part extracted from Rice\_Cammeo\_Osmancik.arff that was obtained at rice+cammeo+and+osmancik.zip from <a href="https://archive.ics.uci.edu/dataset/545/rice+cammeo+and+osmancik">https://archive.ics.uci.edu/dataset/545/rice+cammeo+and+osmancik</a> in the new data folders from UCI

```
Numerical >=0 format: All 3810 values are correct (numeric and >=0).
Perimeter:
 Numerical >=0 format: All 3810 values are correct (numeric and >=0).
Major Axis Length (length):
 Numerical >=0 format: All 3810 values are correct (numeric and >=0).
Minor Axis Length (length):
 Numerical >=0 format: All 3810 values are correct (numeric and >=0).
Eccentricity:
 Numerical format: All 3810 values are numerical.
Convex Area (area):
 Numerical >=0 format: All 3810 values are correct (numeric and >=0).
 Numerical format: All 3810 values are numerical.
 All 3810 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
Osmancik
               2180
 Cammeo
               1630
```

## 6- 360 - Air quality - Computer

Dataset analysed: AirQualityUCI.xlsx from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00360/AirQualityUCI.zip">https://archive.ics.uci.edu/ml/machine-learning-databases/00360/AirQualityUCI.zip</a> (dataset\_file\_url column from Fortydatasets.xlsx) <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00360/AirQualityUCI.zip">https://archive.ics.uci.edu/ml/machine-learning-databases/00360/AirQualityUCI.zip</a> (dataset\_file\_url column from Fortydatasets.xlsx) <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00360/AirQualityUCI.zip">https://archive.ics.uci.edu/ml/machine-learning-databases/00360/AirQualityUCI.zip</a>

Date: Date format: All 9357 date values are valid in the YYYYMMDD format in the range 2004/03/10 to 2005/04/04. format: All 9357 time values are valid in the range 00:00:00 to 23:00:0 Time is an unsupported type, check if it needs cleaning or further analysis True hourly averaged concentration CO in mg/m^3 (averaged): Numerical format: All 9357 values are numerical in the range (-200.0:11.9) PT08.S1 (averaged): Numerical format: All 9357 values are numerical in the range (-200.0:2039.75) True hourly averaged overall Non Metanic HydroCarbons concentration in microg/m^3 (averaged): Numerical format: All 9357 values are numerical in the range (-200:1189) True hourly averaged Benzene concentration in microg/m^3 (averaged): Numerical format: All 9357 values are numerical in the range (-200.0:63.74147644829163) PT08.S2 (averaged): Numerical format: All 9357 values are numerical in the range (-200.0:2214.0) True hourly averaged NOx concentration in ppb (averaged): Numerical format: All 9357 values are numerical in the range (-200.0:1479.0) PT08.S3 (averaged): Numerical format: All 9357 values are numerical in the range (-200.0:2682.75) True hourly averaged NO2 concentration in microg/m^3 (averaged): Numerical format: All 9357 values are numerical in the range (-200.0:339.7) PT08.S4 (averaged): Numerical format: All 9357 values are numerical in the range (-200.0:2775.0) PT08.S5 (averaged): Numerical format: All 9357 values are numerical in the range (-200.0:2522.75) Temperature in °C (temperature): Numerical format: All 9357 values are numerical in the range (-200.0:44.60000038147) Relative Humidity (humidity): Numerical >=0 format: Error(s) found: DQI #15 (Domain Violation - Accuracy): 366 Negative value(s) at index(es): [(524, -200.0), (525, -200.0), (526, -200.0), (701, -200.0), (702, -200.0), (703, -200.0), (704, -200.0), (705, -200.0), (706, -200.0), (707, -200.0), ('...', '...'), (8106, -200.0), (8107, -200.0), (8108, -200.0), (8109, -200.0), (8110, -200.0), (8111, -200.0), (8112, -200.0), (8113, -200.0), (8114, -200.0), (8777, -200.0)] (displaying only the first and last 10 Range of values: (-200.0:88.72500038147) ▼ RH Distinct 4903 Minimum -200 RH Distinct (%) 52,4% Maximum 88,725 Real number (R) Missing Zeros Negative (%) Infinite (%) 0.0% 3.9% Mean 39.483611 Memory size 73.2 KiB AH Absolute Humidity (humidity): Numerical >=0 format: Error(s) found: DQI #15 (Domain Violation - Accuracy): 366 Negative value(s) at index(es): [(524, -200.0), (525, -200.0), (526, -200.0), (701, -200.0), (702, -200.0), (703, -200.0), (704, -200.0), (705, -200.0), (706, -200.0), (707, -200.0), ('...', '...'), (8106, -200.0), (8107, -200.0), (8108, -200.0), (8109, -200.0), (8110, -200.0), (8111, -200.0), (81 Range of values: (-200.0:2.2310357155831864) **▼** AH Distinct 8988 Minimum -200

96.1%

0.0%

0.0%

-6.8376037

Distinct (%)

Infinite
Infinite (%)

Mean

Real number (R)
HIGH CORRELATION

2.231035

0.09

73.2 KiB

Maximum

Memory size

Last run on: 2024-01-23 19:24:17

Ydata Profiling did find negative values in the two last columns, but it did not give any emphasis on it.

Below are lines 524 to 526. Observe that many other values have the content -200, besides the two last Humidity values:

Date	Time	CO (GT )	PT08.S1( CO)	NMHC(G T)	С6H6(G Т)	PT08.S2(NMH C)	NOx(G T)	PT08.S3(NO x)	NO2( GT)	PT08.S4(NO2)	PT08.S5(O 3)	Т	RH	АН
1/04/2004	14:00:00	1.7	-200	222	-200.0	-200	99	-200	72	-200	-200	200	200	200
1/04/2004	15:00:00	1.9	-200	197	-200.0	-200	108	-200	81	-200	-200	200	200	200
1/04/2004	16:00:00	2.3	-200	319	-200.0	-200	131	-200	93	-200	-200	200	200	200

## 7- 242 - Energy efficiency - Computer

Dataset analysed: ENB2012\_data.xlsx from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00242/ENB2012\_data.xlsx">https://archive.ics.uci.edu/ml/machine-learning-databases/00242/ENB2012\_data.xlsx</a> (dataset\_file\_url column from Fortydatasets.xlsx) <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00242/ENB2012\_data.xlsx">https://archive.ics.uci.edu/ml/machine-learning-databases/00242/ENB2012\_data.xlsx</a> (dataset\_file\_url column from Fortydatasets.xlsx) <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00242/ENB2012\_data.xlsx">https://archive.ics.uci.edu/ml/machine-learning-databases/00242/ENB2012\_data.xlsx</a> (dataset\_file\_url column from Fortydatasets.xlsx)

```
Numerical >=0 format: All 768 values are numerical and greater or equal to 0 in the range (0.62:0.98)
X2 (area):
 Numerical >=0 format: All 768 values are numerical and greater or equal to 0 in the range (514.5:808.5)
X3 (area):
 Numerical >=0 format: All 768 values are numerical and greater or equal to 0 in the range (245.0:416.5)
  Numerical >=0 format: All 768 values are numerical and greater or equal to 0 in the range (110.25:220.5)
  Numerical >=0 format: All 768 values are numerical and greater or equal to 0 in the range (3.5:7.0)
X6 (orientation):
  Numerical format: All 768 values are numerical in the range (2:5)
 Numerical >=0 format: All 768 values are numerical and greater or equal to 0 in the range (0.0:0.4)
 Numerical >=0 format: All 768 values are numerical and greater or equal to 0 in the range (0:5)
 Numerical format: All 768 values are numerical in the range (6.01:43.1)
y2 (load):
  Numerical format: All 768 values are numerical in the range (10.9:48.03)
Last run on: 2024-01-23 19:19:31
Overview
            Alerts 10
                          Reproduction
 Alerts
   x1 is highly overall correlated with x2 and 4 other fields
                                                                          High correlation
   x2 is highly overall correlated with x1 and 4 other fields
                                                                          High correlation
   x3 is highly overall correlated with x4 and 1 other fields
                                                                          High correlation
   x4 is highly overall correlated with x1 and 5 other fields
                                                                          High correlation
   x5 is highly overall correlated with x1 and 5 other fields
                                                                          High correlation
   y1 is highly overall correlated with x1 and 4 other fields
                                                                          High correlation
   y2 is highly overall correlated with x1 and 4 other fields
                                                                          High correlation
   x5 is uniformly distributed
   x6 is uniformly distributed
   x8 has 48 (6.2%) zeros
```

YData did not find anything that we care.

# 8- 267 - banknote authentication - Computer

Dataset analysed: data\_banknote\_authentication.txt from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00267/data-banknote-authentication.txt">https://archive.ics.uci.edu/ml/machine-learning-databases/00267/data-banknote-authentication.txt</a> (dataset file url column from Fortydatasets.xlsx)

```
variance of Wavelet Transformed image (variance):
   Numerical format: All 1372 values are numerical in the range (-7.0421:6.8248).
skewness of Wavelet Transformed image (skewness):
```

```
Numerical format: All 1372 values are numerical in the range (-13.7731:12.9516).
curtosis of Wavelet Transformed image (curtosis):
 Numerical format: All 1372 values are numerical in the range (-5.2861:17.9274).
entropy of image (entropy):
 Numerical format: All 1372 values are numerical in the range (-8.5482:2.4495).
class (integer):
 Numerical format: All 1372 values are numerical in the range (0:1).
Last run on: 2024-01-28 15:37:37
Overview
             Alerts 6
                            Reproduction
 Alerts
  Dataset has 11 (0.8%) duplicate rows
   class is highly overall correlated with variance of Wavelet Transformed image
   curtosis of Wavelet Transformed image is highly overall correlated with skewness of
  Wavelet Transformed image
   entropy of image is highly overall correlated with skewness of Wavelet Transformed
                                                                                       High correlation
   skewness of Wavelet Transformed image is highly overall correlated with curtosis of
  Wavelet Transformed image and 1 other fields
```

YData did not find anything that we care.

## 9- 29 - Computer Hardware – Computer

variance of Wavelet Transformed image is highly overall correlated with class

Dataset analysed: machine.data from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/cpu-performance/machine.data">https://archive.ics.uci.edu/ml/machine-learning-databases/cpu-performance/machine.data</a> (dataset\_file\_url column from Fortydatasets.xlsx) It does not have header

```
vendor name (name):
 Name format: All 209 vendor name values are valid.
Frequency Distribution:
 vendor name Frequency
        ibm
        nas
                    19
  honeywell
                    13
        nor
                    13
      sperry
                    13
     siemens
                    12
      amdahl
        cdc
   burroughs
                     8
                     7
         dg
    nixdorf
perkin-elmer
     apollo
       basf
        bti
       wang
    adviser
                     1
 four-phase
                     1
   microdata
                     1
Range of Values: (adviser to wang)
 Model Name format: All 209 Model Name values are valid.
Frequency Distribution:
Model Name Frequency
      100
1100/61-h1
                   1
   1100/81
                   1
   1100/82
                   1
   1100/83
                   1
   1100/84
                   1
   1100/93
                   1
```

1100/94

```
1636-1
   1636-10
                  1
     v8635
    v8650
    v8655
    v8665
    v8670
vax:11/730
vax:11/750
vax:11/780
   vs-100
    vs-90
                  1
Range of Values: (100 to vs-90)
MYCT (integer):
 Numerical format: All 209 values are numerical in the range (17:1500).
MMIN (minimum):
 Numerical format: All 209 values are numerical in the range (64:32000).
 Numerical format: All 209 values are numerical in the range (64:64000).
CACH (integer):
 Numerical format: All 209 values are numerical in the range (0:256).
CHMIN (minimum):
 Numerical format: All 209 values are numerical in the range (0:52).
CHMAX (maximum):
 Numerical format: All 209 values are numerical in the range (0:176).
 Numerical format: All 209 values are numerical in the range (6:1150).
 Numerical format: All 209 values are numerical in the range (15:1238).
Last run on: 2024-01-28 19:22:45
Overview
             Alerts 13
                             Reproduction
```

#### Alerts

CACH is highly overall correlated with CHMAX and 5 other fields	High correlation
CHMAX is highly overall correlated with CACH and 6 other fields	High correlation
CHMIN is highly overall correlated with CHMAX and 4 other fields	High correlation
ERP is highly overall correlated with CACH and 6 other fields	High correlation
MMAX is highly overall correlated with CACH and 5 other fields	High correlation
MMIN is highly overall correlated with CACH and 5 other fields	High correlation
MYCT is highly overall correlated with CACH and 6 other fields	High correlation
PRP is highly overall correlated with CACH and 6 other fields	High correlation
vendor name is highly overall correlated with CHMAX	High correlation
Model Name has unique values	Unique
CACH has 69 (33.0%) zeros	Zeros
CHMIN has 5 (2.4%) zeros	Zeros
CHMAX has 5 (2.4%) zeros	Zeros

YData found that Model Name has unique values. We have not checked that.

# 10- 294 - Combined Cycle Power Plant - Computer

Dataset analysed: CCPP/Folds5x2\_pp.xlsx from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00294/CCPP.zip">https://archive.ics.uci.edu/ml/machine-learning-databases/00294/CCPP.zip</a> (dataset\_file\_url column from Fortydatasets.xlsx)

This was the original description of the Attributes Information:

<sup>&</sup>quot;Features consist of hourly average ambient variables

- 1. Temperature (T) in the range 1.81°C and 37.11°C,
- 2. Ambient Pressure (AP) in the range 992.89-1033.30 milibar,
- 3. Relative Humidity (RH) in the range 25.56% to 100.16%
- 4. Exhaust Vacuum (V) in teh range 25.36-81.56 cm Hg
- 5. Net hourly electrical energy output (EP) 420.26-495.76 MW

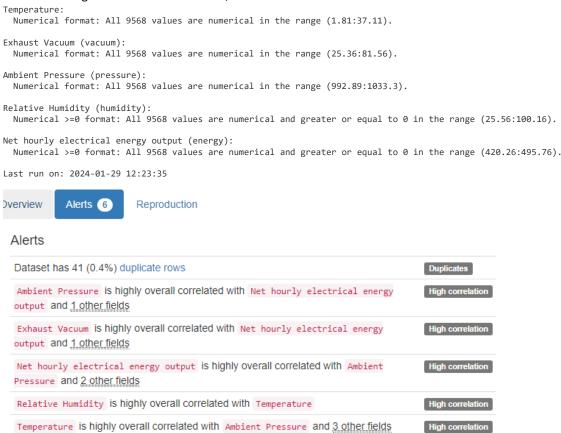
The averages are taken from various sensors located around the plant that record the ambient variables every second. The variables are given without normalization."

And these were the first 5 items:

```
٧
                    AP
                                  PΕ
     AT
0
  14.96 41.76 1024.07 73.17
                              463.26
  25.18 62.96 1020.04 59.08 444.37
1
        39.40 1012.16 92.14
                              488.56
   5.11
  20.86 57.32 1010.24
                              446.48
3
                        76.64
  10.82 37.50
               1009.23 96.62 473.90
```

The order of the dataset are not the ones in the description. Only the first and last are correct. The others are changed. The second column in the data is V, which is the fourth column in the description. The third column in the data is AP, which is the second column in the description, and the fourth column in the dataset is RH, which is the third in the description.

#### After correcting the order of the columns, the result was:



YData did not find anything that we care.

## 11- 229 - Skin Segmentation – Computer

Dataset analysed: Skin\_NonSkin.txt from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00229/Skin">https://archive.ics.uci.edu/ml/machine-learning-databases/00229/Skin</a> NonSkin.txt (dataset file url column from Fortydatasets.xlsx)

```
x1 (b):
   Numerical format: All 245057 values are numerical in the range (0:255).
x2 (g):
   Numerical format: All 245057 values are numerical in the range (0:255).
x3 (r):
   Numerical format: All 245057 values are numerical in the range (0:255).
y (class):
   All 245057 values are correctly categorical.
```

```
Categorical format with 2 unique values:
 Category Frequency
                 194198
                  50859
Last run on: 2024-01-29 13:50:39
             Alerts 8
                           Reproduction
Overview
 Alerts
  Dataset has 20347 (8.3%) duplicate rows
   x1 is highly overall correlated with x2
  x2 is highly overall correlated with x1 and 1 other fields
   x3 is highly overall correlated with x2 and 1 other fields
  y is highly overall correlated with x3
                                                                              High correlation
   x1 has 3237 (1.3%) zeros
  x2 has 2968 (1.2%) zeros
   x3 has 6539 (2.7%) zeros
```

YData did not find anything that we care.

## 12- 246 - 3D Road Network (North Jutland, Denmark) — Computer

Dataset analysed: 3D\_spatial\_network.txt from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00246/3D\_spatial\_network.txt">https://archive.ics.uci.edu/ml/machine-learning-databases/00246/3D\_spatial\_network.txt</a> (dataset\_file\_url column from Fortydatasets.xlsx)

```
OSM_ID (id):
    ID column format: Error(s) found:
DQI #9 (Duplicates - Uniqueness):
421305 Duplicates - Onlydeness).
421305 Duplicate value(s) at index(es): [[0, 144552912], [1, 144552912], [2, 144552912], [3, 144552912], [4, 144552912], [5, 144552912], [6, 144552912], [7, 144552912], [8, 144552912], [9, 144552912], ('...', '...'), [434864, 93323205], [434865, 93323205], [434866, 93323205], [434867, 93323205], [434868, 93323205], [434869, 93323205], [434870, 93323205], [434871, 93323205], [434872, 93323209], [434873, 93323209]] (displaying only the first and last 10 items)
DOI #17 (Wrong Data Type - Consistency):
204080 Inconsistent length in alphanumeric value(s) at index(es): [(19, 42991631), (20, 42991631), (21, 42991631), (22, 42991631), (23, 42991631), (24, 42991631), (25, 42991631), (26, 42991631), (27, 42991632), (28, 42991632), ('...', '...'), (434864, 93323205), (434865, 93323205), (434866, 93323205), (434867, 93323205), (434868, 93323205), (434869, 93323205), (434870, 93323205), (434871, 93323206), (434872, 93323209), (434873, 93323209)] (displaying only the first and last 10 items)
DQI #19 (Uniqueness Violation - Uniqueness):
377545 Uniqueness violation(s) at index(es): [(1, 144552912), (2, 144552912), (3, 144552912), (4, 144552912), (5, 144552912), (6, 144552912), (7, 144552912), (8, 144552912), (9, 144552912), (10, 144552912), ('...', '...'), (434863, 93323205), (434864, 93323205), (434865, 93323205), (434866, 93323205), (434867, 93323205), (434868, 93323205), (434873, 93323209)] (displaying only the first and last 10 items)
Alphanumeric range of values: (100009227 : 99935416)
YDATA found nothing!
     tual range of values: (8.1461259 : 11.1993265)
LATITUDE:
     Latitude format: All 434874 values are numerica
tual range of values: (56.5824856 : 57.750511)
ALTITUDE:
   Numerical format: All 434874 values are numerical in the range (-8.60818370517905:134.441946906076).
Last run on: 2024-01-29 14:03:03
```

# 13- 374 - Appliances energy prediction – Computer

Dataset analysed: energydata\_complete.csv from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00374/energydata">https://archive.ics.uci.edu/ml/machine-learning-databases/00374/energydata</a> complete.csv (dataset\_file\_url column from Fortydatasets.xlsx) has header

```
date time year-month-day hour (date time):
    Datetime format: All 19735 datetime values are valid in the YYYYMMDD format in the range 2016/01/11 17:00 to 2016/05/27 18:00.

Appliances, energy use in Wh (energy):
    Numerical >=0 format: All 19735 values are numerical and greater or equal to 0 in the range (10:1080).

lights, energy use of light fixtures in the house in Wh (energy):
    Numerical >=0 format: All 19735 values are numerical and greater or equal to 0 in the range (0:70).

T1, Temperature in kitchen area, in Celsius (temperature):
```

```
RH 1, Humidity in kitchen area, in % (humidity):
  Numerical >=0 format: All 19735 values are numerical and greater or equal to 0 in the range (27.023333333333336).
T2, Temperature in living room area, in Celsius (temperature):
  Numerical format: All 19735 values are numerical in the range (16.1:29.856666666667).
RH_2, Humidity in living room area, in % (humidity):
  Numerical >=0 format: All 19735 values are numerical and greater or equal to 0 in the range (20.463333333333296:56.0266666666667).
T3, Temperature in laundry room area (temperature):
  Numerical format: All 19735 values are numerical in the range (17.2:29.236).
RH_3, Humidity in laundry room area, in % (humidity):
 Numerical >=0 format: All 19735 values are numerical and greater or equal to 0 in the range (28.766666666667:50.1633333333333).
T4, Temperature in office room, in Celsius (temperature):
  Numerical format: All 19735 values are numerical in the range (15.1:26.2).
RH_4, Humidity in office room, in % (humidity):
  Numerical >=0 format: All 19735 values are numerical and greater or equal to 0 in the range (27.66:51.09).
T5, Temperature in bathroom, in Celsius (temperature):
  Numerical format: All 19735 values are numerical in the range (15.33:25.795).
RH_5, Humidity in bathroom, in % (humidity):

Numerical >=0 format: All 19735 values are numerical and greater or equal to 0 in the range (29.815:96.3216666666666).
T6, Temperature outside the building (temperature):
  Numerical format: All 19735 values are numerical in the range (-6.065:28.29).
RH_6, Humidity outside the building (humidity):
  Numerical >=0 format: All 19735 values are numerical and greater or equal to 0 in the range (1.0:99.9).
T7, Temperature in ironing room , in Celsius (temperature):
  Numerical format: All 19735 values are numerical in the range (15.39:26.0).
RH 7. Humidity in ironing room, in % (humidity):
 Numerical >=0 format: All 19735 values are numerical and greater or equal to 0 in the range (23.2:51.4).
T8, Temperature in teenager room 2, in Celsius (temperature):
 Numerical format: All 19735 values are numerical in the range (16.3066666666667:27.23).
RH_8, Humidity in teenager room 2, in % (humidity):
 Numerical >=0 format: All 19735 values are numerical and greater or equal to 0 in the range (29.6:58.78).
T9, Temperature in parents room, in Celsius (temperature):
  Numerical format: All 19735 values are numerical in the range (14.89:24.5).
RH_9, Humidity in parents room, in % (humidity):
  Numerical >=0 format: All 19735 values are numerical and greater or equal to 0 in the range (29.1666666666667:53.3266666666667).
To, Temperature outside (temperature):
  Numerical format: All 19735 values are numerical in the range (-5.0:26.1).
 Numerical format: All 19735 values are numerical in the range (729.299999999998:772.29999999999).
RH_out, Humidity outside (humidity):
  Numerical >=0 format: All 19735 values are numerical and greater or equal to 0 in the range (24.0:100.0).
Wind speed (speed):
 Numerical format: All 19735 values are numerical in the range (0.0:14.0).
Visibility:
 Numerical format: All 19735 values are numerical in the range (1.0:66.0).
Tdewpoint (dewpoint):
  Numerical format: All 19735 values are numerical in the range (-6.6:15.5).
rv1, Random variable 1, nondimensional (random):
  Numerical format: All 19735 values are numerical in the range (0.0053216819651424:49.9965296825394).
rv2, Random variable 2, nondimensional (random):

Numerical format: All 19735 values are numerical in the range (0.0053216819651424:49.9965296825394).
Last run on: 2024-01-29 14:41:33
           Alerts 27
                        Reproduction
Overview
  Alerts
  RH_1, Humidity in kitchen area, in % is highly overall correlated with RH_2, Humidity in living room area, in % and 6 other fields
```

RH 2, Humidity in living room area, in % is highly overall correlated with RH 1, Humidity in kitchen area, in % and 6 other fields

Numerical format: All 19735 values are numerical in the range (16.79:26.26).

rv2, Random variable 2, nondimensional is highly overall correlated with rv1, Random variable 1, nondimensional	High correlation
date time year-month-day hour has unique values	Unique
rv1, Random variable 1, nondimensional has unique values	Unique
rv2, Random variable 2, nondimensional has unique values	Unique
lights, energy use of light fixtures in the house in Wh has 15252 (77.3%) zeros	Zeros

YData found 3 columns that are Unique. We have not checked that.

## 14- 248 - Buzz in social media - Computer

Dataset analysed: ./regression/TomsHardware/TomsHardware.data from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00248/regression.tar.gz">https://archive.ics.uci.edu/ml/machine-learning-databases/00248/regression.tar.gz</a> (dataset\_file\_url column from Fortydatasets.xlsx) does not have header

```
NCD_0 (number):
 Numerical format: All 28179 values are numerical in the range (0:182).
NCD 1 (number):
 Numerical format: All 28179 values are numerical in the range (0:118).
NCD 2 (number):
 Numerical format: All 28179 values are numerical in the range (0:118).
 Numerical format: All 28179 values are numerical in the range (0:118).
 Numerical format: All 28179 values are numerical in the range (0:118).
NCD 5 (number):
 Numerical format: All 28179 values are numerical in the range (0:118).
NCD 6 (number):
 Numerical format: All 28179 values are numerical in the range (0:154).
NCD 7 (number):
 Numerical format: All 28179 values are numerical in the range (0:88).
 Numerical format: All 28179 values are numerical in the range (0.0:1.0).
 Numerical format: All 28179 values are numerical in the range (0.0:1.0).
BL 2 (burstiness):
 Numerical format: All 28179 values are numerical in the range (0.0:1.0).
BL_3 (burstiness):
 Numerical format: All 28179 values are numerical in the range (0.0:1.0).
 Numerical format: All 28179 values are numerical in the range (0.0:1.0).
BL 5 (burstiness):
 Numerical format: All 28179 values are numerical in the range (0.0:1.0).
BL 6 (burstiness):
 Numerical format: All 28179 values are numerical in the range (0.0:1.0).
BL 7 (burstiness):
 Numerical format: All 28179 values are numerical in the range (0.0:1.0).
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0:217).
NAD_1 (length):
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0:210).
NAD 2 (length):
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0:210).
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0:210).
NAD_4 (length):
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0:210).
NAD_5 (length):
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0:194).
NAD 6 (length):
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0:170).
NAD 7 (length):
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0:185).
AI_0 (increase):
```

```
Numerical format: All 28179 values are numerical in the range (0:158).
AI 1 (increase):
 Numerical format: All 28179 values are numerical in the range (0:146).
AI_2 (increase):
  Numerical format: All 28179 values are numerical in the range (0:149).
AI_3 (increase):
 Numerical format: All 28179 values are numerical in the range (0:161).
AI 4 (increase):
 Numerical format: All 28179 values are numerical in the range (0:169).
AI 5 (increase):
 Numerical format: All 28179 values are numerical in the range (0:149).
AI_6 (increase):
 Numerical format: All 28179 values are numerical in the range (0:156).
AI_7 (increase):
 Numerical format: All 28179 values are numerical in the range (0:160).
NAC 0 (number):
 Numerical format: All 28179 values are numerical in the range (0:1734).
NAC 1 (number):
 Numerical format: All 28179 values are numerical in the range (0:1966).
NAC 2 (number):
 Numerical format: All 28179 values are numerical in the range (0:1734).
 Numerical format: All 28179 values are numerical in the range (0:1734).
NAC_4 (number):
 Numerical format: All 28179 values are numerical in the range (0:1734).
NAC 5 (number):
 Numerical format: All 28179 values are numerical in the range (0:1657).
NAC 6 (number):
 Numerical format: All 28179 values are numerical in the range (0:1403).
 Numerical format: All 28179 values are numerical in the range (0:1707).
 Numerical format: All 28179 values are numerical in the range (0:235271).
ND 1 (number):
 Numerical format: All 28179 values are numerical in the range (0:197284).
ND 2 (number):
 Numerical format: All 28179 values are numerical in the range (0:225099).
ND 3 (number):
 Numerical format: All 28179 values are numerical in the range (0:207633).
ND 4 (number):
 Numerical format: All 28179 values are numerical in the range (0:272256).
ND 5 (number):
 Numerical format: All 28179 values are numerical in the range (0:272256).
ND 6 (number):
 Numerical format: All 28179 values are numerical in the range (0:330561).
ND_7 (number):
 Numerical format: All 28179 values are numerical in the range (0:272256).
CS_0 (contribution):
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0.0:1.0).
CS 1 (contribution):
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0.0:1.0).
CS 2 (contribution):
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0.0:1.0).
CS_3 (contribution):
 \overline{\text{Numerical}} >= 0 format: All 28179 values are numerical and greater or equal to 0 in the range (0.0:1.0).
CS_4 (contribution):
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0.0:1.0).
CS 5 (contribution):
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0.0:1.0).
CS 6 (contribution):
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0.0:1.0).
CS_7 (contribution):
```

```
Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0.0:1.0).
AT 0 (interaction):
 Numerical format: All 28179 values are numerical in the range (0.0:105.0).
AT_1 (interaction):
  Numerical format: All 28179 values are numerical in the range (0.0:98.0).
AT 2 (interaction):
 Numerical format: All 28179 values are numerical in the range (0.0:104.0).
AT 3 (interaction):
 Numerical format: All 28179 values are numerical in the range (0.0:106.0).
AT 4 (interaction):
 Numerical format: All 28179 values are numerical in the range (0.0:107.0).
AT_5 (interaction):
 Numerical format: All 28179 values are numerical in the range (0.0:107.0).
AT_6 (interaction):
 Numerical format: All 28179 values are numerical in the range (0.0:104.0).
AT 7 (interaction):
 Numerical format: All 28179 values are numerical in the range (0.0:106.0).
NA 0 (number):
 Numerical format: All 28179 values are numerical in the range (0:313).
 Numerical format: All 28179 values are numerical in the range (0:313).
 Numerical format: All 28179 values are numerical in the range (0:313).
NA 3 (number):
 Numerical format: All 28179 values are numerical in the range (0:313).
NA 4 (number):
 Numerical format: All 28179 values are numerical in the range (0:313).
NA 5 (number):
 Numerical format: All 28179 values are numerical in the range (0:322).
 Numerical format: All 28179 values are numerical in the range (0:264).
NA 7 (number):
 Numerical format: All 28179 values are numerical in the range (0:309).
ADL 0 (length):
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0.0:150.0).
ADL_1 (length):
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0.0:150.0).
ADL 2 (length):
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0.0:150.0).
ADL 3 (length):
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0.0:150.0).
ADL 4 (length):
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0.0:150.0).
ADL 5 (length):
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0.0:150.0).
ADL_6 (length):
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0.0:150.0).
ADL_7 (length):
 Numerical >=0 format: All 28179 values are numerical and greater or equal to 0 in the range (0.0:150.0).
-- Attention Level (measured with number of authors) (AS(NA))
      (columns [80,87])
          +----+
```

feature					
AS(NA)_0	0	0.120	0.002	0.008	
AS(NA)_1	0	0.129	0.002	0.008	
AS(NA)_2	0	0.129	0.002	0.008	
AS(NA)_3	0	0.159	0.002	0.008	
AS(NA)_4	0	0.166	0.002	0.008	
AS(NA)_5	0	0.166	0.002	0.008	
AS(NA)_6	:	0.155			:

```
AS_NA)_0 (number):
  Numerical format: All 28179 values are numerical in the range (0.0:0.120455).
AS_NA)_1 (number):
  Numerical format: All 28179 values are numerical in the range (0.0:0.12885).
AS NA) 2 (number):
 Numerical format: All 28179 values are numerical in the range (0.0:0.12885).
AS NA) 3 (number):
  Numerical format: All 28179 values are numerical in the range (0.0:0.159019).
AS_NA)_4 (number):
  Numerical format: All 28179 values are numerical in the range (0.0:0.16569).
AS_NA)_5 (number):
  Numerical format: All 28179 values are numerical in the range (0.0:0.16569).
AS_NA)_6 (number):
  Numerical format: All 28179 values are numerical in the range (0.0:0.155071).
AS NA) 7 (number):
 Numerical format: All 28179 values are numerical in the range (0.0:0.147081).
 -- Attention Level (measured with number of contributions) (AS(NAC))
      (columns [88,95])
          | feature | min | max
                                    mean
           AS(NAC)_0 | 0
                            0.107 | 0.002 | 0.006
          | AS(NAC)_1 | 0
                            0.130 | 0.002 | 0.006
          | AS(NAC)_2 | 0
                            | 0.136 | 0.002 | 0.006
          | AS(NAC) 3 | 0
                            | 0.153 | 0.002 | 0.006 |
           AS(NAC)_4 | 0
                            | 0.153 | 0.002 | 0.006 |
            AS(NAC)_5 | 0
                            | 0.147 | 0.002 | 0.006 |
            AS(NAC)_6 | 0
                            | 0.179 | 0.002 | 0.006 |
           AS(NAC)_7 | 0
                            0.188 | 0.002 | 0.006
AS_NAC)_0 (number):
  Numerical format: All 28179 values are numerical in the range (0.0:0.107367).
AS_NAC)_1 (number):
  Numerical format: All 28179 values are numerical in the range (0.0:0.130337).
AS NAC) 2 (number):
  Numerical format: All 28179 values are numerical in the range (0.0:0.135633).
AS_NAC)_3 (number):
  Numerical format: All 28179 values are numerical in the range (0.0:0.153209).
AS NAC) 4 (number):
 Numerical format: All 28179 values are numerical in the range (0.0:0.153209).
AS NAC) 5 (number):
 Numerical format: All 28179 values are numerical in the range (0.0:0.147003).
AS_NAC)_6 (number):
 Numerical format: All 28179 values are numerical in the range (0.0:0.179334).
AS NAC) 7 (number):
  Numerical format: All 28179 values are numerical in the range (0.0:0.187696).
Feature to predict (feature):
  Error(s) found:
DQI #10 (Structural Conflicts - Consistency, Uniqueness):
Data seems not categorical or has too many categories (> 100). Sample values: [4.5, 3.5, 2.0, 1.5, 5.0, 2.5, 1.0, 0.5, 3.0, 6.0] Categorical format with 8895 unique values:
 Category Frequency
     2.0
      1.5
                182
      0.5
                175
      2.5
                169
      1.0
                157
      4.0
                143
      3.0
                142
      3.5
                134
                128
      6.0
      4.5
                126
 206317.5
```

| AS(NA)\_7 | 0 | 0.147 | 0.002 | 0.007 |

211342.0 1 211967.0 1 214715.0 1 215597.5 1 219497.5 1 223389.5 1 237281.0 1 242383.0 1 265916.5 1

Last run on: 2024-01-30 14:49:19

Feature to predict
Real number (R)

#### HIGH CORRELATION

Distinct	8895	Minimum	0
Distinct (%)	31.6%	Maximum	265916.5
Missing	0	Zeros	125
Missing (%)	0.0%	Zeros (%)	0.4%
Infinite	0	Negative	0
Infinite (%)	0.0%	Negative (%)	0.0%
Mean	3486.441	Memory size	220.3 KiB

YData confirmed in its analysis that the Feature to predict column had really 8895 distinct values, which seems to be a major problem in such a column. Our code found that to be a problem for a Categorical column should not have over 100 distinct values. Besides that YData found many High Correlation and Zeros.

### 15- 343 - Occupancy Detection – Computer

Dataset analysed: datatest.txt from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00357/occupancy\_data.zip">https://archive.ics.uci.edu/ml/machine-learning-databases/00357/occupancy\_data.zip</a> (dataset file url column from Fortydatasets.xlsx)

```
date time year-month-day hour (date time):
  Datetime format: All 2665 datetime values are valid in the YYYYMMDD format in the range 2015/02/02 14:19 to 2015/02/04 10:43.
Temperature, in Celsius (temperature):
  Numerical format: All 2665 values are numerical in the range (20.2:24.408333333333).
Relative Humidity, % (humidity):
  Numerical >=0 format: All 2665 values are numerical and greater or equal to 0 in the range (22.1:31.4725).
Light, in Lux (light):
  Numerical >=0 format: All 2665 values are numerical and greater or equal to 0 in the range (0.0:1697.25).
CO2, in ppm (co2):
  Numerical >=0 format: All 2665 values are numerical and greater or equal to 0 in the range (427.5:1402.25).
Humidity Ratio, Derived quantity from temperature and relative humidity, in kgwater-vapor/kg-air (temperature):
 Numerical format: All 2665 values are numerical in the range (0.0033033144722347:0.0053777588397133).
Occupancy, 0 or 1, 0 for not occupied, 1 for occupied status (status):
  All 2665 values are correctly categorical.
Categorical format with 2 unique values:
 Category Frequency
                 1693
                  972
Last run on: 2024-01-30 18:56:59
          Alerts 8
 Alerts
  CO2, in ppm is highly overall correlated with Humidity Ratio, Derived quantity from High correlation
  temperature and relative humidity, in kgwater-vapor and 4 other fields
  Humidity Ratio, Derived quantity from temperature and relative humidity, in
                                                                         High correlation
  kgwater-vapor is highly overall correlated with co2, in ppm and 4 other fields
  Light, in Lux is highly overall correlated with co2, in ppm and 4 other fields
                                                                         High correlation
  Occupancy, 0 or 1, 0 for not occupied, 1 for occupied status is highly overall
  correlated with co2, in ppm and 4 other fields
  Relative Humidity, % is highly overall correlated with CO2, in ppm and 4 other fields
  Temperature, in Celsius is highly overall correlated with CO2, in ppm and 4 other
  date time year-month-day hour has unique values
                                                                         Unique
```

YData found that the first column is Unique. We have not checked that.

Light, in Lux has 1615 (60.6%) zeros

## 16- 128 - KDD Cup 1999 Data - Computer

https://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html, instead of the original dataset file url column from

Dataset analysed: corrected (without file ending), renamed to corrected.txt, from

http://kdd.ics.uci.edu/databases/kddcup99/corrected.gz obtained at

Fortydatasets.xlsx which was at https://archive.ics.uci.edu/ml/machine-learning-databases/kddcup99-mld/corrected.gz but it is not valid anymore. It does not have heade duration: Numerical >=0 format: All 311029 values are numerical and greater or equal to 0 in the range (0:57715). protocol type (type): All 311029 values are correctly categorical. Categorical format with 3 unique values: Category Frequency icmp tcp 119357 udp 26703 service: All 311029 values are correctly categorical. Categorical format with 65 unique values: Category Frequency 164352 ecr i private http smtp 8268 3972 pop\_3 domain\_u 3160 ftp\_data 2223 other 2185 telnet 2077 ftp 837 csnet ns 50 50 klogin supdup 50 hostnames 48 pm\_dump 16 X11 15 tim\_i icmp 2 tftp\_u 1 All 311029 values are correctly categorical. Categorical format with 11 unique values: Category Frequency S0 18012 RSTO 1393 **RSTR** 872 S3 289 SH 84 S1 27 S2 22 ОТН 4 RST0S0 Numerical >=0 format: All 311029 values are numerical and greater or equal to 0 in the range (0:62825648). Numerical >=0 format: All 311029 values are numerical and greater or equal to 0 in the range (0:5203179). land (symbolic): All 311029 values are correctly categorical. Categorical format with 2 unique values: Category Frequency 311020 1 wrong\_fragment (continuous): Numerical >=0 format: All 311029 values are numerical and greater or equal to 0 in the range (0:3). urgent (continuous): Numerical >=0 format: All 311029 values are numerical and greater or equal to 0 in the range (0:3). hot (continuous): Numerical >=0 format: All 311029 values are numerical and greater or equal to 0 in the range (0:101). num failed logins (num):

Numerical format: All 311029 values are numerical in the range (0:4).

```
logged in (symbolic):
 All 311029 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
              257384
       1
               53645
num compromised (num):
 Numerical format: All 311029 values are numerical in the range (0:796).
root shell (continuous):
 Numerical >=0 format: All 311029 values are numerical and greater or equal to 0 in the range (0:1).
su attempted (attempted):
 Numerical >=0 format: All 311029 values are numerical and greater or equal to 0 in the range (0:2).
 Numerical format: All 311029 values are numerical in the range (0:878).
num_file_creations (num):
 Numerical format: All 311029 values are numerical in the range (0:100).
num shells (num):
 Numerical format: All 311029 values are numerical in the range (0:5).
num access files (num):
 Numerical format: All 311029 values are numerical in the range (0:4).
num_outbound_cmds (num):
 Numerical format: All 311029 values are numerical in the range (0:0).
is_host_login (symbolic):
 All 311029 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
              311017
       1
is_guest_login (symbolic):
 All 311029 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
        a
              310275
       1
                754
count:
 Numerical >=0 format: All 311029 values are numerical and greater or equal to 0 in the range (0:511).
srv count (count):
 Numerical >=0 format: All 311029 values are numerical and greater or equal to 0 in the range (0:511).
 Numerical format: All 311029 values are numerical in the range (0.0:1.0).
 Numerical format: All 311029 values are numerical in the range (0.0:1.0).
rerror rate (rate):
 Numerical format: All 311029 values are numerical in the range (0.0:1.0).
srv rerror rate (rate):
 Numerical format: All 311029 values are numerical in the range (0.0:1.0).
same_srv_rate (rate):
 Numerical format: All 311029 values are numerical in the range (0.0:1.0).
diff srv rate (rate):
 Numerical format: All 311029 values are numerical in the range (0.0:1.0).
srv diff host rate (rate):
 Numerical format: All 311029 values are numerical in the range (0.0:1.0).
dst host count (count):
 Numerical >=0 format: All 311029 values are numerical and greater or equal to 0 in the range (0:255).
dst host srv count (count):
 Numerical >=0 format: All 311029 values are numerical and greater or equal to 0 in the range (0:255).
dst_host_same_srv_rate (rate):
 Numerical format: All 311029 values are numerical in the range (0.0:1.0).
dst_host_diff_srv_rate (rate):
 Numerical format: All 311029 values are numerical in the range (0.0:1.0).
dst host same src port rate (rate):
 Numerical format: All 311029 values are numerical in the range (0.0:1.0).
dst host srv diff host rate (rate):
 Numerical format: All 311029 values are numerical in the range (0.0:1.0).
```

```
dst host serror rate (rate):
 Numerical format: All 311029 values are numerical in the range (0.0:1.0).
dst_host_srv_serror_rate (rate):
 Numerical format: All 311029 values are numerical in the range (0.0:1.0).
dst_host_rerror_rate (rate):
 Numerical format: All 311029 values are numerical in the range (0.0:1.0).
dst_host_srv_rerror_rate (rate):
 Numerical format: All 311029 values are numerical in the range (0.0:1.0).
class:
 All 311029 values are correctly categorical.
Categorical format with 38 unique values:
      Category Frequency
       smurf.
                   60593
       normal.
                   58001
     neptune.
snmpgetattack.
                    7741
     mailbomb.
                    5000
 guess_passwd.
                    4367
                    2406
    snmpguess.
                    1633
       satan.
 warezmaster.
                    1602
                    1098
        back.
       xlock.
                       4
       xsnoop.
    ftp_write.
                       3
   loadmodule.
        perl.
                       2
         phf.
    sqlattack.
    udpstorm.
        worm.
                       2
        imap.
Last run on: 2024-02-21 21:03:45
Dataset
                                     Reproduction
Overview
                 Alerts 43
```

#### Alerts



It shows column num\_outbound\_cmds with constant value "". We have not checked that. For us it is a numerical with values in the range (0:0).

```
      dst_host_srv_diff_host_rate
      is highly skewed (γ1 = 21.57807214)

      duration
      has 298054 (95.8%) zeros

      Zeros
```

## 17- 303 - Perfume Data - Computer

Dataset analysed: perfume\_data.xlsx from <a href="https://archive.ics.uci.edu/static/public/303/perfume+data.zip">https://archive.ics.uci.edu/static/public/303/perfume+data.zip</a> instead of the original dataset\_file\_url column from Fortydatasets.xlsx which was at <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00303/perfume\_data.xlsx">https://archive.ics.uci.edu/ml/machine-learning-databases/00303/perfume\_data.xlsx</a> but it is not valid anymore.

```
Perfume_name (name):
 Name format: Error(s) found:
DQI #5 (Extraneous Data - Consistency, Uniqueness):
1 Extraneous data value(s) at index(es): [(16, 'constrected2')]
Frequency Distribution:
    Perfume_name Frequency
        RoseMusk
      TeaTreeOil
                          1
          ajayeb
                          1
           ajmal
                          1
          amreai
                          1
            aood
```

```
asgar ali
        bukhoor
       burberrry
carolina_herrera
    constrected
    constrected2
      dehenalaod
          junaid
          kausar
                          1
 oudh_ma'alattar
                          1
      raspberry
                          1
            rose
                          1
      solidmusk
                          1
     strawberry
                          1
Range of Values: (RoseMusk to strawberry)
Take_1 (take):
 Numerical format: All 20 values are numerical in the range (46014:85056).
 Numerical format: All 20 values are numerical in the range (46014:85056).
Take 3 (take):
 Numerical format: All 20 values are numerical in the range (46014:85056).
Take 4 (take):
 Numerical format: All 20 values are numerical in the range (46014:85056).
 Numerical format: All 20 values are numerical in the range (46014:85056).
 Numerical format: All 20 values are numerical in the range (46015:85056).
Take 7 (take):
 Numerical format: All 20 values are numerical in the range (46015:85056).
Take 8 (take):
 Numerical format: All 20 values are numerical in the range (46015:85056).
Take_9 (take):
 Numerical format: All 20 values are numerical in the range (46015:85056).
 Numerical format: All 20 values are numerical in the range (46015:85056).
Take_11 (take):
 Numerical format: All 20 values are numerical in the range (46015:85056).
Take 12 (take):
 Numerical format: All 20 values are numerical in the range (46015:85056).
Take_13 (take):
 Numerical format: All 20 values are numerical in the range (46015:85056).
Take_14 (take):
 Numerical format: All 20 values are numerical in the range (46015:85056).
Take 15 (take):
 Numerical format: All 20 values are numerical in the range (46015:85056).
Take 16 (take):
 Numerical format: All 20 values are numerical in the range (46015:85056).
Take 17 (take):
 Numerical format: All 20 values are numerical in the range (46015:85056).
Take_18 (take):
 Numerical format: All 20 values are numerical in the range (46015:85056).
Take_19 (take):
 Numerical format: All 20 values are numerical in the range (46015:85056).
Take 20 (take):
 Numerical format: All 20 values are numerical in the range (46015:85056).
Take 21 (take):
 Numerical format: All 20 values are numerical in the range (46015:85056).
Take_22 (take):
 Numerical format: All 20 values are numerical in the range (46015:85056).
Take_23 (take):
 Numerical format: All 20 values are numerical in the range (46015:85056).
Take 24 (take):
 Numerical format: All 20 values are numerical in the range (46015:85056).
Take 25 (take):
 Numerical format: All 20 values are numerical in the range (46015:85056).
Take_26 (take):
```

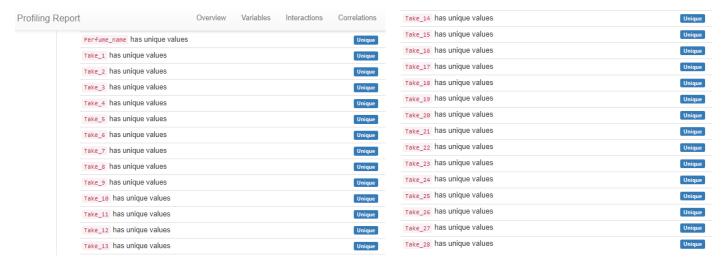
```
Numerical format: All 20 values are numerical in the range (46015:85056).

Take_27 (take):
   Numerical format: All 20 values are numerical in the range (46015:85056).

Take_28 (take):
   Numerical format: All 20 values are numerical in the range (46015:85056).

Last run on: 2024-02-22 13:35:16
```

#### YData found that all columns have Unique values, but did not find the problem in Perfume\_name



### 18- 225 - Restaurant & consumer data – Computer

Dataset analysed: geoplaces2.csv from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00232/RCdata.zip">https://archive.ics.uci.edu/ml/machine-learning-databases/00232/RCdata.zip</a> (dataset\_file\_url column from Fortydatasets.xlsx) has header but it's format is not UTF-8.

```
placeID (id):
    ID column format: All 130 ID values are unique
Alphanumeric range of values: (132560 : 135109)

Alerts

fax has constant value ""
    Constant

dress_code is highly imbalanced (68.2%)
    Imbalance

url is highly imbalanced (63.7%)

Rambience is highly imbalanced (63.7%)
    Imbalance

other_services is highly imbalanced (68.6%)

placeID has unique values

and valid, and thus suitable for use as a Primary Key.

Constant

Imbalance

Imbalance

Unique
```

the geom meter has unique values

YData did not find the many '?'s, but found two Unique values that we also found, as ID and in the categorical column below that has all 130 values as unique, but we have not alerted it. And it found that fax column has only one value, considered Constant, but considered for them a different value than the '?' we found.

```
latitude:
 Latitude format: All 130 values are numerical and valid in the range [-90, 90].
Actual range of values: (18.859803 : 23.7602683)
longitude:
 Longitude format: All 130 values are numerical and valid in the range [-180, 180].
Actual range of values: (-101.0286 : -99.1265059)
the_geom_meter (nominal):
 Error(s) found:
OQI #10 (Structural Conflicts - Consistency, Uniqueness):
Data seems not categorical or has too many categories
                                                        (> 100). Sample values: ['0101000020957F000088568DE356715AC138C0A525FC464A41',
'0101000020957F00001AD016568C4858C1243261274BA54B41',
                                                       '0101000020957F0000649D6F21634858C119AE9BF528A34B41'
0101000020957F00005D67BCDDED8157C1222A2DC8D84D4941
                                                       '0101000020957F00008EBA2D06DC8157C194E03B7B504E4941
'0101000020957F00001B552189B84A58C15A2AAEFD2CA24B41',
                                                       '0101000020957F00008A20F615808157C16272FFCRF84F4941
'0101000020957F00008A2A0747DF4758C11FB31D2A31A84B41
                                                       '0101000020957F0000A478418BBA8057C133851EB22C4E4941
0101000020957F0000A29FAF95CD4958C1FEEEBB73A9914B41']
Categorical format with 130 unique values:
                                          Category Frequency
0101000020957F000000DD3546816E5AC119D4BD17FD544A41
0101000020957F000003B195E25F8457C1C535BD04614B4941
0101000020957F000004457BB7AA8657C15F10835CD9444941
0101000020957F000005810F19B84858C136805B2745A74B41
0101000020957F000000B6735CA004858C108FD525CB2A44B41
0101000020957F000000F14BF6B2C8657C1963CCB8E5C464941
```

```
0101000020957F000011DFB4E3EE4858C13F78758452AB4B41
0101000020957F000011E92CCE714B58C19BF8C0CA75924B41
 0101000020957F000013696871A24558C11EE432FB04A14B41
                                                                                                            1
 0101000020957F0000F9F19DDC3B4858C191B265A83BA44B41
0101000020957F0000FA1A0E5A9B4858C17C884C4173AE4B41
0101000020957F0000FBF7171F056F5AC1F8A6C0A5AF554A41
0101000020957F0000FC60BDA8E88157C1B2C357D6DA4E4941
                                                                                                            1
0101000020957F0000FC799354656C5AC1233FCC70F7584A41
                                                                                                            1
0101000020957F0000FC8866CF17785AC14CD7055B02514A41
                                                                                                            1
0101000020957F0000FD78FFF7AD4458C1896A33F029A04B41
                                                                                                            1
0101000020957F0000FDF8D26EE08157C1FEDB6A1FDB4E4941
                                                                                                            1
0101000020957F0000FE987FAF936D5AC1E1D486215E524A41
                                                                                                            1
0101000020957F0000FEC3FB453E4B58C19B4A463617994B41
 DQI #5 (Extraneous Data - Consistency, Uniqueness):
  3 Extraneous data value(s) at index(es): [(10, 'Restaurante 75'), (66, 'Cenaduria El Rincón de Tlaquepaque'), (103, 'Carnitas Mata
 Calle 16 de Septiembre')]
DQI #15 (Domain Violation - Accuracy):
DQI #15 (Domain Violation - Accuracy):
57 Capitalization/Format issue(s) at index(es): [(1, 'puesto de tacos'), (3, 'little pizza Emilio Portes Gil'), (4, 'carnitas_mata'),
(5, 'Restaurant los Compadres'), (6, 'Taqueria EL amigo'), (7, 'shi ro ie'), (9, 'la Estrella de Dimas'), (12, 'El angel
Restaurante'), (15, 'Tortas y hamburguesas el gordo'), (17, 'rockabilly'), ('...', '...'), (109, 'carnitas mata calle Emilio Portes
Gil'), (110, 'crudalia'), (111, 'tacos de barbacoa enfrente del Tec'), (115, 'Cafeteria cenidet'), (119, 'Cafeteria y Restaurant El
Pacifico'), (122, 'tacos abi'), (123, 'la perica hamburguesa'), (124, 'McDonalds Centro'), (128, 'Restaurant Bar Coty y Pablo'), (129,
'sirloin stockade')] (displaying only the first and last 10 items)
Frequency Distribution:
                                                             name Frequency
                                 Gorditas Dona Tota
                      Abondance Restaurante Bar
                                        Arrachela Grill
                                        Cabana Huasteca
                                             Cafe Chaires
                                                                                   1
                                    Cafeteria cenidet
      Cafeteria y Restaurant El Pacifico
                                                    Carls Jr
Carnitas Mata Calle 16 de Septiembre
               Carreton de Flautas y Migadas
                                                                                   1
                                        puesto de tacos
                                              rockabilly
                                                   shi ro ie
                                       sirloin stockade
                                                  tacos abi
     tacos de barbacoa enfrente del Tec
                               tacos de la estacion
                                                                                    1
                                   tacos los volcanes
                                            tortas hawai
                                                            vips
Range of Values: (Abondance Restaurante Bar to vips)
   Street format: Error(s) found:
DQI #5 (Extraneous Data - Consistency, Uniqueness):
28 Extraneous street data value(s) at index(es): [(7, '?'), (19, '?'), (23, '?'), (27, '?'), (31, '?'), (43, '?'), (45, '?'), (62, '?'), (70, '?'), (72, '?'), ('...'), (95, '?'), (96, '?'), (107, '?'), (110, '?'), (111, '?'), (118, '?'), (122, '?'), (127, '?'), (129, '?')] (displaying only the first and last 10 items)
DQI #15 (Domain Violation - Accuracy):
DQI #15 (Domain Violation - Accuracy):

26 Capitalization/Format issue(s) at index(es): [(1, 'esquina santos degollado y leon guzman'), (3, 'calle emilio portes gil'), (4, 'lic. Emilio portes gil'), (5, 'Camino a Simon Diaz 155 Centro'), (8, 'tampico'), (16, 'carr. mexico'), (17, 'agustin de iturbide'), (26, 'avenida salvador montiel '), (29, 'r.b. anaya esq. florencia'), (34, 'la. de Lozada 1'), ('...', '...'), (80, 'sevilla y olmedo 715 a'), (81, 'circuito oriente esq. carretera 57'), (85, 'Av. Saan Luis entre moctezuma y salinas'), (89, 'Zaragoza entre Francisco Zarco y Lopez Velarde'), (93, 'avenida hivno nacional'), (98, 'Paseo de las Fuentes'), (101, 'Av. de los Pintores'), (106, 'frente al tecnologico'), (114, 'Himno nacional esq. Blvd. Juarez'), (124, 'Rayon sn col. Centro')] (displaying only the first and last 10 items)
city:
  City format: Error(s) found:
DQI #5 (Extraneous Data - Consistency, Uniqueness):

18 Extraneous data value(s) at index(es): [(7, '?'), (27, '?'), (31, '?'), (43, '?'), (45, '?'), (62, '?'), (70, '?'), (72, '?'), (83, '?'), (87, '?'), (95, '?'), (96, '?'), (107, '?'), (110, '?'), (111, '?'), (118, '?'), (127, '?'), (129, '?')]
DQI #15 (Domain Violation - Accuracy):
23 Capitalization/Format issue(s) at index(es): [(1, 's.l.p.'), (3, 'victoria'), (4, 'victoria'), (8, 'victoria'), (17, 'san luis potosi'), (19, 'victoria'), (26, 'cuernavaca'), (29, 'slp'), (39, 'san luis potosi'), (41, 'victoria'), ('...', '...'), (76, 'victoria'), (80, 'san luis potosi'), (81, 'san luis potosi'), (90, 's.l.p'), (93, 'san luis potosi'), (103, 'victoria'), (106, 'victoria'), (109, 'victoria'), (122, 'victoria'), (123, 'victoria')] (displaying only the first and last 10 items)
Frequency Distribution (showing top and bottom 10 of 15 categories):
                   City Frequency
San Luis Potosi
                                            64
                                            18
         Cuernavaca
                                            15
            victoria
                                            12
san luis potosi
                                             6
             Jiutepec
 Ciudad Victoria
```

1

0101000020957F00000F624C60B54958C11757339BCCA24B41

```
Cd Victoria
    Cd. Victoria
       cuernavaca
              s.1.p
              s.1.p.
 san luis potos
                  slp
                                       1
state:
  State format: Error(s) found:
DQI #5 (Extraneous Data - Consistency, Uniqueness):
                                                                                                                                                    '?'), (56, '?'), (70, '?'), (72, '?'), (83, '?'), (129, '?')]
18 Extraneous data value(s) at index(es): [(7, '?'), (27'?'), (87, '?'), (95, '?'), (96, '?'), (107, '?'), (110,
                                                                                                       (31, '?'), (43, '?'),
111, '?'), (118, '?'),
DQI #15 (Domain Violation - Accuracy):
20 Capitalization 'Accurate Accurate (s) at index(es): [(1, 's.l.p.'), (3, 'tamaulipas'), (17, 'san luis potosi'), (19, 'tamaulipas'), (26, 'morelos'), (29, 'slp'), (39, 'san luis potosi'), (46, 'mexico'), (66, 'san luis potosi'), (73, 'tamaulipas'), (76, 'tamaulipas'), (80, 'san luis potosi'), (81, 'slp'), (90, 'mexico'), (93, 'san luis potosi'), (103, 'tamaulipas'), (106, 'tamaulipas'), (122, 'tamaulipas'), (123, 'tamaulipas')]
Frequency Distribution (showing top and bottom 10 of 13 categories):
               State Frequency
                  SLP
                                      50
                                      19
            Morelos
                                      18
San Luis Potosi
                                      14
       tamaulipas
        Tamaulipas
san luis potosi
                                       4
              S.L.P.
                                       2
              mexico
                                       2
                  slp
                                       2
            morelos
                                       1
              s.1.p.
                                       1
 san luis potos
                                       1
country:
  Country format: Error(s) found:
DQI #5 (Extraneous Data - Consistency, Uniqueness):
28 Extraneous data value(s) at index(es): [(3, '?'), (7, '?'), (27, '?'), (31, '?'), (39, '?'), (43, '?'), (45, '?'), (46, '?') ('?'), (70, '?'), ('...', '...'), (98, '?'), (103, '?'), (107, '?'), (109, '?'), (110, '?'), (111, '?'), (118, '?'), (123, '?'), (129, '?')] (displaying only the first and last 10 items)
DQI #15 (Domain Violation - Accuracy):
13 Capitalization/Format issue(s) at index(es): [(1, 'mexico'), (17, 'mexico'), (19, 'mexico'), (26, 'mexico'), (29, 'mexico'), (66, 'mexico'), (76, 'mexico'), (80, 'mexico'), (81, 'mexico'), (85, 'mexico'), (93, 'mexico'), (106, 'mexico'), (122, 'mexico')]
Frequency Distribution:
Country
              Frequency
 Mexico
                         28
                         13
fax:
  Phone format: Error(s) found:
DQI #15 (Domain Violation - Accuracy):
 130 Incorrect telephone number format issue(s) at index(es): [(0, '?'), (1, '?'), (2, '?'), (3, '?'), (4, '?'), (5, '?'), (6, '?'), (7, '?'), (8, '?'), (9, '?'), ('...', '...'), (120, '?'), (121, '?'), (122, '?'), (123, '?'), (124, '?'), (125, '?'), (126, '?'), (127, '?'), (128, '?'), (129, '?')] (displaying only the first and last 10 items)
 fax has constant value ""
zin:
  Postal Code format: Error(s) found:
DQI #17 (Wrong Data Type - Consistency):
74 Non-alphanumeric value(s) at index(es): [(0, '?'), (3, '?'), (4, '?'), (7, '?'), (8, '?'), (15, '?'), (16, '?'), (17, '?'), (18, '?'), (19, '?'), ('...', '...'), (115, '?'), (117, '?'), (118, '?'), (120, '?'), (122, '?'), (123, '?'), (125, '?'), (127, '?'), (128, '?'), (129, '?')] (displaying only the first and last 10 items)
alcohol (nominal):
  All 130 values are correctly categorical.
Categorical format with 3 unique values:
              Category Frequency
No_Alcohol_Served
                                         87
            Wine-Beer
                                         34
             Full_Bar
smoking_area (nominal):
   All 130 values are correctly categorical.
Categorical format with 5 unique values:
        Category Frequency
              none
                                  70
not permitted
                                  25
```

Soledad

```
section
                          24
    permitted
                           9
  only at bar
dress_code (code):
  All 130 values are correctly categorical.
Categorical format with 3 unique values:
Category Frequency
informal
                  118
  casual
                    10
  formal
accessibility (nominal):
  All 130 values are correctly categorical.
Categorical format with 3 unique values:
         Category Frequency
no_accessibility
       completely
                              45
        partially
                               9
price (nominal):
  All 130 values are correctly categorical.
Categorical format with 3 unique values:
Category Frequency
  medium
                    60
     low
                    45
    high
                    25
  URL format: Error(s) found:
ONE FORMACL CRIMIC (3) Found.

DQI #15 (Domain Violation - Accuracy):

117 Invalid URL format issue(s) at index(es): [(1, '?'), (2, '?'), (3, '?'), (4, '?'), (5, '?'), (6, '?'), (7, '?'), (8, '?'), (9, '?'), (10, '?'), ('...', '...'), (119, '?'), (120, '?'), (121, '?'), (122, '?'), (123, '?'), (124, 'no'), (125, '?'), (127, '?'), (129, '?')] (displaying only the first and last 10 items)
Rambience (nominal):
 All 130 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
familiar
                 121
   quiet
                     9
franchise (ranch):
  All 130 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
                  108
                    22
area (nominal):
  All 130 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
  closed
                  115
    open
                    15
other_services (service):
  All 130 values are correctly categorical.
Categorical format with 3 unique values:
Category Frequency
    none
 variety
Internet
                      4
Last run on: 2024-02-22 15:41:41
```

## 19- 17 - Breast Cancer Wisconsin (Diagnostic) – Life

Dataset analysed: wdbc.data from <a href="https://archive.ics.uci.edu/static/public/17/breast+cancer+wisconsin+diagnostic.zip">https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/breast-cancer-wisconsin.data</a> that was in the dataset\_file\_url column from Fortydatasets.xlsx. Curiously this same dataset is being used for a similar dataset, 15 - Breast Cancer Wisconsin (Original) – Life, that is presented later on this document. It does not have header.

ID number (id): ID column format: All 569 ID values are unique and valid, and thus suitable for use as a Primary Key. Alphanumeric range of values: (842302 : 92751) Alerts 7 Reproduction Alerts ID number has unique values Unique Cell Nucleus 1 - g) concavity has 13 (2.3%) zeros Zeros Cell Nucleus 1 - h) concave points has 13 (2.3%) zeros Zeros Cell Nucleus 2 - g) concavity has 13 (2.3%) zeros Cell Nucleus 2 - h) concave points has 13 (2.3%) zeros Cell Nucleus 3 - g) concavity has 13 (2.3%) zeros Cell Nucleus 3 - h) concave points has 13 (2.3%) zeros

YData shows the column 'ID number' as Unique. We consider it not only unique but suitable for use as a Primary Key.

```
All 569 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
                212
Cell Nucleus 1 - a) radius (radius):
 Numerical >=0 format: All 569 values are numerical and greater or equal to 0 in the range (6.981:28.11).
Cell Nucleus 1 - b) texture (texture):
 Numerical format: All 569 values are numerical in the range (9.71:39.28).
Cell Nucleus 1 - c) perimeter (perimeter):
 Numerical >=0 format: All 569 values are numerical and greater or equal to 0 in the range (43.79:188.5).
Cell Nucleus 1 - d) area (area):
 Numerical >=0 format: All 569 values are numerical and greater or equal to 0 in the range (143.5:2501.0).
Cell Nucleus 1 - e) smoothness (smoothness):
 Numerical >=0 format: All 569 values are numerical and greater or equal to 0 in the range (0.05263:0.1634).
Cell Nucleus 1 - f) compactness (compactness):
 Numerical >=0 format: All 569 values are numerical and greater or equal to 0 in the range (0.01938:0.3454).
Cell Nucleus 1 - g) concavity (concavity):
 Numerical format: All 569 values are numerical in the range (0.0:0.4268).
Cell Nucleus 1 - h) concave points (points):
 Numerical >=0 format: All 569 values are numerical and greater or equal to 0 in the range (0.0:0.2012).
Cell Nucleus 1 - i) symmetry (symmetry):
 Numerical format: All 569 values are numerical in the range (0.106:0.304).
Cell Nucleus 1 - j) fractal dimension (dimension):
 Numerical format: All 569 values are numerical in the range (0.04996:0.09744).
Cell Nucleus 2 - a) radius (radius):
 Numerical >=0 format: All 569 values are numerical and greater or equal to 0 in the range (0.1115:2.873).
Cell Nucleus 2 - b) texture (texture):
 Numerical format: All 569 values are numerical in the range (0.3602:4.885).
Cell Nucleus 2 - c) perimeter (perimeter):
 Numerical >=0 format: All 569 values are numerical and greater or equal to 0 in the range (0.757:21.98).
Cell Nucleus 2 - d) area (area):
 Numerical >=0 format: All 569 values are numerical and greater or equal to 0 in the range (6.802:542.2).
Cell Nucleus 2 - e) smoothness (smoothness):
 Numerical >=0 format: All 569 values are numerical and greater or equal to 0 in the range (0.001713:0.03113).
Cell Nucleus 2 - f) compactness (compactness):
 Numerical >=0 format: All 569 values are numerical and greater or equal to 0 in the range (0.002252:0.1354).
Cell Nucleus 2 - g) concavity (concavity):
  Numerical format: All 569 values are numerical in the range (0.0:0.396).
Cell Nucleus 2 - h) concave points (points):
 Numerical >=0 format: All 569 values are numerical and greater or equal to 0 in the range (0.0:0.05279).
Cell Nucleus 2 - i) symmetry (symmetry):

Numerical format: All 569 values are numerical in the range (0.007882:0.07895).
```

```
Cell Nucleus 2 - j) fractal dimension (dimension):
 Numerical format: All 569 values are numerical in the range (0.0008948:0.02984).
Cell Nucleus 3 - a) radius (radius):
 Numerical >=0 format: All 569 values are numerical and greater or equal to 0 in the range (7.93:36.04).
Cell Nucleus 3 - b) texture (texture):
 Numerical format: All 569 values are numerical in the range (12.02:49.54).
Cell Nucleus 3 - c) perimeter (perimeter):
 Numerical >=0 format: All 569 values are numerical and greater or equal to 0 in the range (50.41:251.2).
Cell Nucleus 3 - d) area (area):
 Numerical >=0 format: All 569 values are numerical and greater or equal to 0 in the range (185.2:4254.0).
Cell Nucleus 3 - e) smoothness (smoothness):
 Numerical >=0 format: All 569 values are numerical and greater or equal to 0 in the range (0.07117:0.2226).
Cell Nucleus 3 - f) compactness (compactness):
 Numerical >=0 format: All 569 values are numerical and greater or equal to 0 in the range (0.02729:1.058).
Cell Nucleus 3 - g) concavity (concavity):
 Numerical format: All 569 values are numerical in the range (0.0:1.252).
Cell Nucleus 3 - h) concave points (points):
 Numerical >=0 format: All 569 values are numerical and greater or equal to 0 in the range (0.0:0.291).
Cell Nucleus 3 - i) symmetry (symmetry):
 Numerical format: All 569 values are numerical in the range (0.1565:0.6638).
Cell Nucleus 3 - j) fractal dimension (dimension):
 Numerical format: All 569 values are numerical in the range (0.05504:0.2075).
Last run on: 2024-02-23 18:27:01
                   850 - Raisin Dataset - Life
         20-
Dataset analysed: Raisin_Dataset\Raisin_Dataset.xlsx from <a href="https://archive.ics.uci.edu/ml/machine-learning-">https://archive.ics.uci.edu/ml/machine-learning-</a>
databases/00617/Raisin Dataset.zip (dataset file url column from Fortydatasets.xlsx) has header
 Numerical >=0 format: All 900 values are numerical and greater or equal to 0 in the range (25387:235047).
Perimeter:
 Numerical >=0 format: All 900 values are numerical and greater or equal to 0 in the range (225.629541:997.2919406).
MajorAxisLength (length):
 Numerical >=0 format: All 900 values are numerical and greater or equal to 0 in the range (143.7108718:492.2752785).
 Numerical format: All 900 values are numerical in the range (0.348729642:0.96212444).
 Numerical format: All 900 values are numerical in the range (26139:278217).
ConvexArea (area):
 Numerical >=0 format: All 900 values are numerical and greater or equal to 0 in the range (0.379856115:0.835454545).
Extent (ratio):
 Numerical format: All 900 values are numerical in the range (619.074:2697.753).
Class:
 All 900 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
  Besni
                450
 Kecimen
                450
Last run on: 2024-02-23 16:02:26
        Alerts 7
                  Reproduction
Overview
 Alerts
  class is uniformly distributed
 Perimeter has unique values
```

YData shows 6 columns as Unique. We have not checked that.

MajorAxisLength has unique values

MinorAxisLength has unique values

ConvexArea has unique values

#### 21- 1 - Abalone – Life

Dataset analysed: abalone.data from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/abalone.data">https://archive.ics.uci.edu/ml/machine-learning-databases/abalone.data</a> (dataset\_file\_url column from Fortydatasets.xlsx) It does not have header

```
Sex:
 All 4177 values are correctly categorical.
Categorical format with 3 unique values:
Category Frequency
               1528
       Т
               1342
               1307
Length:
 Numerical >=0 format: All 4177 values are numerical and greater or equal to 0 in the range (0.075:0.815).
Diameter:
 Numerical >=0 format: All 4177 values are numerical and greater or equal to 0 in the range (0.055:0.65).
 Numerical format: All 4177 values are numerical in the range (0.0:1.13).
 Numerical >=0 format: All 4177 values are numerical and greater or equal to 0 in the range (0.002:2.8255).
Shucked weight (weight):
 Numerical >=0 format: All 4177 values are numerical and greater or equal to 0 in the range (0.001:1.488).
Viscera weight (weight):
 Numerical >=0 format: All 4177 values are numerical and greater or equal to 0 in the range (0.0005:0.76).
Shell weight (weight):
 Numerical >=0 format: All 4177 values are numerical and greater or equal to 0 in the range (0.0015:1.005).
 Numerical format: All 4177 values are numerical in the range (1:29).
Last run on: 2024-02-23 16:26:24
```

#### YData did not find any Alerts in this Dataset

## 22- 15 - Breast Cancer Wisconsin (Original) – Life

Dataset analysed: breast-cancer-wisconsin.data from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/breast-cancer-wisconsin.data">https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/breast-cancer-wisconsin.data</a> (dataset\_file\_url column from Fortydatasets.xlsx) It does not have header.

Overview



Reproduction

#### Alerts

Dataset has 8 (1.1%) duplicate rows

Duplicates

YData only found Duplicate entire rows. We found a duplicate ID column, and other problem in the first ID column. We also found '?' in 'Bare Nuclei'

```
Sample code number (id):
    TD column format: Error(s) found:
    DQI #9 (Duplicates - Uniqueness):
    100 Duplicate value(s) at index(es): [[4, 1017023], [8, 1033078], [9, 1033078], [29, 1070935], [30, 1070935], [42, 1100524], [47, 1105524], [61, 1115293], [62, 1116116], [64, 1116192], ('...', '...'), [661, 1339781], [672, 1354840], [673, 1354840], [683, 466906], [684, 466906], [689, 654546], [690, 654546], [691, 695091], [697, 897471], [698, 897471]] (displaying only the first and last 10 items)

DQI #17 (Wrong Data Type - Consistency):
    188 Inconsistent length in alphanumeric value(s) at index(es): [(243, 128059), (246, 144888), (247, 145447), (248, 167528), (249, 169356), (250, 183913), (251, 191250), (259, 242970), (260, 255644), (261, 263538), ('...', '...'), (689, 654546), (690, 654546), (691, 695091), (692, 714039), (693, 763235), (694, 776715), (695, 841769), (696, 888820), (697, 897471), (698, 897471)] (displaying only the first and last 10 items)

DQI #19 (Uniqueness violation(s) at index(es): [(9, 1033078), (30, 1070935), (82, 1143978), (109, 1171710), (116, 1173347), (121, 1174057), (195, 1212422), (208, 1218860), (252, 1017023), (253, 1100524), ('...', '...'), (618, 1061990), (632, 1238777), (639, 1277792), (644, 1299596), (661, 1339781), (673, 1354840), (684, 466906), (690, 654546), (691, 695091), (698, 897471)] (displaying only the first and last 10 items)

Alphanumeric range of values: (1000025 : 95719)

Clump Thickness (thickness):
    Numerical format: All 699 values are numerical in the range (1:10).
```

```
Uniformity of Cell Size (size):
  Numerical >=0 format: All 699 values are numerical and greater or equal to 0 in the range (1:10).
Uniformity of Cell Shape (uniformity):
  Numerical format: All 699 values are numerical in the range (1:10).
Marginal Adhesion (adhesion):
  Numerical format: All 699 values are numerical in the range (1:10).
Single Epithelial Cell Size (size):
  Numerical >=0 format: All 699 values are numerical and greater or equal to 0 in the range (1:10).
Bare Nuclei (nuclei):
Numerical format: Error(s) found:

DQI #17 (Wrong Data Type - Consistency):

16 Non-numeric value(s) at index(es): [(23, '?'), (40, '?'), (139, '?'), (145, '?'), (158, '?'), (164, '?'), (235, '?'), (249, '?'), (275, '?'), (292, '?'), (294, '?'), (315, '?'), (321, '?'), (411, '?'), (617, '?')]
Range of values: (1.0:10.0).
Bland Chromatin (chromatin):
  Numerical format: All 699 values are numerical in the range (1:10).
Normal Nucleoli (nucleoli):
  Numerical format: All 699 values are numerical in the range (1:10).
Mitoses:
  Numerical format: All 699 values are numerical in the range (1:10).
  All 699 values are correctly categorical.
Categorical format with 2 unique values:
 Category Frequency
                   458
         4
                   241
Last run on: 2024-02-23 18:29:53
```

## 23- 73 - Mushroom - Life

Dataset analysed: agaricus-lepiota.data from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/mushroom/agaricus-lepiota.data">https://archive.ics.uci.edu/ml/machine-learning-databases/mushroom/agaricus-lepiota.data</a> (dataset\_file\_url column from Fortydatasets.xlsx) It does not have header.

All fields are categorical and have been altered manually in the AllColumnsfromFortyDatasets file.

```
poisonous (categorical):
 All 8124 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
       р
               3916
cap-shape (categorical):
  All 8124 values are correctly categorical.
Categorical format with 6 unique values:
Category Frequency
               3656
               3152
                828
       b
                452
                 32
cap-surface (categorical):
  All 8124 values are correctly categorical.
Categorical format with 4 unique values:
Category Frequency
               3244
               2556
       S
               2320
       g
cap-color (color):
  All 8124 values are correctly categorical.
Categorical format with 10 unique values:
Category Frequency
       n
               2284
       g
               1840
               1500
       e
               1072
       У
               1040
                168
                144
       р
```

```
bruises? (categorical):
  All 8124 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
f 4748
        t
                  3376
odor (categorical):
  All 8124 values are correctly categorical.
Categorical format with 9 unique values:
Category Frequency
                  3528
        n
        s
                   576
                    576
        у
                    400
        1
                    400
        р
                   256
                   192
        m
                    36
gill-attachment (categorical):
  All 8124 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
                  7914
                   210
gill-spacing (categorical):
  All 8124 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
                  6812
        C
                  1312
        W
gill-size (categorical):
  All 8124 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
        b
                  5612
        n
                  2512
gill-color (color):
  All 8124 values are correctly categorical.
Categorical format with 12 unique values:
Category Frequency
        b
        р
                  1202
        n
                  1048
        g
                   752
        h
                   732
        u
                   492
        k
                   408
                    96
        е
        У
                     86
                     64
        0
                     24
stalk-shape (categorical):
  All 8124 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
        t
                  4608
        е
                  3516
stalk-root (categorical):
    Error(s) found:
DQI #4 (Ambiguous Data - Accuracy, Consistency):
2480 Unacceptable value(s) at index(es): [(3984, '?'), (4023, '?'), (4076, '?'), (4100, '?'), (4104, '?'), (4196, '?'), (4200, '?') (4283, '?'), (4291, '?'), (4326, '?'), ('...', '...'), (8113, '?'), (8115, '?'), (8116, '?'), (8117, '?'), (8118, '?'), (8120, '?'), (8121, '?'), (8122, '?'), (8123, '?')] (displaying only the first and last 10 items)

Categorical format with 5 unique values:
Category Frequency
                  3776
        ?
                  2480
        e
                  1120
        С
                   556
                   192
stalk-surface-above-ring (categorical):
  All 8124 values are correctly categorical.
```

16

Categorical format with 4 unique values:

u

```
Category Frequency
               5176
               2372
                552
                 24
stalk-surface-below-ring (categorical):
 All 8124 values are correctly categorical.
Categorical format with 4 unique values:
Category Frequency
               4936
               2304
       f
                600
                284
      У
stalk-color-above-ring (color):
 All 8124 values are correctly categorical.
Categorical format with 9 unique values:
Category Frequency
               4464
       р
               1872
       g
                576
                448
      n
      b
                432
                192
      0
                 96
       е
                 36
       C
                 8
      У
stalk-color-below-ring (color):
 All 8124 values are correctly categorical.
Categorical format with 9 unique values:
Category Frequency
               4384
       р
               1872
                576
       g
                512
      n
      b
                432
      0
                192
       е
                 36
                 24
veil-type (type):
 All 8124 values are correctly categorical.
Categorical format with 1 unique values:
Category Frequency
               8124
      р
Alerts
 veil-type has constant value ""
```

YData just found that this column has only one value. We showed it but have not considered an Alert. Besides that it did not find that 'stalk-root' had '?' values

#### Other Alerts they show are High Correlation and Imbalance

```
veil-color (color):
 All 8124 values are correctly categorical.
Categorical format with 4 unique values:
Category Frequency
               7924
       n
                 96
       О
                 96
                  8
ring-number (categorical):
 All 8124 values are correctly categorical.
Categorical format with 3 unique values:
Category Frequency
              7488
      0
       t
                600
ring-type (type):
 All 8124 values are correctly categorical.
Categorical format with 5 unique values:
Category Frequency
               3968
      р
               2776
       е
               1296
       1
       f
                 48
                 36
```

```
spore-print-color (color):
  All 8124 values are correctly categorical.
Categorical format with 9 unique values:
Category Frequency
               2388
       n
               1968
       k
               1872
       h
               1632
                 72
                 48
       h
                 48
       0
                 48
       u
       У
population (categorical):
  All 8124 values are correctly categorical.
Categorical format with 6 unique values:
Category Frequency
               4949
               1712
       ς
               1248
                400
       n
                384
       а
                340
       C
habitat (categorical):
  All 8124 values are correctly categorical.
Categorical format with 7 unique values:
Category Frequency
               3148
       g
               2148
               1144
       1
                832
       u
                368
                292
       m
                192
Last run on: 2024-02-23 20:33:28
```

### 24- 14 - Breast Cancer – Life

Dataset analysed: breast-cancer.data from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer/breast-cancer.data">https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer.data</a> (dataset\_file\_url column from Fortydatasets.xlsx) It does not have header.



YData just found whole Duplicate lines, High correlation and Imbalance. Not the 2 cases of '?' values.

```
All 286 values are correctly categorical.
Categorical format with 2 unique values:
           Category Frequency
no-recurrence-events
                            201
  recurrence-events
 All 286 values are correctly categorical.
Categorical format with 6 unique values:
Category Frequency
  50-59
                 96
   40-49
                 90
   60-69
                 57
   30-39
                 36
   70-79
                 6
   20-29
```

menopause

All 286 values are correctly categorical.

```
Category Frequency
premeno
                 150
    ge40
                 129
    1t40
tumor-size (0-):
  All 286 values are correctly categorical.
Categorical format with 11 unique values:
Category Frequency
   30-34
                  60
   25-29
                  54
   20-24
                  50
   15-19
                  30
   10-14
                  28
   40-44
                  22
   35-39
                  19
     0-4
                   8
   50-54
                   8
     5-9
                   4
   45-49
                   3
inv-nodes (node):
 All 286 values are correctly categorical.
Categorical format with 7 unique values:
Category Frequency
     0-2
     3-5
                  36
     6-8
                  17
    9-11
                  10
   15-17
                   6
   12-14
                   3
   24-26
                   1
node-caps (node):
Error(s) found:

DQI #4 (Ambiguous Data - Accuracy, Consistency):

8 Unacceptable value(s) at index(es): [(145, '?'), (163, '?'), (164, '?'), (183, '?'), (184, '?'), (233, '?'), (263, '?'), (264, '?')]
Categorical format with 3 unique values:
Category Frequency
      no
                 222
     yes
                  56
deg-malig (degree):
  All 286 values are correctly categorical.
Categorical format with 3 unique values:
Category Frequency
                  130
         3
                   85
        1
                   71
breast:
  All 286 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
   left
                 152
   right
                 134
breast-quad (breast):
 Error(s) found:
DQI #4 (Ambiguous Data - Accuracy, Consistency):

1 Unacceptable value(s) at index(es): [(206, '?')]
Categorical format with 6 unique values:
 Category Frequency
 left_low
                  110
 left_up
                   97
 right_up
                   33
right_low
                   24
  central
                   21
irradiat (yes):
 All 286 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
                 218
      no
     yes
                  68
Last run on: 2024-02-25 11:33:55
```

Categorical format with 3 unique values:

#### 25- 236 - seeds – Life

Dataset analysed: seeds\_dataset.txt from <u>archive.ics.uci.edu/ml/machine-learning-databases/00236/seeds\_dataset.txt</u> (dataset\_file\_url column from Fortydatasets.xlsx) t does not have header.

```
2.221
15.26
          14.84
                    0.871
                               5.763
                                         3.312
                                                              5.22
14.88
          14.57
                    0.8811
                              5.554
                                                   1.018
                                                              4.956
                                                                        1
                                         3.333
14.29
          14.09
                    0.905
                               5.291
                                                   2.699
                                                              4.825
                                         3.337
                                                                        1
                    0.8955
                                                   2.259
                                                              4.805
13.84
          13.94
                               5.324
                                         3.379
                                                                        1
16.14
          14.99
                    0.9034
                               5.658
                                         3.562
                                                   1.355
                                                              5.175
14.38
          14.21
                    0.8951
                               5.386
                                         3.312
                                                    2.462
                                                              4.956
14.69
          14.49
                    0.8799
                               5.563
                                         3.259
                                                    3.586
                                                              5.219
14.11
          14.1
                    0.8911
                               5.42
                                         3.302
          15.46
                    0.8747
                               6.053
                                         3.465
                                                   2.04
                                                              5.877
16.63
                                                                        1
16.44
          15.25
                    0.888
                               5.884
                                         3.505
                                                    1.969
                                                              5.533
                                                                        1
15.26
          14.85
                    0.8696
                               5.714
                                         3.242
                                                   4.543
                                                              5.314
```

#### The code was altered so that two tabs now are just one.

```
area A (area):
  Numerical >=0 format: All 210 values are numerical and greater or equal to 0 in the range (10.59:21.18).
perimeter P (perimeter):
  Numerical >=0 format: All 210 values are numerical and greater or equal to 0 in the range (12.41:17.25).
compactness C = 4*pi*A (compactness):
  Numerical >=0 format: All 210 values are numerical and greater or equal to 0 in the range (0.8081:0.9183).
length of kernel (length):
  Numerical >=0 format: All 210 values are numerical and greater or equal to 0 in the range (4.899:6.675).
width of kernel (width):
  Numerical >=0 format: All 210 values are numerical and greater or equal to 0 in the range (2.63:4.033).
asymmetry coefficient (coefficient):
  Numerical format: All 210 values are numerical in the range (0.7651:8.456).
length of kernel groove (length):
  Numerical >=0 format: All 210 values are numerical and greater or equal to 0 in the range (4.519:6.55).
  All 210 values are correctly categorical.
Categorical format with 3 unique values:
 Category Frequency
         1
                    70
                    70
         2
         3
                    70
Last run on: 2024-02-25 12:14:03
           Alerts 8
                        Reproduction
Overview
  Alerts
   area A is highly overall correlated with class and 5 other fields
   class is highly overall correlated with area A and 1 other fields
                                                                          High correlation
   compactness C = 4*pi*A is highly overall correlated with area A and 3 other fields
   length of kernel is highly overall correlated with area A and 3 other fields
   length of kernel groove is highly overall correlated with area A and 3 other fields
   perimeter P is highly overall correlated with area A and 4 other fields
   width of kernel is highly overall correlated with area A and 4 other fields
```

YData did not find anything that we care.

class is uniformly distributed

#### 26- 111 - Zoo - Life

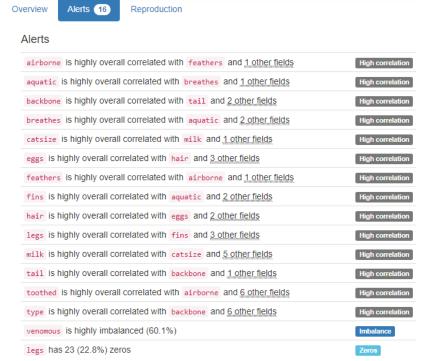
Dataset analysed: zoo.data from <u>archive.ics.uci.edu/ml/machine-learning-databases/zoo/zoo.data</u> (dataset\_file\_url column from Fortydatasets.xlsx) t does not have header.

```
animal name (name):
Name format: All 101 animal name values are valid.

Frequency Distribution:
animal name Frequency
frog 2
aardvark 1
antelope 1
bass 1
bear 1
```

```
boar
                   1
   buffalo
      calf
       carp
   catfish
                   1
    tuatara
      tuna
                    1
    vampire
      vole
   vulture
   wallaby
                   1
      wasp
                   1
      wolf
                    1
      worm
      wren
Range of Values: (aardvark to wren)
hair (boolean):
Binary format: All 101 binary values are valid.
Frequency Distribution:
Value Frequency
    0
              58
              43
    1
feathers (boolean):
 Binary format: All 101 binary values are valid.
Frequency Distribution:
Value Frequency
    0
              81
    1
               20
eggs (boolean):
 Binary format: All 101 binary values are valid.
Frequency Distribution:
Value Frequency
     0
               42
milk (boolean):
 Binary format: All 101 binary values are valid.
Frequency Distribution:
Value Frequency
    0
              60
               41
    1
airborne (boolean):
 Binary format: All 101 binary values are valid.
Frequency Distribution:
Value Frequency
    0
               24
    1
aquatic (boolean):
 Binary format: All 101 binary values are valid.
Frequency Distribution:
Value Frequency
    0
    1
               36
predator (boolean):
 Binary format: All 101 binary values are valid.
Frequency Distribution:
Value Frequency
              56
    1
               45
toothed (boolean):
 Binary format: All 101 binary values are valid.
Frequency Distribution:
Value Frequency
    1
              61
    0
              40
backbone (boolean):
 Binary format: All 101 binary values are valid.
```

```
Frequency Distribution:
 Value Frequency
     0
               18
breathes (boolean):
  Binary format: All 101 binary values are valid.
Frequency Distribution:
 Value Frequency
               80
     1
     0
               21
venomous (boolean):
  Binary format: All 101 binary values are valid.
Frequency Distribution:
 Value Frequency
    0
               93
     1
                8
fins (boolean):
  Binary format: All 101 binary values are valid.
Frequency Distribution:
 Value Frequency
     1
               17
legs (numeric):
  Numerical format: All 101 values are numerical in the range (0:8).
tail (boolean):
  Binary format: All 101 binary values are valid.
Frequency Distribution:
 Value Frequency
     0
               26
domestic (boolean):
   Binary format: All 101 binary values are valid.
Frequency Distribution:
 Value Frequency
    0
               88
               13
     1
catsize (boolean):
  Binary format: All 101 binary values are valid.
Frequency Distribution:
 Value Frequency
     0
     1
               44
type (numeric):
  Numerical format: All 101 values are numerical in the range (1:7).
Last run on: 2024-02-25 16:52:25
```



YData did not find anything that we care.

Elevation:

### 27- 31 - Covertype – Life

Dataset analysed: covtype.data from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/covtype/covtype.data.gz">https://archive.ics.uci.edu/ml/machine-learning-databases/covtype.data.gz</a> (dataset\_file\_url column from Fortydatasets.xlsx) It does not have header.

```
Numerical format: All 581012 values are numerical in the range (1859:3858).
Aspect (speed):
  Numerical format: All 581012 values are numerical in the range (0:360).
  Numerical format: All 581012 values are numerical in the range (0:66).
Horizontal_Distance_To_Hydrology (distance):
  Numerical format: All 581012 values are numerical in the range (0:1397).
Vertical_Distance_To_Hydrology (distance):
  Numerical format: All 581012 values are numerical in the range (-173:601).
Horizontal_Distance_To_Roadways (distance):
  Numerical format: All 581012 values are numerical in the range (0:7117).
Hillshade 9am (index):
  Numerical format: All 581012 values are numerical in the range (0:254).
Hillshade Noon (index):
  Numerical format: All 581012 values are numerical in the range (0:254).
Hillshade_3pm (index):
  Numerical format: All 581012 values are numerical in the range (0:254).
Horizontal_Distance_To_Fire_Points (distance):
  Numerical format: All 581012 values are numerical in the range (0:7173).
  Binary format: All 581012 binary values are valid.
    quency Distribution:
 Value
        Frequency
           320216
Wilderness_Area_2 (binary):
Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value
        Frequency
     a
           551128
     1
            29884
Wilderness_Area_3 (binary):
  Binary format: All 581012 binary values are valid.
```

```
Frequency Distribution:
 Value Frequency
           327648
           253364
Wilderness_Area_4 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
           544044
     0
            36968
     1
Soil_Type_1 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
    0
         577981
             3031
     1
Soil_Type_2 (binary):
   Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
     0
          573487
            7525
Soil_Type_3 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
     0
           576189
            4823
Soil_Type_4 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
           568616
     0
            12396
     1
Soil_Type_5 (binary):
   Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
    0
        579415
     1
             1597
Soil_Type_6 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
       574437
     0
            6575
Soil_Type_7 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency 0 580907
              105
     1
Soil_Type_8 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
    0
         580833
     1
              179
Soil_Type_9 (binary):
Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
```

```
1147
Soil_Type_10 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
     0
           548378
     1
            32634
Soil_Type_11 (binary):
   Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
     0
           568602
     1
            12410
Soil_Type_12 (binary):
   Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
     0
           551041
     1
            29971
Soil_Type_13 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
           563581
     0
            17431
     1
Soil_Type_14 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
     0
           580413
     1
              599
Soil_Type_15 (binary):
   Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
     0
          581009
Soil_Type_16 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
           578167
     0
     1
             2845
Soil_Type_17 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
     0
           577590
     1
             3422
Soil_Type_18 (binary):
   Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
         579113
             1899
Soil_Type_19 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
           576991
     0
     1
             4021
```

0

579865

```
Soil_Type_20 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
           571753
     1
             9259
Soil_Type_21 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
           580174
     0
Soil_Type_22 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
           547639
     0
           33373
     1
Soil Type 23 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
     0
        523260
     1
           57752
Soil_Type_24 (binary):
   Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
           559734
     1
            21278
Soil_Type_25 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
          580538
    0
             474
     1
Soil_Type_26 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
     0
          578423
            2589
     1
Soil_Type_27 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
     0
           579926
     1
            1086
Soil_Type_28 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
           580066
     0
Soil_Type_29 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
           465765
     0
           115247
     1
Soil_Type_30 (binary):
  Binary format: All 581012 binary values are valid.
```

```
Frequency Distribution:
 Value Frequency
           550842
            30170
Soil_Type_31 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
           555346
     0
            25666
     1
Soil_Type_32 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
    0
         528493
     1
            52519
Soil_Type_33 (binary):
   Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
     0
           535858
            45154
Soil_Type_34 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
     0
           579401
            1611
Soil_Type_35 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
           579121
     0
     1
             1891
Soil_Type_36 (binary):
   Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
    0
         580893
     1
              119
Soil_Type_37 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
       580714
     0
              298
Soil_Type_38 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
           565439
     0
            15573
     1
Soil_Type_39 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
    0
         567206
     1
            13806
Soil_Type_40 (binary):
  Binary format: All 581012 binary values are valid.
Frequency Distribution:
 Value Frequency
```

```
Cover_Type (integer):
 Numerical format: All 581012 values are numerical in the range (1:7).
Last run on: 2024-02-25 17:56:27
Overview
               Alerts 57
                                   Reproduction
```

### Alerts

0

572262 8750



YData did not find anything that we care.

#### 28-19 - Car Evaluation - N/A

Dataset analysed: car.data from https://archive.ics.uci.edu/ml/machine-learningdatabases/car/car.data (dataset\_file\_url column from Fortydatasets.xlsx) It does not have header.

```
buying (high):
  All 1728 values are correctly categorical.
Categorical format with 4 unique values:
Category Frequency
                432
    high
                432
    low
                432
     med
   vhigh
                432
maint (high):
  All 1728 values are correctly categorical.
Categorical format with 4 unique values:
Category Frequency
   high
                432
    low
                432
                432
    med
   vhigh
                432
doors (more):
 All 1728 values are correctly categorical.
Categorical format with 4 unique values:
Category Frequency
                432
       3
                432
       4
                432
   5more
                432
persons (more):
 All 1728 values are correctly categorical.
Categorical format with 3 unique values:
Category Frequency
       4
                576
    more
                576
lug_boot (big):
  All 1728 values are correctly categorical.
Categorical format with 3 unique values:
```

```
Category Frequency
                576
     big
                 576
     med
   small
safety (high):
  All 1728 values are correctly categorical.
Categorical format with 3 unique values:
Category Frequency
   high
                576
                 576
     low
     med
                 576
Class Values (class):
  All 1728 values are correctly categorical.
Categorical format with 4 unique values:
Category Frequency
   unacc
                1210
     acc
                 384
    good
                 69
   vgood
                  65
Overview
             Alerts 6
                            Reproduction
 Alerts
  buying is uniformly distributed
                                                                                 Uniform
  maint is uniformly distributed
                                                                                 Uniform
  doors is uniformly distributed
                                                                                 Uniform
   persons is uniformly distributed
                                                                                 Uniform
  lug_boot is uniformly distributed
   safety is uniformly distributed
```

YData did not find anything that we care.

### 29- 10 - Automobile - N/A

Dataset analysed: imports-85.data from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/autos/imports-85.data">https://archive.ics.uci.edu/ml/machine-learning-databases/autos/imports-85.data</a> (dataset\_file\_url column from Fortydatasets.xlsx) It does not have header.

```
symboling (symbol):
    String format: Error(s) found:

DQI #17 (Non-String Data Type - Consistency):

205 Non-string value(s) at index(es): [(0, 3), (1, 3), (2, 1), (3, 2), (4, 2), (5, 2), (6, 1), (7, 1), (8, 1), (9, 0), ('...', '...'),

(195, -1), (196, -2), (197, -1), (198, -2), (199, -1), (200, -1), (201, -1), (202, -1), (203, -1), (204, -1)] (displaying only the
String range (lexicographical): (-1 : 3)
normalized-losses (loss):
Numerical format: Error(s) found:

DQI #17 (Wrong Data Type - Consistency):

41 Non-numeric value(s) at index(es): [(0, '?'), (1, '?'), (2, '?'), (5, '?'), (7, '?'), (9, '?'), (14, '?'), (15, '?'), (16, '?'), (17, '?'), (127, '?'), (128, '?'), (129, '?'), (130, '?'), (131, '?'), (181, '?'), (189, '?'), (191, '?'), (193, '?')] (displaying only the first and last 10 items)
Range of values: (65.0:256.0).
   All 205 values are correctly categorical.
Categorical format with 22 unique values:
     Category Frequency
         toyota
         nissan
                                   18
           mazda
                                   17
           honda
                                   13
  mitsubishi
                                   13
         subaru
                                   12
  volkswagen
                                   12
         peugot
                                   11
                                    9
           dodge
             audi
      plymouth
            saab
        porsche
                                    5
           isuzu
alfa-romero
```

```
jaguar
                    3
    renault
   mercury
fuel-type (type):
 All 205 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
    gas
               185
 diesel
                20
aspiration:
 All 205 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
    std
   turbo
                37
num-of-doors (four):
 Error(s) found:
DQI #4 (Ambiguous Data - Accuracy, Consistency):
2 Unacceptable value(s) at index(es): [(27, '?'), (63, '?')]
Categorical format with 3 unique values:
Category Frequency
   four
               114
    two
                89
body-style (style):
 All 205 values are correctly categorical.
Categorical format with 5 unique values:
   Category Frequency
     sedan
                    96
 hatchback
                    70
                    25
     wagon
   hardtop
                    8
convertible
                    6
drive-wheels (drive):
 All 205 values are correctly categorical.
Categorical format with 3 unique values:
Category Frequency
     fwd
               120
     rwd
                76
    4wd
                 9
engine-location (location):
 All 205 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
  front
               202
wheel-base (continuous):
 Numerical >=0 format: All 205 values are numerical and greater or equal to 0 in the range (86.6:120.9).
 Numerical >=0 format: All 205 values are numerical and greater or equal to 0 in the range (141.1:208.1).
width:
 Numerical >=0 format: All 205 values are numerical and greater or equal to 0 in the range (60.3:72.3).
 Numerical >=0 format: All 205 values are numerical and greater or equal to 0 in the range (47.8:59.8).
curb-weight (weight):
 Numerical >=0 format: All 205 values are numerical and greater or equal to 0 in the range (1488:4066).
engine-type (type):
 All 205 values are correctly categorical.
Categorical format with 7 unique values:
Category Frequency
    ohc
               148
   ohcf
                15
    ohcv
                13
   dohc
                 12
      1
                12
   rotor
                 4
   dohcv
                 1
num-of-cylinders (four):
 All 205 values are correctly categorical.
Categorical format with 7 unique values:
Category Frequency
    four
```

3

chevrolet

```
six
                  24
    five
                  11
   eight
                   5
                   4
     two
   three
                   1
  twelve
engine-size (size):
  Numerical >=0 format: All 205 values are numerical and greater or equal to 0 in the range (61:326).
fuel-system (system):
  All 205 values are correctly categorical.
Categorical format with 8 unique values:
Category Frequency
    mpfi
                 94
    2bb1
                  66
     idi
                  20
    1bbl
                 11
    spdi
                  9
    4bbl
                  3
     mfi
                  1
    spfi
                  1
bore:
 Numerical >=0 format: Error(s) found:
DQI #17 (Wrong Data Type - Consistency):
 4 Non-numeric value(s) at index(es): [(55, '?'), (56, '?'), (57, '?'), (58, '?')]
Range of values: (2.54:3.94).
stroke:
 Numerical >=0 format: Error(s) found:
DQI #17 (Wrong Data Type - Consistency):
 4 Non-numeric value(s) at index(es): [(55, '?'), (56, '?'), (57, '?'), (58, '?')]
Range of values: (2.07:4.17).
compression-ratio (ratio):
  Numerical format: All 205 values are numerical in the range (7.0:23.0).
horsepower:
  Numerical >=0 format: Error(s) found:
DQI #17 (Wrong Data Type - Consistency):
 2 Non-numeric value(s) at index(es): [(130, '?'), (131, '?')]
Range of values: (48.0:288.0).
peak-rpm (rpm):
  Numerical >=0 format: Error(s) found:
DQI #17 (Wrong Data Type - Consistency):
 2 Non-numeric value(s) at index(es): [(130, '?'), (131, '?')]
Range of values: (4150.0:6600.0).
city-mpg (mpg):
  Numerical >=0 format: All 205 values are numerical and greater or equal to 0 in the range (13:49).
highway-mpg (mpg):
  Numerical >=0 format: All 205 values are numerical and greater or equal to 0 in the range (16:54).
price:
 Numerical format: Error(s) found:
DQI #17 (Wrong Data Type - Consistency):
4 Non-numeric value(s) at index(es): [(9, '?'), (44, '?'), (45, '?'), (129, '?')]
Range of values: (5118.0:45400.0).
Last run on: 2024-02-25 22:05:41
            Alerts 4
                         Reproduction
Overview
  Alerts
   fuel-type is highly imbalanced (53.9%)
   engine-location is highly imbalanced (89.0%)
   num-of-cylinders is highly imbalanced (57.6%)
   symboling has 67 (32.7%) zeros
```

### 30- 9 - Auto MPG - N/A

Dataset analysed: auto-mpg.data from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/auto-mpg/auto-mpg.data">https://archive.ics.uci.edu/ml/machine-learning-databases/auto-mpg/auto-mpg.data</a> (dataset file url column from Fortydatasets.xlsx) It does not have header.



### YData just found High Correlation. It did not find '?', Year and Name Data Quality Issues.

```
Numerical >=0 format: All 398 values are numerical and greater or equal to 0 in the range (9.0:46.6).
cylinders (cylinder):
   Numerical >=0 format: All 398 values are numerical and greater or equal to 0 in the range (3:8).
displacement:
   Numerical >=0 format: All 398 values are numerical and greater or equal to 0 in the range (68.0:455.0).
   Numerical >=0 format: Error(s) found:
DQI #17 (Wrong Data Type - Consistency)
 6 Non-numeric value(s) at index(es): [(32, '?'), (126, '?'), (330, '?'), (336, '?'), (354, '?'), (374, '?')]
Range of values: (46.0:230.0).
weight:
   Numerical >=0 format: All 398 values are numerical and greater or equal to 0 in the range (1613.0:5140.0).
acceleration:
   Numerical format: All 398 values are numerical in the range (8.0:24.8).
model year (year):
 OQI #15 (Domain Violation - Accuracy):
398 Value(s) outside range [1800, 2100] at index(es): [(0, 70), (1, 70), 70), (9, 70), ('...', '...'), (388, 82), (389, 82), (390, 82), (391, 82), 82)] (displaying only the first and last 10 items)
                                                                                                                 (2, 70), (3, 70), (4, 70), (5, 70), (6, 70), (7, 70), (392, 82), (393, 82), (394, 82), (395, 82), (396, 82),
Actual range of values: (70 : 82)
origin:
  All 398 values are correctly categorical.
Categorical format with 3 unique values:
 Category Frequency
                          249
            3
                           79
            2
                           70
car name (name):
Name format: Error(s) found:
DQI #5 (Extraneous Data - Consistency, Uniqueness):
 116 Extraneous data value(s) at index(es): [(1, 'buick skylark 320'), (5, 'ford galaxie 500'), (11, "plymouth 'cuda 340"), (18, 'datsun pl510'), (19, 'volkswagen 1131 deluxe sedan'), (20, 'peugeot 504'), (21, 'audi 100 ls'), (22, 'saab 99e'), (23, 'bmw 2002'), (25, 'ford f250'), ('...', '...'), (354, 'renault 18i'), (357, 'datsun 200sx'), (358, 'mazda 626'), (359, 'peugeot 505s turbo diesel'), (362, 'datsun 810 maxima'), (369, 'chevrolet cavalier 2-door'), (370, 'pontiac j2000 se hatchback'), (385, 'datsun 310 gx'), (391, 'dodge charger 2.2'), (397, 'chevy s-10')] (displaying only the first and last 10 items)
    equency Distribution:
```

```
car name Frequenc
        ford pinto
       amc matador
     ford maverick
    toyota corolla
       amc gremlin
        amc hornet
chevrolet chevette
  chevrolet impala
       peugeot 504
     toyota corona
       volvo 144ea
   volvo 145e (sw)
       volvo 244dl
       volvo 245
volvo 264gl
      volvo diesel
vw dasher (diesel)
```

Category Frequency

Range of Values: (amc ambassador brougham to vw rabbit custom)

# 31- 40 - Flags - N/A

Dataset analysed: flag.data from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/flags/flag.data">https://archive.ics.uci.edu/ml/machine-learning-databases/flags/flag.data</a> (dataset\_file\_url column from Fortydatasets.xlsx) It does not have header.

```
name:
 Name format: All 194 name values are valid.
Frequency Distribution:
          name Frequency
    Afghanistan
       Albania
                        1
       Algeria
                        1
American-Samoa
                        1
       Andorra
                        1
        Angola
      Anguilla
Antigua-Barbuda
     Argentina
     Argentine
       Uruguay
       Vanuatu
  Vatican-City
      Venezuela
       Vietnam
 Western-Samoa
     Yugoslavia
                        1
         Zaire
        Zambia
      Zimbabwe
Range of Values: (Afghanistan to Zimbabwe)
 All 194 values are correctly categorical.
Categorical format with 6 unique values:
Category Frequency
                  52
                  39
        5
        3
                  35
       1
                  31
        6
                  20
zone (1=):
 All 194 values are correctly categorical.
Categorical format with 4 unique values:
Category Frequency
       1
                  91
       4
                  58
        2
                  29
        3
                  16
 Numerical >=0 format: All 194 values are numerical and greater or equal to 0 in the range (0:22402).
 Numerical >=0 format: All 194 values are numerical and greater or equal to 0 in the range (0:1008).
 All 194 values are correctly categorical.
Categorical format with 10 unique values:
Category Frequency
      10
                  46
                  43
       1
                  30
                  21
        3
                  17
        4
                   6
        5
                   4
                   4
       9
                   4
 All 194 values are correctly categorical.
Categorical format with 8 unique values:
```

```
0
                  40
                  36
                  27
                  15
                  8
        4
                  1
bars (number):
 Numerical format: All 194 values are numerical in the range (0:5).
stripes (number):
 Numerical format: All 194 values are numerical in the range (0:14).
colours (number):
 Numerical format: All 194 values are numerical in the range (1:8).
 All 194 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
                153
       0
                 41
green (0 if):
 All 194 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
        0
               103
       1
                 91
blue (0 if):
 All 194 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
                 99
       1
       0
                 95
gold (0 if):
 All 194 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
       a
                103
       1
                 91
white (0 if):
 All 194 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
       0
                 48
black (0 if):
 All 194 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
       0
                142
       1
                 52
orange (0 if):
 All 194 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
       0
                168
       1
                 26
mainhue (colour):
 All 194 values are correctly categorical.
Categorical format with 8 unique values:
Category Frequency
    red
   blue
                40
  green
                31
   white
                22
   gold
                19
  black
                 5
 orange
                 4
  brown
                 2
circles (circle):
 Numerical >=0 format: All 194 values are numerical and greater or equal to 0 in the range (0:4).
crosses (number):
 Numerical format: All 194 values are numerical in the range (0:2).
```

60

1

```
saltires (number):
 Numerical format: All 194 values are numerical in the range (0:1).
quarters (number):
 Numerical format: All 194 values are numerical in the range (0:4).
sunstars (number):
 Numerical format: All 194 values are numerical in the range (0:50).
 All 194 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
              183
       0
triangle (1 if):
 All 194 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency 0 167
       1
                 27
icon (1 if):
 All 194 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
       0
             145
       1
                 49
animate (animation):
 All 194 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
              155
       0
       1
                 39
text (1 if):
 All 194 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
           178
       0
       1
                 16
topleft (colour):
 All 194 values are correctly categorical.
Categorical format with 7 unique values:
Category Frequency
   blue
                43
  white
                41
   green
                32
  black
                12
   gold
                 6
 orange
botright (colour):
 All 194 values are correctly categorical.
Categorical format with 8 unique values:
Category Frequency
    red
   blue
                47
  green
                40
   white
                17
  black
                 9
   gold
                 9
  brown
                 2
 orange
                 1
Last run on: 2024-02-26 16:59:40
```



### Alerts

Age:

bars is highly imbalanced (59.1%)	Imbalance
circles is highly imbalanced (66.4%)	Imbalance
crosses is highly imbalanced (59.9%)	Imbalance
saltires is highly imbalanced (55.4%)	Imbalance
quarters is highly imbalanced (62.2%)	Imbalance
crescent is highly imbalanced (68.6%)	Imbalance
text is highly imbalanced (58.9%)	Imbalance
name has unique values	Unique
area has 34 (17.5%) zeros	Zeros
population has 56 (28.9%) zeros	Zeros
religion has 40 (20.6%) zeros	Zeros
stripes has 110 (56.7%) zeros	Zeros
sunstars has 114 (58.8%) zeros	Zeros

YData did not find anything that we care.

### 32- 336 - Chronic Kidney Disease - N/A

Dataset analysed: @DATA part extracted from chronic\_kidney\_disease\_full.arff that was obtained at Chronic\_Kidney\_Disease.rar from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00336/Chronic\_Kidney\_Disease.rar">https://archive.ics.uci.edu/ml/machine-learning-databases/00336/Chronic\_Kidney\_Disease.rar</a> (dataset\_file\_url column from Fortydatasets.xlsx) It does not have header.

YData Profiling just brought High Correlation and Imbalance alerts, no '?' information. We found dozens of cases with '?'.

```
Age format: Error(s) found:
DQI #17 (Wrong Data Type - Consistency):
9 Non-numeric value(s) at index(es): [(30, '?'), (73, '?'), (112, '?'), (116, '?'), (117, '?'), (169, '?'), (191, '?'), (203, '?'),
Actual range of values: (2.0 : 90.0)
Blood Pressure (blood):
   Numerical format: Error(s) found:
DQI #17 (Wrong Data Type - Consistency):
12 Non-numeric value(s) at index(es): [(7, '?'), (75, '?'), (132, '?'), (138, '?'), (161, '?'), (164, '?'), (185, '?'), (188, '?'), (215, '?'), (293, '?'), (316, '?')]
Range of values: (50.0:180.0).
Specific Gravity (nominal):
   Error(s) found:
DQI #4 (Ambiguous Data - Accuracy, Consistency):

47 Unacceptable value(s) at index(es): [(13, '?'), (17, '?'), (21, '?'), (28, '?'), (30, '?'), (37, '?'), (50, '?'), (57, '?'), (59, '?'), (78, '?'), ('...', '...'), (228, '?'), (231, '?'), (236, '?'), (245, '?'), (268, '?'), (280, '?'), (295, '?'), (322 '?'), (346, '?')] (displaying only the first and last 10 items)

Categorical format with 6 unique values:
Category Frequency
    1.020
                         106
    1,010
                           84
                          81
    1.025
    1.015
     1.005
Albumin (nominal):
   Error(s) found:
DQI #4 (Ambiguous Data - Accuracy, Consistency):

46 Unacceptable value(s) at index(es): [(13, '?'), (17, '?'), (21, '?'), (30, '?'), (37, '?'), (50, '?'), (57, '?'), (59, '?'), (78, '?'), (81, '?'), ('...', '...'), (228, '?'), (231, '?'), (236, '?'), (238, '?'), (245, '?'), (268, '?'), (280, '?'), (295, '?'), (322, '?'), (346, '?')] (displaying only the first and last 10 items)
Categorical format with 7 unique values:
Category Frequency
           0
                         199
                          46
                           44
           3
                           43
           4
           5
Sugar (nominal):
   Error(s) found:
DQI #4 (Ambiguous Data - Accuracy, Consistency):
```

```
49 Unacceptable value(s) at index(es): [(13, '?'), (17, '?'), (21, '?'), (30, '?'), (37, '?'), (50, '?'), (57, '?'), (59, '?'), (78, '?'), (81, '?'), ('...', '...'), (228, '?'), (231, '?'), (236, '?'), (238, '?'), (245, '?'), (268, '?'), (280, '?'), (295, '?'), (322, '?'), (346, '?')] (displaying only the first and last 10 items)
Categorical format with 7 unique values:
Category Frequency
                          290
                            49
            2
                            18
           3
                            14
           1
                            13
           4
                            13
           5
                             3
Red Blood Cells (nominal):
   Error(s) found:
DQI #4 (Ambiguous Data - Accuracy, Consistency):
152 Unacceptable value(s) at index(es): [(0, '?'), (1, '?'), (5, '?'), (6, '?'), (10, '?'), (12, '?'), (13, '?'), (15, '?'), (16, '?'), (17, '?'), ('...', '...'), (245, '?'), (280, '?'), (280, '?'), (290, '?'), (309, '?'), (309, '?'), (322, '?'), (349, '?'), (381, '?')] (displaying only the first and last 10 items)
Categorical format with 3 unique values:
Category Frequency
  normal
                         201
                          152
abnormal.
                            47
Pus Cell (nominal):
  Error(s) found:
DQI #4 (Ambiguous Data - Accuracy, Consistency):
65 Unacceptable value(s) at index(es): [(5, '?'), (13, '?'), (17, '?'), (21, '?'), (28, '?'), (30, '?'), (34, '?'), (37, '?'), (39, '?'), (43, '?'), ('...', '...'), (245, '?'), (280, '?'), (290, '?'), (295, '?'), (309, '?'), (322, '?'), (349, '?'), (350, '?'), (381, '?')] (displaying only the first and last 10 items)
Categorical format with 3 unique values:
Category Frequency
   normal
                         259
abnormal
                            76
                            65
Pus Cell clumps (nominal):
  Error(s) found:
DQI #4 (Ambiguous Data - Accuracy, Consistency):
4 Unacceptable value(s) at index(es): [(290, '?'), (300, '?'), (316, '?'), (328, '?')]
Categorical format with 3 unique values:
  Category Frequency
notpresent
                             354
                               42
   present
                                 4
Bacteria (nominal):
   Error(s) found:
DQI #4 (Ambiguous Data - Accuracy, Consistency):
4 Unacceptable value(s) at index(es): [(290, '?'), (300, '?'), (316, '?'), (328, '?')]
Categorical format with 3 unique values:
   Category Frequency
notpresent
                             374
    present
                               22
Blood Glucose Random (blood):
  Numerical format: Error(s) found:
44 Non-numeric value(s) at index(es): [(1, '?'), (21, '?'), (23, '?'), (24, '?'), (29, '?'), (38, '?'), (41, '?'), (47, '?'), (52, '?'), (54, '?'), ('...', '...'), (209, '?'), (215, '?'), (234, '?'), (276, '?'), (283, '?'), (312, '?'), (315, '?'), (332, '?'), (378, '?')] (displaying only the first and last 10 items)
Range of values: (22.0:490.0).
Blood Urea (blood):
   Numerical format: Error(s) found:
DQI #17 (Wrong Data Type - Consistency):
19 Non-numeric value(s) at index(es): [(23, '?'), (54, '?'), (55, '?'), (64, '?'), (67, '?'), (113, '?'), (134, '?'), (161, '?'), (165, '?'), (209, '?'), (215, '?'), (220, '?'), (276, '?'), (283, '?'), (312, '?'), (315, '?'), (334, '?'), (378, '?')]
Range of values: (1.5:391.0).
Serum Creatinine (creatinine):
Numerical format: Error(s) found:

DQI #17 (Wrong Data Type - Consistency):

17 Non-numeric value(s) at index(es): [(23, '?'), (55, '?'), (64, '?'), (67, '?'), (113, '?'), (161, '?'), (165, '?'), (215, '?'), (216, '?'), (220, '?'), (232, '?'), (276, '?'), (283, '?'), (312, '?'), (315, '?'), (334, '?')]
Range of values: (0.4:76.0).
   Numerical format: Error(s) found:
DQI #17 (Wrong Data Type - Consistency):
87 Non-numeric value(s) at index(es): [(0, '?'), (1, '?'), (2, '?'), (4, '?'), (7, '?'), (8, '?'), (10, '?'), (19, '?'), (23, '?'), (28, '?'), ('...', '...'), (232, '?'), (234, '?'), (235, '?'), (237, '?'), (240, '?'), (283, '?'), (303, '?'), (315, '?'), (363, '?')] (displaying only the first and last 10 items)
Range of values: (4.5:163.0).
Potassium:
   Numerical format: Error(s) found:
```

```
DQI #17 (Wrong Data Type - Consistency):
88 Non-numeric value(s) at index(es): [(0, '?'), (1, '?'), (2, '?'), (4, '?'), (7, '?'), (8, '?'), (10, '?'), (19, '?'), (21, '?') (23, '?'), ('...', '...'), (232, '?'), (234, '?'), (235, '?'), (237, '?'), (240, '?'), (283, '?'), (303, '?'), (315, '?'), (363, '?')] (displaying only the first and last 10 items)
Range of values: (2.5:47.0).
Hemoglobin:
 Numerical format: Error(s) found:
DQI #17 (Wrong Data Type - Consistency):
52 Non-numeric value(s) at index(es): [(23, '?'), (28, '?'), (30, '?'), (34, '?'), (41, '?'), (57, '?'), (60, '?'), (61, '?'), (66, '?'), (67, '?'), ('...', '...'), (230, '?'), (233, '?'), (247, '?'), (273, '?'), (319, '?'), (324, '?'), (328, '?'), (330, '?'), (365, '?')] (displaying only the first and last 10 items)
Range of values: (3.1:17.8).
Packed Cell Volume (volume):
  Numerical >=0 format: Error(s) found:
DQI #17 (Wrong Data Type - Consistency):
71 Non-numeric value(s) at index(es): [(13, '?'), (16, '?'), (17, '?'), (23, '?'), (28, '?'), (30, '?'), (34, '?'), (38, '?'), (41, '?'), (45, '?'), ('...', '...'), (224, '?'), (238, '?'), (232, '?'), (233, '?'), (247, '?'), (273, '?'), (319, '?'), (324, '?'), (365, '?')] (displaying only the first and last 10 items)
Range of values: (9.0:54.0).
White Blood Cell Count (count):
  Numerical >=0 format: Error(s) found:
DQI #17 (Wrong Data Type - Consistency):
106 Non-numeric value(s) at index(es): [(6, '?'), (10, '?'), (13, '?'), (16, '?'), (17, '?'), (23, '?'), (28, '?'), (29, '?'), (30, '?'), (33, '?'), ('...', '...'), (238, '?'), (247, '?'), (273, '?'), (274, '?'), (287, '?'), (302, '?'), (319, '?'), (324, '?'), (330, '?')] (displaying only the first and last 10 items)
Range of values: (2200.0:26400.0).
Red Blood Cell Count (count):
  Numerical >=0 format: Error(s) found:
DQI #17 (Wrong Data Type - Consistency):
131 Non-numeric value(s) at index(es): [(1, '?'), (2, '?'), (6, '?'), (10, '?'), (13, '?'), (16, '?'), (17, '?'), (23, '?'), (28, '?'), (29, '?'), ('...', '...'), (238, '?'), (239, '?'), (247, '?'), (273, '?'), (274, '?'), (287, '?'), (302, '?'), (319, '?'), (324, '?'), (330, '?')] (displaying only the first and last 10 items)
Range of values: (2.1:8.0).
Hypertension (nominal):
Categorical format with 3 unique values:
Category Frequency
       no
                      251
                      147
      yes
                        2
Diabetes Mellitus (nominal):
  Error(s) found:
DQI #4 (Ambiguous Data - Accuracy, Consistency):
2 Unacceptable value(s) at index(es): [(288, '?'), (297, '?')]
Categorical format with 5 unique values:
Category Frequency
       yes
                      134
        no
Coronary Artery Disease (disease):
  Error(s) found:
DQI #4 (Ambiguous Data - Accuracy, Consistency):
 2 Unacceptable value(s) at index(es): [(288, '?'), (297, '?')]
Categorical format with 4 unique values:
Category Frequency
       yes
                        34
        no
        5
Appetite:
  Error(s) found:
DQI #4 (Ambiguous Data - Accuracy, Consistency):
 1 Unacceptable value(s) at index(es): [(294,
Categorical format with 3 unique values:
Category Frequency
     good
                      317
     poor
                       82
                        1
         ?
Pedal Edema (nominal):
  Error(s) found:
DQI #4 (Ambiguous Data - Accuracy, Consistency):
 1 Unacceptable value(s) at index(es): [(294,
Categorical format with 3 unique values:
Category Frequency
       no
```

```
yes
Anemia:
 Error(s) found:
DQI #4 (Ambiguous Data - Accuracy, Consistency):
 1 Unacceptable value(s) at index(es): [(294,
Categorical format with 3 unique values:
Category Frequency
     nο
                339
     yes
                 60
Class:
 All 400 values are correctly categorical.
Categorical format with 3 unique values:
Category Frequency
  notckd
    ckd
```

Last run on: 2024-02-26 20:59:59

## 33- 50 - Image Segmentation - N/A

Dataset analysed: segmentation.data from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/image/segmentation.data">https://archive.ics.uci.edu/ml/machine-learning-databases/image/segmentation.data</a> (dataset\_file\_url column from Fortydatasets.xlsx) It has not only a header but also some comments before the header. Therefore it was necessary to copy the file to local directory to extract extraneous lines.

```
;;; -*- Mode:Common-Lisp; Base:10 -*-
;;; *-* Last-edit: 11/21/90 15:28:30 by Brodley; *-*
 REGION-CENTROID-COL, REGION-CENTROID-ROM, REGION-PIXEL-COUNT, SHORT-LINE-DENSITY-5, SHORT-LINE-DENSITY-2, VEDGE-MEAN, VEDGE-SD, HEDGE-MEAN, HEDGE-SD, INTENSITY-MEAN, RAWRED-MEAN, RAWBLUE-MEAN, RAWGREEN-MEAN, EXGREEN-MEAN, EXGREEN-MEAN, VEDGE-MEAN, EXGREEN-MEAN, RAWGREEN-MEAN, RAWGREEN-MEAN, RAWGREEN-MEAN, EXGREEN-MEAN, EXGREEN-MEAN, VEDGE-MEAN, VEDGE-MEAN, VEDGE-MEAN, VEDGE-SD, HEDGE-MEAN, RAWGREEN-MEAN, VEDGE-MEAN, VEDGE-MEAN, VEDGE-MEAN, RAWGREEN-MEAN, RAWGREEN-MEA
 BRICKFACE, 140, 0.125, 0.9, 0.0, 0.0, 0.277779, 0.06296301, 0.66666675, 0.31111118, 6.185185, 7.333335, 7.6666665, 3.5555556, 3.4444444, 4.4444447, -7.888889, 7.777777, 0.5456349, -1.1218182
 BRICKFACE, 188.0, 133.0, 9,0.0,0.0,0.33333334,0.26666674,0.5,0.077777736,6.6666665,8.333334,7.7777777,3.8888888,5.0,3.333333,-8.333333,8.444445,0.53858024,-0.92481726
BRICKFACE, 105.0,139.0,9,0.0,0.0,0.27777782,0.107407436,0.83333325,0.52222216,6.111111,7.5555553,7.2222223,3.5555556,4.3333333,-7.6666665,7.5555553,0.5326279,-0.96594584
class:
   All 210 values are correctly categorical.
Categorical format with 7 unique values:
  Category
                       Frequency
BRICKFACE
      CEMENT
                                      30
    FOLIAGE
                                      30
        GRASS
                                       30
          PATH
                                       30
            SKY
                                      30
      WINDOW
region-centroid-col (column):
    Numerical format: All 210 values are numerical in the range (1.0:252.0).
region-centroid-row (row):
    Numerical format: All 210 values are numerical in the range (11.0:250.0).
region-pixel-count (count):
    Numerical >=0 format: All 210 values are numerical and greater or equal to 0 in the range (9:9).
short-line-density-5 (density):
    Numerical >=0 format: All 210 values are numerical and greater or equal to 0 in the range (0.0:0.111111111).
short-line-density-2 (density):
    Numerical >=0 format: All 210 values are numerical and greater or equal to 0 in the range (0.0:0.22222222).
vedge-mean (mean):
    Numerical format: All 210 values are numerical in the range (0.0:25.5).
vegde-sd (standard deviation):
    Numerical >=0 format: All 210 values are numerical and greater or equal to 0 in the range (0.0:572.9964).
    Numerical format: All 210 values are numerical in the range (0.0:44.722225).
hedge-sd (standard deviation):
    Numerical >=0 format: All 210 values are numerical and greater or equal to 0 in the range (0.0:1386.3292).
intensity-mean (intensity):
    Numerical format: All 210 values are numerical in the range (0.0:143.44444).
rawred-mean (mean):
    Numerical format: All 210 values are numerical in the range (0.0:136.88889).
    Numerical format: All 210 values are numerical in the range (0.0:150.88889).
rawgreen-mean (mean):
```

```
Numerical format: All 210 values are numerical in the range (0.0:142.55556).
exred-mean (mean):
  Numerical format: All 210 values are numerical in the range (-48.22222:5.7777777).
exblue-mean (mean):
  Numerical format: All 210 values are numerical in the range (-9.666667:78.77778).
exgreen-mean (mean):
  Numerical format: All 210 values are numerical in the range (-30.555555:21.88889).
value-mean (mean):
 Numerical format: All 210 values are numerical in the range (0.0:150.88889).
saturatoin-mean (mean):
 Numerical format: All 210 values are numerical in the range (0.0:1.0).
 Numerical format: All 210 values are numerical in the range (-2.5309503:2.8649306).
Last run on: 2024-02-26 22:03:00
verview
           Alerts 18
                          Reproduction
Alerts
 region-pixel-count has constant value ""
 short-line-density-5 is highly imbalanced (61.1%)
 short-line-density-2 is highly imbalanced (80.4%)
 class is uniformly distributed
                                                                              Uniform
 vedge-mean has 5 (2.4%) zeros
 vegde-sd has 5 (2.4%) zeros
 hedge-mean has 5 (2.4%) zeros
 hedge-sd has 5 (2.4%) zeros
 intensity-mean has 5 (2.4%) zeros
 rawred-mean has 10 (4.8%) zeros
                                                                               Zeros
 rawblue-mean has 5 (2.4%) zeros
 rawgreen-mean has 9 (4.3%) zeros
 exred-mean has 5 (2.4%) zeros
 exblue-mean has 5 (2.4%) zeros
 exgreen-mean has 5 (2.4%) zeros
 value-mean has 5 (2.4%) zeros
 saturatoin-mean has 5 (2.4%) zeros
 hue-mean has 5 (2.4%) zeros
```

YData shows column region-pixel-count with constant value "". We have not checked that. For us it is a numerical with values in the range (9:9).

# region-pixel-count Categorical CONSTANT Distinct 1 9 210

## 34- 162 - Forest Fires – Physical

Dataset analysed: forestfiles.csv from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/forest-fires/forestfires.csv">https://archive.ics.uci.edu/ml/machine-learning-databases/forest-fires/forestfires.csv</a> (dataset file url column from Fortydatasets.xlsx) has header

```
X - x-axis spatial coordinate within the Montesinho park map (x):
Numerical format: All 517 values are numerical in the range (1:9).
Y - y-axis spatial coordinate within the Montesinho park map (y):
Numerical format: All 517 values are numerical in the range (2:9).
month - month of the year (month):
Month format: All 517 month values are valid.
Frequency Distribution:
Month Frequency
```

```
March
     July
 ebruary
  October
    April
 December
  January
     Mav
  Weekday
           Frequency
   Sunday
  Friday
                  85
 Saturday
  Monday
  Tuesday
FFMC - FFMC index from the FWI system (index):
 Numerical format: All 517 values are numerical in the range (18.7:96.2).
DMC - DMC index from the FWI system (index):
 Numerical format: All 517 values are numerical in the range (1.1:291.3).
DC - DC index from the FWI system (index):
 Numerical format: All 517 values are numerical in the range (7.9:860.6).
ISI - ISI index from the FWI system (index):
 Numerical format: All 517 values are numerical in the range (0.0:56.1).
temp - temperature in Celsius degrees (temperature):
 Numerical format: All 517 values are numerical in the range (2.2:33.3).
RH - relative humidity in % (humidity):
 Numerical >=0 format: All 517 values are numerical and greater or equal to 0 in the range (15:100).
wind - wind speed in km (speed):
 Numerical format: All 517 values are numerical in the range (0.4:9.4).
rain - outside rain in mm (rain):
 Numerical >=0 format: All 517 values are numerical and greater or equal to 0 in the range (0.0:6.4).
area - the burned area of the forest (area):
 Numerical >=0 format: All 517 values are numerical and greater or equal to 0 in the range (0.0:1090.84).
           Alerts 3
                         Reproduction
verview
 Alerts
 Dataset has 4 (0.8%) duplicate rows
 rain - outside rain in mm has 509 (98.5%) zeros
  area - the burned area of the forest has 247 (47.8%) zeros
```

YData did not find anything that we care.

# 35- 235 - Individual household electric power consumption – Physical

Dataset analysed: household\_power\_consumption.txt from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00235/household\_power\_consumption.zip">https://archive.ics.uci.edu/ml/machine-learning-databases/00235/household\_power\_consumption.zip</a> (dataset\_file\_url column from Fortydatasets.xlsx) has header and <a href="https://orchive.ics.uci.edu/ml/machine-learning-databases/00235/household\_power\_consumption.zip">https://orchive.ics.uci.edu/ml/machine-learning-databases/00235/household\_power\_consumption.zip</a> (dataset\_file\_url column from Fortydatasets.xlsx) has header and <a href="https://orchive.ics.uci.edu/ml/machine-learning-databases/00235/household\_power\_consumption.zip">https://orchive.ics.uci.edu/ml/machine-learning-databases/00235/household\_power\_consumption.zip</a> (dataset\_file\_url column from Fortydatasets.xlsx) has header and <a href="https://orchive.ics.uci.edu/ml/machine-learning-databases/">https://orchive.ics.uci.edu/ml/machine-learning-databases/</a> (dataset\_file\_url column from Fortydatasets.xlsx) has header and <a href="https://orchive.ics.uci.edu/ml/machine-learning-databases/">https://orchive.ics.uci.edu/ml/machine-learning-databases/</a> (dataset\_file\_url column from Fortydatasets.xlsx) has header and <a href="https://orchive.ics.uci.edu/ml/machine-learning-databases/">https://orchive.ics.uci.edu/ml/machine-learning-databases/</a> (dataset\_file\_url column from Fortydatasets.xlsx) has header and <a href="https://orchive.ics.uci.edu/ml/machine-learning-databases/">https://orchive.ics.uci.edu/ml/machine-learning-databases/</a> (dataset\_file\_url column from Fortydatasets.xlsx) has header databases (databases/">https://orchive.ics.uci.edu/ml/machine-learning-databases/</a> (databases/">https://orchive.ics.uci.edu/ml/machine-learning-databases/</a> (databases/">https://orchive.ics.uci.edu/ml/machine-learning-databases/</a> (databases/">https://orchive.ics.uci.edu/ml/machine-learning-databases/</a> (databases/">https://orchive.ics.uci.edu/ml/machine-learning-databases/</a> (dat

Start of Analysis on: 2024-02-29 14:08:46 date:

Date format: All 2075259 date values are valid in the DDMMYYYY format in the range 16/12/2006 to 26/11/2010.

time:

```
Overview Alerts 7 Reproduction
```

### Alerts

```
global_active_power is an unsupported type, check if it needs cleaning or further analysis
             global_reactive_power is an unsupported type, check if it needs cleaning or further analysis
             voltage is an unsupported type, check if it needs cleaning or further analysis
             global_intensity is an unsupported type, check if it needs cleaning or further analysis
             sub metering 1 is an unsupported type, check if it needs cleaning or further analysis
             sub metering 2 is an unsupported type, check if it needs cleaning or further analysis
             sub metering 3 is an unsupported type, check if it needs cleaning or further analysis
global_active_power (power):
       Numerical format: Error(s) found:
Notifier Let 1 to mat. Error (s) round.

DQI #17 (Wrong Data Type - Consistency):
25979 Non-numeric value(s) at index(es): [(6839, '?'), (6840, '?'), (19724, '?'), (19725, '?'), (41832, '?'), (61909, '?'), (98254, '?'), (98255, '?'), (142588, '?'), (190497, '?'), ('...', '...'), (1990180, '?'), (1990181, '?'), (1990182, '?'), (1990183, '?'), (1990184, '?'), (1990185, '?'), (1990186, '?'), (1990187, '?'), (1990188, '?'), (2027411, '?')] (displaying only the first and last 10 to 10 t
 items)
Range of values: (0.076:11.122).
global_reactive_power (power):
       Numerical format: Error(s) found:
DQI #17 (Wrong Data Type - Consistency):

25979 Non-numeric value(s) at index(es): [(6839, '?'), (6840, '?'), (19724, '?'), (19725, '?'), (41832, '?'), (61909, '?'), (98254, '?'), (98255, '?'), (142588, '?'), (190497, '?'), ('...', '...'), (1990180, '?'), (1990181, '?'), (1990182, '?'), (1990183, '?'), (1990184, '?'), (1990185, '?'), (1990186, '?'), (1990187, '?'), (1990188, '?')] (displaying only the first and last 10 to the first and last 1
items)
Range of values: (0.0:1.39).
voltage:
         Numerical format: Error(s) found:
Numerical Profiles: Profiles: Terror(s) Tourid:

DQI #17 (Wrong Data Type - Consistency):

25979 Non-numeric value(s) at index(es): [(6839, '?'), (6840, '?'), (19724, '?'), (19725, '?'), (41832, '?'), (61909, '?'), (98254, '?'), (198255, '?'), (142588, '?'), (190497, '?'), ('...', '...'), (1990180, '?'), (1990181, '?'), (1990182, '?'), (1990183, '?'), (1990184, '?'), (1990185, '?'), (1990186, '?'), (1990187, '?'), (1990188, '?')] (displaying only the first and last 10
 items)
Range of values: (223.2:254.15).
global intensity (intensity):
       Numerical format: Error(s) found:
NOI #17 (Wrong Data Type - Consistency):

25979 Non-numeric value(s) at index(es): [(6839, '?'), (6840, '?'), (19724, '?'), (19725, '?'), (41832, '?'), (61909, '?'), (98254, '?'), (98255, '?'), (142588, '?'), (190497, '?'), ('...', '...'), (1990180, '?'), (1990181, '?'), (1990182, '?'), (1990183, '?'), (1990184, '?'), (1990185, '?'), (1990186, '?'), (1990187, '?'), (1990188, '?')] (displaying only the first and last 10
 items)
Range of values: (0.2:48.4).
sub_metering_1 (metering):
       Numerical format: Error(s) found:
Note: 15 Not
 items)
Range of values: (0.0:88.0).
 sub_metering_2 (metering):
       Numerical format: Error(s) found:
DQI #17 (Wrong Data Type - Consistency):

25979 Non-numeric value(s) at index(es): [(6839, '?'), (6840, '?'), (19724, '?'), (19725, '?'), (41832, '?'), (61909, '?'), (98254, '?'), (98255, '?'), (142588, '?'), (190497, '?'), ('...', '...'), (1990180, '?'), (1990181, '?'), (1990182, '?'), (1990183, '?'), (1990184, '?'), (1990185, '?'), (1990186, '?'), (1990187, '?'), (1990187, '?')] (displaying only the first and last 10
 items)
Range of values: (0.0:80.0).
sub_metering_3 (metering):
        Numerical format: Error(s) found:
 DQI #1 (Missing Data - Completeness):
Digital (Missing Data - Completeness).

25979 Blank/Empty/Null/NaN value(s) at index(es): [(6839, ''), (6840, ''), (19724, ''), (19725, ''), (41832, ''), (61909, ''), (98254, ''), (98255, ''), (142588, ''), (190497, ''), ('...', '...'), (1990180, ''), (1990181, ''), (1990182, ''), (1990183, ''), (1990184, ''), (1990185, ''), (1990186, ''), (1990187, ''), (1990188, ''), (2027411, '')] (displaying only the first and last 10 items)
Range of values: (0.0:31.0).
```

## 36- 165 - Concrete Compressive Strength – Physical

Dataset analysed: Concrete\_Data.xls from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/concrete/compressive/Concrete\_Data.xls">https://archive.ics.uci.edu/ml/machine-learning-databases/concrete/compressive/Concrete\_Data.xls</a> (dataset\_file\_url column from Fortydatasets.xlsx) has header

```
Numerical >=0 format: All 1030 values are numerical and greater or equal to 0 in the range (102.0:540.0).
  Numerical format: All 1030 values are numerical in the range (0.0:359.4).
Flv Ash (ash):
 Numerical >=0 format: All 1030 values are numerical and greater or equal to 0 in the range (0.0:200.1).
 Numerical format: All 1030 values are numerical in the range (121.75:247.0).
Superplasticizer (quantitative):
  Numerical format: All 1030 values are numerical in the range (0.0:32.2).
Coarse Aggregate (quantitative):
  Numerical format: All 1030 values are numerical in the range (801.0:1145.0).
Fine Aggregate (quantitative):
  Numerical format: All 1030 values are numerical in the range (594.0:992.6).
                        Day (day):
  Day format: All 1030 values are
 ctual range of values: (1 : 365)
Concrete compressive strength -- quantitative -- MPa -- Output Variable (quantitative):
  Numerical format: All 1030 values are numerical in the range (2.331807832:82.5992248).
Last run on: 2024-02-29 14:35:20
Dverview
            Alerts 8
                          Reproduction
 Alerts
  Dataset has 11 (1.1%) duplicate rows
  Age -- quantitative -- Day is highly overall correlated with Concrete compressive
                                                                                    High correlation
  strength -- quantitative -- MPa -- Output Variable
                                                                                    High correlation
  Concrete compressive strength -- quantitative -- MPa -- Output Variable iS
  highly overall correlated with Age -- quantitative -- Day
  Superplasticizer is highly overall correlated with Water
  Water is highly overall correlated with Superplasticizer
  Blast Furnace Slag has 466 (45.2%) zeros
  Fly Ash has 566 (55.0%) zeros
```

Superplasticizer has 379 (36.8%) zeros

YData did not find anything that we care.

138 - Robot Execution Failures — archive.ics.uci.edu/ml/machine-learning-databases/robotfailure-mld/lp1.data

Physical had data set in different lines for the same instance of data. There is no way to read such dataset. See two instances below:

```
normal
                                                                                                                                                              0
                                                                                                         -3
                                                    -1
                                                                               61
                                                                               63
                                                                                                         -2
-3
-3
-3
-3
                                                                                                                                                              0
0
0
0
                                                                               63
63
63
63
                                                                               61
61
                                                                                                         -3
-3
-3
                                                                               64
                                                                               60
normal
                                                                               63
                                                                                                                                                              0
```

```
    0
    -4
    63
    1
    0
    0

    0
    -1
    59
    -2
    0
    -1

    -3
    3
    57
    -8
    -3
    -1

    -1
    3
    70
    -10
    -2
    -1

    0
    -3
    61
    0
    0
    0
    0

    0
    -2
    53
    -1
    -2
    0

    0
    -3
    66
    1
    4
    0

    -3
    3
    58
    -10
    -5
    0

    -1
    -1
    66
    -4
    -2
    0

    0
    1
    66
    -6
    -3
    -1

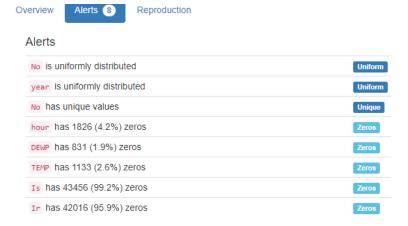
    1
    66
    -6
    -3
    -1
```

52 – Ionosphere - Physical https://archive.ics.uci.edu/ml/machine-learning-databases/ionosphere/ionosphere.data

This dataset does not exhibit names/labels of columns/attributes.

### 37- 381 - Beijing PM2.5 Data - Physical

Dataset analysed: PRSA\_data\_2010.1.1-2014.12.31.csv from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00381/PRSA\_data\_2010.1.1-2014.12.31.csv">https://archive.ics.uci.edu/ml/machine-learning-databases/00381/PRSA\_data\_2010.1.1-2014.12.31.csv</a> (dataset\_file\_url column from Fortydatasets.xlsx) has header



YData found that the column 'No' has unique values. We just found that it is a numerical format with nothing wrong. Besides that we found that the column pm2.5 should have only numerical values but we found over 2000 values with the content 'NA'. They also found that but did not consider an alert.

```
Numerical format: All 43824 values are numerical in the range (1:43824).

year:
Year format: All 43824 values are numerical and valid in the range [1800, 2100]
Actual range of values: (2010 : 2014)

month:
Month format: All 43824 month values are valid.

Frequency Distribution:
Month Frequency
January 3720
March 3720
May 3720
July 3720
August 3720
October 3720
October 3720
December 3720
April 3600
```

```
ay:
Day format: All 43824 values are numerical and valid in the range [1, 366].
ctual range of values: (1 : 31)
```

eptember November

ebruary

3600

Range of values: (0.0:994.0).

Numerical (between 0 and 24) format: All 43824 values are numerical and valid in the range [0, 24]. Actual range of values: (0: 23)

```
pm2.5 (concentration):
    Numerical >=0 format: Error(s) found:
DQI #17 (Wrong Data Type - Consistency):
2067 Non-numeric value(s) at index(es): [(0, 'NA'), (1, 'NA'), (2, 'NA'), (3, 'NA'), (4, 'NA'), (5, 'NA'), (6, 'NA'), (7, 'NA'), (8, 'NA'), (9, 'NA'), ('...', '...'), (43283, 'NA'), (43544, 'NA'), (43545, 'NA'), (43546, 'NA'), (43547, 'NA'), (43548, 'NA'), (43550, 'NA'), (43551, 'NA'), (43552, 'NA')] (displaying only the first and last 10 items)
```

```
DEWP (dew point):
 Numerical format: All 43824 values are numerical in the range (-40:28).
TEMP (temperature):
 Numerical format: All 43824 values are numerical in the range (-19.0:42.0).
 Numerical format: All 43824 values are numerical in the range (991.0:1046.0).
  All 43824 values are correctly categorical.
            format with 4 unique values:
         Frequency
 ategory
              15296
               9387
 Numerical format: All 43824 values are numerical in the range (0.45:585.6).
Ts (hours):
 Numerical >=0 format: All 43824 values are numerical and greater or equal to 0 in the range (0:27).
Ir (hours):
 Numerical >=0 format: All 43824 values are numerical and greater or equal to 0 in the range (0:36).
Last run on: 2024-02-29 15:32:41
```

### 38- 320 - Student Performance - Social

Dataset analysed: student-por.csv from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00320/.student.zip\_old">https://archive.ics.uci.edu/ml/machine-learning-databases/00320/.student.zip\_old</a> (dataset\_file\_url column from Fortydatasets.xlsx), but, as the original file ends with zip\_old, the file had to be renamed to only .zip to extract the final dataset <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00320/.student.zip\_old">https://archive.ics.uci.edu/ml/machine-learning-databases/00320/.student.zip\_old</a> (dataset\_file\_url column from Fortydatasets.xlsx), but, as the original file ends with zip\_old, the file had to be renamed to only .zip to extract the final dataset

```
School - student's school (binary):
Binary format: All 649 binary values are valid.

Frequency Distribution:
Value Frequency
GP 423
MS 226

Sex - student's sex (binary):
Binary format: All 649 binary values are valid.

Frequency Distribution:
Value Frequency
F 383
M 266

age - student's age (age):
Age format: All 649 values are numerical and valid in the range [0, 130].

Actual range of values: (15 : 22)

address - student's home address type (binary):
Binary format: All 649 binary values are valid.

Frequency Distribution:
Value Frequency
U 452
R 197

famsize - family size (binary):
Binary format: All 649 binary values are valid.

Frequency Distribution:
Value Frequency
GT3 457
LE3 192

Pstatus - parent's cohabitation status (binary):
Binary format: All 649 binary values are valid.
```

```
Numerical format: All 649 values are numerical in the range (0:4).
Fedu - father's education (numeric):
  Numerical format: All 649 values are numerical in the range (0:4).
Mjob - mother's job (job):
  All 649 values are correctly categorical.
Categorical format with 5 unique values:
Category Frequency
  other
                258
services
                136
 at home
                135
 teacher
  health
                 48
Fjob - father's job (job):
 All 649 values are correctly categorical.
Categorical format with 5 unique values:
Category Frequency
  other
                367
services
                181
 at home
                 42
 teacher
                 36
  health
                 23
reason - reason to choose this school (reason):
 All 649 values are correctly categorical.
Categorical format with 4 unique values:
  Category Frequency
    course
                  285
     home
                  149
reputation
                  143
                   72
    other
guardian - student's guardian (guardian):
 All 649 values are correctly categorical.
Categorical format with 3 unique values:
Category Frequency
  mother
                455
  father
                153
  other
                 41
traveltime - home to school travel time (numeric):
  Numerical format: All 649 values are numerical in the range (1:4).
studytime - weekly study time (numeric):
  Numerical format: All 649 values are numerical in the range (1:4).
failures - number of past class failures (number):
  Numerical format: All 649 values are numerical in the range (0:3).
 choolsup - extra educational support (binary):
Binary format: All 649 binary values are valid
  equency Distribution:
 alue Frequency
  no
             581
famsup - family educational support (binary):
  Binary format: All 649 binary values are valid.
Frequency Distribution:
Value Frequency
 yes
             398
  no
             251
paid - extra paid classes within the course subject (binary):
 Binary format: All 649 binary values are valid.
Frequency Distribution:
Value Frequency
  no
              39
 yes
activities - extra-curricular activities (binary):
 Binary format: All 649 binary values are valid.
Frequency Distribution:
Value Frequency
             334
  no
 yes
             315
```

Medu - mother's education (numeric):

```
nursery - attended nursery school (binary):
  Binary format: All 649 binary values are valid.
Frequency Distribution:
Value Frequency
 yes
   no
             128
higher - wants to take higher education (binary):
  Binary format: All 649 binary values are valid.
Frequency Distribution:
Value Frequency
             580
 yes
  no
internet - Internet access at home (binary):
  Binary format: All 649 binary values are valid.
Frequency Distribution:
Value Frequency
             498
 yes
             151
  no
romantic - with a romantic relationship (binary):
 Binary format: All 649 binary values are valid.
Frequency Distribution:
Value Frequency
             410
  no
 yes
             239
famrel - quality of family relationships (numeric):
 Numerical format: All 649 values are numerical in the range (1:5).
freetime - free time after school (numeric):
 Numerical format: All 649 values are numerical in the range (1:5).
goout - going out with friends (numeric):
  Numerical format: All 649 values are numerical in the range (1:5).
Dalc - workday alcohol consumption (alcohol):
  Numerical format: All 649 values are numerical in the range (1:5).
Walc - weekend alcohol consumption (alcohol):
  Numerical format: All 649 values are numerical in the range (1:5).
health - current health status (numeric):
 Numerical format: All 649 values are numerical in the range (1:5).
absences - number of school absences (number):
 Numerical format: All 649 values are numerical in the range (0:32).
G1 - first period grade (grade):
 Numerical >=0 format: All 649 values are numerical and greater or equal to 0 in the range (0:19).
G2 - second period grade (grade):
 Numerical >=0 format: All 649 values are numerical and greater or equal to 0 in the range (0:19).
G3 - final grade (grade):
 Numerical >=0 format: All 649 values are numerical and greater or equal to 0 in the range (0:19).
Last run on: 2024-02-29 21:26:03
            Alerts 7
Overview
                         Reproduction
 Alerts
  failures - number of past class failures is highly imbalanced (59.9%)
  schoolsup - extra educational support is highly imbalanced (51.6%)
  paid - extra paid classes within the course subject is highly imbalanced (67.2%)
  higher - wants to take higher education is highly imbalanced (51.1%)
  absences - number of school absences has 244 (37.6%) zeros
  G2 - second period grade has 7 (1.1%) zeros
```

Zeros

YData did not find anything that we care.

G3 - final grade has 15 (2.3%) zeros

## 39- 275 - Bike Sharing Dataset - Social

Dataset analysed: hour.csv from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/00275/Bike-Sharing-Dataset.zip">https://archive.ics.uci.edu/ml/machine-learning-databases/00275/Bike-Sharing-Dataset.zip</a> (dataset\_file\_url column from Fortydatasets.xlsx) has header instant:

```
instant:
 Numerical >=0 format: All 17379 values are numerical and greater or equal to 0 in the range (1:17379).
  All 17379 values are correctly categorical.
Categorical format with 4 unique values:
Category Frequency
                4496
                4409
        1
                4242
        4
                4232
  All 17379 values are correctly categorical.
Categorical format with 2 unique values:
Category Frequency
                8734
        0
                8645
mnth (month):
 Month format: All 17379 month values are valid.
Frequency Distribution:
    Month Frequency
     May
                1488
     July
                1488
December
                1483
   August
                1475
                1473
   March
  October
                1451
                1440
    June
    April
                1437
September
                1437
                1437
 November
  January
 February
                1341
  Numerical (between 0 and 24)
  tual range of values: (0 : 23)
        Frequency
 Value
  Weekday
  Friday
 Thursday
   Sunday
  Tuesday
  Binary format: All 17379 binary values are valid.
         Distribution:
        Frequency
            1186
weathersit (weathersit):
 All 17379 values are correctly categorical.
```

Categorical format with 4 unique values:
Category Frequency
1 11413
2 4544
3 1419

4

```
Normalized format: All
atemp (normalized):
 Normalized format: All 17379 values are numerical and valid in the range [0, 1].
Actual range of values: (0.0 : 1.0)
hum (normalized):
 Normalized format: All 17379 values are numerical and valid in the range [0, 1].
Actual range of values: (0.0 : 1.0)
windspeed (normalized):
 Normalized format: All 17379 values are numerical and valid in the range [0, 1].
Actual range of values: (0.0 : 0.8507)
casual (count):
 Numerical >=0 format: All 17379 values are numerical and greater or equal to 0 in the range (0:367).
 Numerical`>=0 format: All 17379 values are numerical and greater or equal to 0 in the range (0:886).
cnt (count):
 Numerical >=0 format: All 17379 values are numerical and greater or equal to 0 in the range (1:977)
Last run on: 2024-03-01 12:49:47.
Overview
            Alerts (19)
                          Reproduction
 Alerts
```

atemp is highly overall correlated with casual and 2 other fields	High correlation
casual is highly overall correlated with atemp and 3 other fields	High correlation
cnt is highly overall correlated with casual and 2 other fields	High correlation
hr is highly overall correlated with cnt and 1 other fields	High correlation
instant is highly overall correlated with season and 1 other fields	High correlation
mnth is highly overall correlated with season	High correlation
registered is highly overall correlated with casual and 2 other fields	High correlation
season is highly overall correlated with atemp and 2 other fields	High correlation
temp is highly overall correlated with atemp and 1 other fields	High correlation
weekday is highly overall correlated with workingday	High correlation
workingday is highly overall correlated with weekday	High correlation
yr is highly overall correlated with instant	High correlation
holiday is highly imbalanced (81.2%)	Imbalance
instant is uniformly distributed	Uniform
instant has unique values	Unique
hr has 726 (4.2%) zeros	Zeros
weekday has 2502 (14.4%) zeros	Zeros
windspeed has 2180 (12.5%) zeros	Zeros
casual has 1581 (9.1%) zeros	Zeros

YData shows the column 'instant' as Unique. We have not checked that.

### 40- 13 - Balloons - Social

Dataset analysed: adult+stretch.data from <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/balloons/adult+stretch.data">https://archive.ics.uci.edu/ml/machine-learning-databases/balloons/adult+stretch.data</a> (dataset\_file\_url column from Fortydatasets.xlsx) <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/balloons/adult+stretch.data">https://archive.ics.uci.edu/ml/machine-learning-databases/balloons/adult+stretch.data</a> (dataset\_file\_url column from Fortydatasets.xlsx) <a href="https://archive.ics.uci.edu/ml/machine-learning-databases/balloons/adult+stretch.data">https://archive.ics.uci.edu/ml/machine-learning-databases/balloons/adult+stretch.data</a> (dataset\_file\_url column from Fortydatasets.xlsx)

```
Color:
All 19 values are correctly categorical.

Categorical format with 2 unique values:
Category Frequency
PURPLE 10
YELLOW 9

size (large):
All 19 values are correctly categorical.
```



YData did not find anything that we care.