# Conformal Association Rule Mining (CARM)

Ilia Nouretdinov, James Gammerman,
Centre for Reliable Machine Learning,
Royal Holloway, University of London, Egham, Surrey, UK

Royal Holloway, University of London

*i.r.nouretdinov@rhul.ac.uk*

*jgammerman@gmail.com*

September 14, 2023

# Introduction

- ► CARM is an integration of two frameworks: *conformal prediction* and *association rule mining*.

- ► It enables detection of errors within a set of binary labels, with the usual CP guarantees on validity.

- ► XAI angle

- ► As an extension, we analyse the errors using *probabilistic prediction* to suggest corrections.

# Origin story: COPA 2019 poster

# Association Rule Mining

- A rule-based ML method for discovering relationships between items in a dataset.
- E.g for supermarkets: {bread} → {butter}
- In general terms: *{antecedent}* → *{consequent}*
- In a ML setting, an example rule might be:

  IF feature $F = a$ THEN feature $G = b$

  $G$ could alternatively refer to a label.
- **Implication**: Given a set of rules, errors can be identified as deviations from rules that hold true for most of the data

# ARM concepts

- **Support**: the proportion of examples where the rule holds.

- **Confidence**: the conditional probability of the rule's consequent given its antecedent.

- Most common way to generate rules: *apriori algorithm*
  - Requires user-defined thresholds for minimum support and confidence

# Apriori algorithm: lack of statistical rigor

- ▶ No notion of statistical significance

- ▶ Arbitrary choice of thresholds can lead to spurious and missed rules

- ▶ Solution: use a different approach to ARM...

# Our approach to ARM

- ▶ As per Hamalainen et al. (2009)
- ▶ Idea: subject any possible rule to a binomial test
- ▶ Each example is treated as an independent Bernoulli trial, whose outcome is either 1 or 0
    - ▶ 1 means the rule $(F = a) \to (G = b)$ occurs
    - ▶ 0 means the rule doesn't occur
- ▶ $H_0$: no association between antecedent and consequent
- ▶ Interpretation of $p$-values: the 'weight' of a given rule
    - ▶ a small value provides evidence against $H_0$

# Non-conformity measure

- ▶ So we have...
  - ▶ A complete pool of rules with associated weights $p_r$
  - ▶ A set of examples $Z$ which either comply with those rules or not
- ▶ So for a given example $z \in Z$, we have an expectation for its label based on each rule $r$ (assuming the antecedent of $r$ is true for $z$)
- ▶ If the expected label doesn't match the observed label, then $z$ is an exception to that rule, i.e. non-conforming.
- ▶ Non-conformity score $\alpha$ for $z$:
  - ▶ Initialise $\alpha$ at 0, then add $-\log p_r$ for every broken rule
  - ▶ accumulates degree of non-conformity

# Converting NCM to conformal *p*-values

- ▶ NB. Different kind of *p*-values!
- ▶ **Input**: The data sequence $Z' = ((x_1, y_1), \ldots, (x_N, y_N))$, a non-conformity measure $\mathcal{A}$, and a threshold $\varepsilon$.
- ▶ **Non-conformity scores**: For each data point $i = 1, \ldots, N$, compute:

  $$\alpha_i = \mathcal{A}\left((x_i, y_i), Z' \setminus \{(x_i, y_i)\}\right)$$

- ▶ *p*-**value Calculation**: Determine $p_N$ using:

  $$p_N = \frac{\left|\{j : y_j = y_N, \alpha_j \geq \alpha_N\}\right|}{\left|\{j : y_j = y_N\}\right|}$$

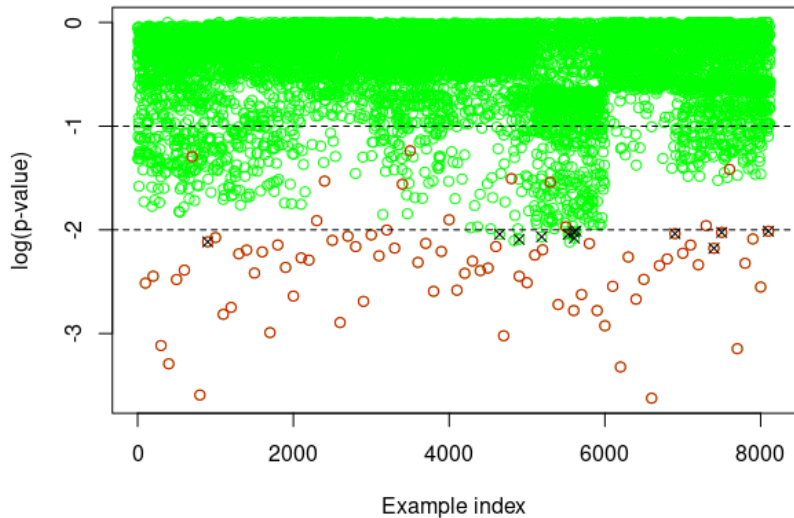- ▶ **Error Detection**: If $p_N < \varepsilon$, flag the example $z_N$ as an error.

# Validity

- All conformal predictors automatically come with the property of *validity*, i.e. the probability of incorrectly rejecting $H_0$ is at most $\varepsilon$.

- In our setting of error detection, $\varepsilon$ becomes an upper bound on the probability of a false positive (false alarm).
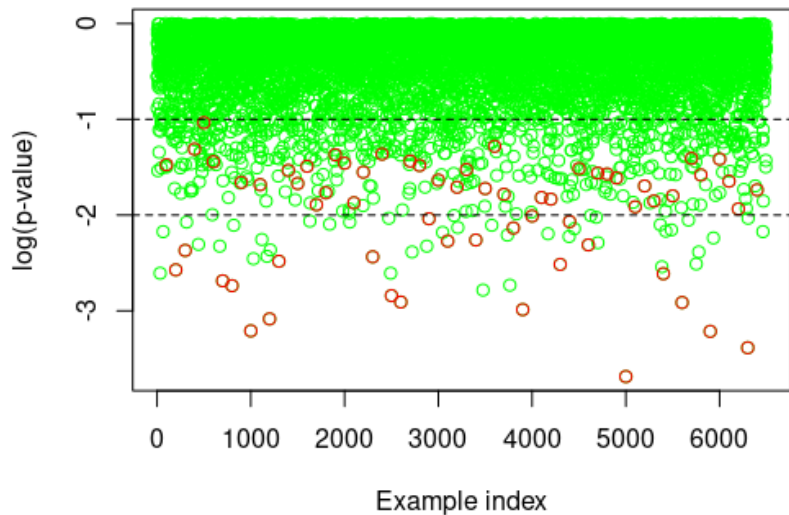
# Datasets & pre-processing

1. **Mushrooms**. Label: edible or posionous
2. **Wine Quality.** Label: red or white
3. **Adult Income**. Label: salary greater/less than $50K$

▶ Introduced known errors into the labels for 1% of examples
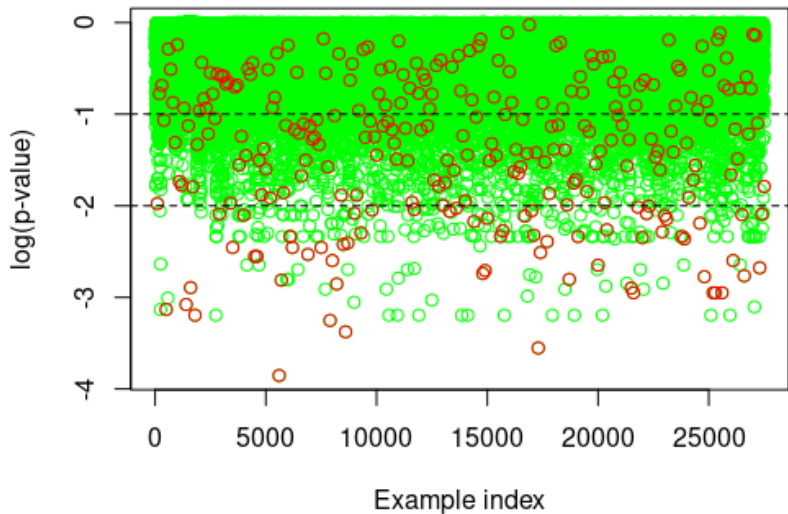▶ Remaining features converted to binary form by one-hot encoding

# CARM *p*-values (log scale) - Mushrooms

# CARM *p*-values (log scale) - Wine

# CARM *p*-values (log scale) - Adult

# Evaluation of CARM

Precision at three *p*-value thresholds

| Dataset | $p < 0.1\%$ | $p < 1\%$ | $p < 10\%$ |
| --- | --- | --- | --- |
| Mushrooms | 100% | 87% | 10% |
| Wine | 100% | 37% | 10% |
| Adult Income | 29% | 23% | 6% |

Recall at three *p*-value thresholds

| Dataset | $p < 0.1\%$ | $p < 1\%$ | $p < 10\%$ |
| --- | --- | --- | --- |
| Mushrooms | 9% | 87% | 100% |
| Wine | 8% | 37% | 100% |
| Adult Income | 2% | 23% | 62% |

# Example report

| Index in data | 6600 | |
|---|---|---|
| **Observed label** | 1 (= edible) | |
| *p*-value | 0.000237 | |
| **Broken rules** | (listed if significant at the threshold $10^{-9}$) | **Confidence** |
| *One-feature rules* | IF cap-shape = k THEN Y = 0 | 0.72 |
| | IF cap-surface = y THEN Y = 0 | 0.54 |
| | IF bruises = f THEN Y = 0 | 0.69 |
| | IF odor = y THEN Y = 0 | 0.99 |
| | IF gill-spacing = c THEN Y = 0 | 0.56 |
| | IF gill-size = n THEN Y = 0 | 0.88 |
| | IF gill-color = b THEN Y = 0 | 0.99 |
| | IF stalk-root = ? THEN Y = 0 | 0.71 |
| | IF stalk-surface-above-ring = k THEN Y = 0 | 0.93 |
| | IF stalk-surface-below-ring = k THEN Y = 0 | 0.93 |
| | IF stalk-color-above-ring = p THEN Y = 0 | 0.69 |
| | IF ring-type = e THEN Y = 0 | 0.63 |
| | IF spore-print-color = w THEN Y = 0 | 0.75 |
| | IF population = v THEN Y = 0 | 0.70 |
| | IF habitat = p THEN Y = 0 | 0.88 |
| *Two-feature rules* | IF gill-attachment = f AND ring-number = o THEN Y = 0 | 0.47 |
| | IF veil-color = w AND ring-number = o THEN Y = 0 | 0.47 |

# Probabilistic Prediction

- ► When a data example is identified as an error, it is important to conduct a thorough investigation that includes also a suggested correction for it.
- ► For this purpose, we can employ the *Venn-ABERS (VA) framework*.

- ► VA works in same assumptions as CP.
- ► Like CP, VA framework can be linked to an underlying method as well.
- ► The difference is in output: where CP produces *p-values*, VA outputs *lower* and *upper probabilities*.

# CARM & VA results for 10 most suspicious examples

| Index | Label. | p-value | VA probs | VA prediction | True label | Comment |
|-------|--------|---------|----------|---------------|------------|---------|
| 6600 | 1 | 0.00024 | 0.0074–0.0078 | 0 | 0 | 1 corrected to 0 |
| 800 | 0 | 0.00026 | 0.990–0.991 | 1 | 1 | 0 corrected to 1 |
| 6200 | 1 | 0.00048 | 0.0074 – 0.0078 | 0 | 0 | 1 corrected to 0 |
| 400 | 0 | 0.00051 | 0.990 – 0.991 | 1 | 1 | 0 corrected to 1 |
| 7700 | 1 | 0.00071 | 0.0074 – 0.0078 | 0 | 0 | 1 corrected to 0 |
| 300 | 0 | 0.00077 | 0.990 – 0.991 | 1 | 1 | 0 corrected to 1 |
| 4700 | 1 | 0.00095 | 0.00780 – 0.98 | 0 | 0 | 1 corrected to 0 |
| 1700 | 0 | 0.001 | 0.990 – 0.991 | 1 | 1 | 0 corrected to 1 |
| 6000 | 1 | 0.0012 | 0.00737 – 0.00780 | 0 | 0 | 1 corrected to 0 |
| 2600 | 0 | 0.0013 | 0.990 – 0.991 | 1 | 1 | 0 corrected to 1 |

▶ *p*-values for alternative labels are all close to 1, i.e. no longer suspicious.

▶ VA predictions indicate high confidence in alternative labels.

# Conclusions and further work

- ▶ Demonstrated integration of modified ARM with the CP framework for *explainable* error detection in data labels.
- ▶ Validity property limits false alarms during error detection.
- ▶ Association rules enhances interpretability and serve as basis for probabilistic analysis using Venn-ABERS

Further work:

- ▶ Exploration of more complex rules.
- ▶ Extension of the methodology for multi-class labels
- ▶ Extension of methodology to correcting features, not just labels