

## Abstract

- We propose a novel technique based on a combination of association rule learning and conformal prediction in its Mondrian form.
- As an application, we use data about (anonymised) business customers of a multinational energy company, Centrica plc.
- There are multiple fields in Centrica's SAP database indicating if a customer is an Industrial Corporation or Small/Medium-sized Enterprise. We consider these as labels.
- Often these labels are incorrect or inconsistent across the SAP system, which has a financial impact on the company. The aim of this work is to use machine learning to identify potential errors and propose corrections.

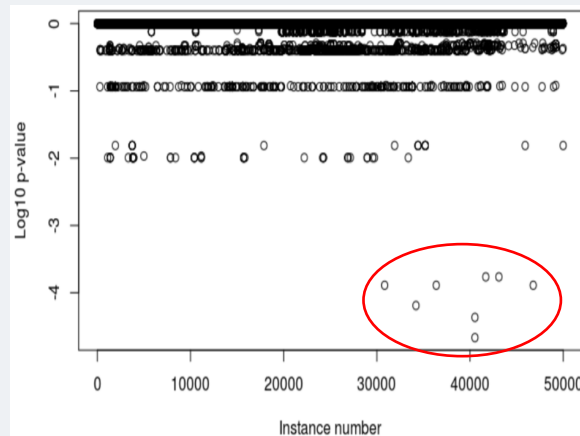
## Data details and preprocessing

- Original dataset: 2.4m rows.  
Sample: First 50,000 records (from ~37,000 companies). All relevant features are categorical.
- Features:** *abwck*, *formkey*, *z-edl-inv*, *zahlkond*, *sparte*, *zz-mbd-flag*, *zz-mba-flag*, *zz-mb-flag*
- These give information on the customers e.g. payment terms, type of invoice, type of energy (gas vs electricity), microbusiness or not
- After one-hot-encoding to convert the data to a binary representation, the number of features rises from 8 to 43
- Labels:** *kofiz-sd*, *kofiz*, *bpkind*, *zccustomer-type*  
These indicate whether an example is assigned as INC or SME, but sometimes contradict each other:
  - kofiz-sd*: 7% INC, 93% SME
  - kofiz*: 6% INC, 94% SME
  - bpkind*: 5% INC, 95% SME
  - zccustomer-type*: 7% INC, 93% SME
- We engineer a new feature indicating 'overall label', which is the average of the four labels (Let INC=1, SME=0). Rounding is used where necessary to make it binary.

## Methodology

- Our approach:
  - Come up with a set of rules relating features to labels.
  - Identify which data examples break these rules the most often, i.e. the most non-conforming examples, using conformal prediction.
- We consider only rules of the following types:
  - (+) IF *i*th feature = 1 AND *j*th feature = 1 THEN label=1
  - (-) IF *i*th feature = 1 AND *j*th feature = 1 THEN label=0
- A rule is considered valid if the following conditions are met:
  - there are at least 20 supporting examples of the features being found together
  - average label over the supporting examples is above 0.8 for (+)-type rules, or below 0.05 for (-)-type rules
- The **non-conformity score (NCS)** of a data instance is: *the number of rules that the instance breaks by having the wrong label*
- The **p-value** of a data instance #*i* is: *the number of data instances with the same (rounded) label as #*i* but with a greater NCS than #*i**; i.e. If the p-value is close to 0, the label of that instance is likely to be wrong.

## Results



- Most anomalous examples circled in red – worthy of investigation

## Interpretation of results

Fragment of investigation report for an individual anomaly:

```
[1] "anomalous example number"
[1] 1107
[1] "p-value"
[1] 0.0101567
[1] "the features in original format"
      abwck      formkey      z_edl_inv
[1] FALSE "Z BI INDIVIDUAL INVOICE"      "Z"
      zahlkond      sparte      zz_mbd_flag
[1] "ZD02"      "02"      "X"
      zz_mba_flag      zz_mb_flag
[1] "X"      "X"
[1] "the features in binary format"
[1] 1 0 1 0 1 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 1 0 1 1 0
[1] "the labels in original format"
      kofiz_sd      kofiz      bpkind zccustomer_type
[1] "02"      "02"      "ZSME"      "MJ"
[1] "the average label (1=INC, 0=SME)"
[1] 0.75
[1] "the non-conformity score"
[1] 53
[1] "the rules broken by the example"
[1] "IF"      "abwck"      "="      "FALSE"      "AND"      "abwck"      "="
[2] "FALSE"      "THEN Y=SME"      "abwck"      "-"
[3] "IF"      "abwck"      "abwck"      "-"
[4] "FALSE"      "AND"      "formkey"      "-"
[7] "-"      "Z BI INDIVIDUAL INVOICE"      "THEN Y=SME"
[1] "IF"      "abwck"      "="      "FALSE"      "AND"      "z_edl_inv"      "-"
[2] "Z"      "THEN Y=CMR"
```

- This example has 3 out of 4 labels as INC, but the algorithm suggests that it is actually an SME. This has been validated.

## Conclusions and Future Work

- We have developed a novel technique combining conformal prediction and association rule mining to detect anomalies and applied it to find possible errors in the database of a large corporation.
- Future directions:
  - Further develop NCM function, such as elimination of the parameters and making p-values more sensitive.
  - Add a prediction step after anomaly detection:
    - To check whether the alternative label is anomalous as well (indicated by a low CP p-value) → shows features may be unreliable
    - To get a probabilistic (Venn-ABERS) prediction of the label.
  - Use the recently developed **probabilistic input** version of conformal framework, for a deeper analysis of contradictions between labels in the input data.

## Contact information and acknowledgements

Corresponding authors:  
Iliia Nouretdinov, [i.r.nouretdinov@rhul.ac.uk](mailto:i.r.nouretdinov@rhul.ac.uk)  
James Gammerman, [james.Gammerman@centrica.co.uk](mailto:james.Gammerman@centrica.co.uk)

Work funded by Centrica plc.