# pRSL: Interpretable Multi–label Stacking by Learning Probabilistic Rules

**Michael Kirchhof**[1]    **Lena Schmid**[1]    **Christopher Reining**[2]    **Michael ten Hompel**[2]    **Markus Pauly**[1]

[1]Department of Statistics, TU Dortmund University, Dortmund, Germany
[2]Chair of Material Handling and Warehousing, TU Dortmund University, Dortmund, Germany

## Abstract

A key task in multi–label classification is modeling the structure between the involved classes. Modeling this structure by probabilistic and interpretable means enables application in a broad variety of tasks such as zero–shot learning or learning from incomplete data. In this paper, we present the probabilistic rule stacking learner (pRSL) which uses probabilistic propositional logic rules and belief propagation to combine the predictions of several underlying classifiers. We derive algorithms for exact and approximate inference and learning, and show that pRSL reaches state–of–the–art performance on various benchmark datasets.

In the process, we introduce a novel multicategorical generalization of the noisy–or gate. Additionally, we report simulation results on the quality of loopy belief propagation algorithms for approximate inference in bipartite noisy–or networks.

## 1 INTRODUCTION

Classifiers that predict a single class have been studied and improved in detail in the past years. However, recent applications in complex real–world environments make it necessary to respect a whole set of different classes and labels [Reining et al., 2019] [Raies and Bajic, 2018]. The success of multi–label classifiers has shown that exploiting the structure between those classes is key to exceeding the performance of isolated classifiers [Pakrashi et al., 2016].

As an example, take the human activity recognition task portrayed in Figure 1 [Niemann et al., 2020] where we want to detect the actions of a warehouse worker. Three embedded machine learning sensors provide probabilistic information about the worker's stance and hand position and a first guess on whether the worker is performing overhead work. Intuitively, we would combine those beliefs to wrap our heads



Figure 1: Human Activity Recognition: The three embedded ML sensors tell us: $P(\text{overhead work}) = 0.5$, $P(\text{standing}) = 0.95$, $P(\text{hands high}) = 0.7$. Given the rules: *if standing and hands high then often overhead work*, *if normal work then almost always hands centered or low*, *not standing iff legs moving*, what is the new probability $P(\text{overhead work})$ that the worker does overhead work?

around the situation the worker is in. To decide how likely the worker is working overhead, we may use rules that tell us whether some beliefs contradict or support other beliefs. We can easily formulate such rules like the three exemplary ones given in Figure 1. However, the challenging question is how to mathematically apply those non–deterministic rules to the uncertain input beliefs to calculate updated beliefs.

Our model, pRSL, approaches this multi–label classification problem via Bayesian belief propagation along probabilistically generalized propositional logic rules. Thereby, it stays interpretable and allows for a wide variety of relations between the classes. Besides expert–given rules, it can learn rules from data that achieve competitive performance.

This paper is structured as follows: In Section 2, we discuss existing approaches to multi–label classification and belief combination. In Section 3, we define pRSL and explain it along the above example. We propose algorithms for exact and approximate inference, learning rules from possibly incomplete data, and suitable regularization. Section 4 concerns the inference methods' approximation quality and benchmarks of pRSL's against the state–of–the–art. The paper closes with a discussion in Section 5.

## 2 RELATED WORK

This section gives an overview of recent approaches in multi–label learning with an emphasis on models that are stacked on top of ground learners. Such models take the ground learners' probabilistic estimates as input and refine them via modeling the structure between the classes in different ways: We group them into attribute–class, knowledge graph, Bayesian network, and probabilistic rule based methods.

**Attribute–Class:** Attribute–class approaches emerged from classical "flat" classification [Lampert et al., 2013] to give a better insight into the reasoning and utilize knowledge of previously learned classes. They add a layer of semantically interpretable attributes in–between the raw inputs and the output classes. In Atzmon and Chechik [2018] this layer is connected to the classes via and–or formulas describing each class while grouping similar attributes. Liu et al. [2020] model the inter–class structure as a graph where classes sharing similar attributes are close to each other. Attribute–class approaches allow solving zero–shot problems [Xian et al., 2018] in which a new class has to be recognized that has never been shown by example, but only described by its attributes.

**Knowledge Graphs:** Other approaches do not split labels into attributes and classes but utilize knowledge graphs instead, thus allowing for broader relations between the labels. Wu et al. [2018] use a graph that embodies hierarchical and co–occurrence relations between labels which allows learning even when not all ground truth labels are provided. Lee et al. [2018] and Liu et al. [2020] apply a form of belief propagation where initial beliefs on labels are shared between one another to obtain a joint solution. Knowledge graphs are easy to interpret and applied to large scale problems, but often require ground truth graphs and allow only a few types of heuristic relations between the labels.

**Bayesian Networks:** Bayesian networks allow extending these relations to the probabilistic setting. While they have been applied as stand–alone approaches [Wang and Li, 2013], they experience a recent interest as stacking approaches on top of ground learners [Chen et al., 2020]. For example, Shen et al. [2018] employ a Bayesian network as an ensemble learner that combines beliefs of several CNNs. A limitation of a Bayesian network is that its directed graph must be acyclic, thus raising design choices that can be difficult outside obvious causal relations.

**Probabilistic Rules:** When discussing probabilistic rule learners, a distinction has to be made regarding the goal. Rule mining [Mencía and Janssen, 2016] [Pham and Aksoy, 1995] generally evolves around analytical insight, while in multi–label classification rules are learned to improve classification accuracy. Although crisp rules have been explored to define classes in programming APIs [Krupka et al., 2017], probabilistic rules take uncertainties into account by not completely ruling out contradicting labels, but only weighting them down as in Ding et al. [2015]. Moreover, Rapp et al. [2020] recently proposed a boosting–based method for finding soft if–then rules. Using propositional logic as the starting point for multi–label classification allows formulating complex rules but loses some of the interpretability of the aforementioned graph models.

Our approach borrows principles from all previously described methods: Our model may be stacked onto any ground learners for initial beliefs on the labels or no classifier may be provided to enable zero–shot classification. We use Bayesian networks as underlying framework due to their probabilistic interpretation and the ability to perform fast belief propagation. However, labels are not connected directly to one another but only via rules in a bipartite manner to ensure acyclicity. Rules can take the form of any propositional logic expression and are extended to the probabilistic setting to allow for vast relations between the labels. Like in knowledge graphs, prior knowledge may be incorporated as expert–given rules, but rules can also be learned from data.

## 3 METHODS

### 3.1 BAYESIAN NETWORKS

A Bayesian network [Pearl, 1988] is a compact representation of the joint distribution of $N$ random variables $X_1, \ldots, X_N$. In our case they are categorical. Each $X_n$ stochastically depends only on a subset of other variables denoted $\mathrm{Pa}(X_n)$ as specified in their conditional probability tables $P(X_n|\mathrm{Pa}(X_n)), n = 1, \ldots, N$. These dependencies form a directed acyclic graph, where each node is a random variable and each edge shows a dependence. Observations are fed into the network by updating priors or providing hard evidence on the states some variables take. The new information is then propagated through the network to update the beliefs on all variables. In the sequel, we use Bayesian networks as the basis for our novel stacking learner as they feature a transparent structure and probabilistic foundation.

### 3.2 PROBABILISTIC RULE STACKING LEARNER

The probabilistic rule stacking learner (pRSL) models the multi–label stacking problem using three kinds of nodes: Classifiers, labels, and rules. The model structure is visualized in Figure 2 and described in the following paragraphs.

Classifiers $C_j, j = 1, \ldots, J$, are plugged into the root of pRSL. Each $C_j$ is an arbitrary machine learning model that accepts some input data $x_j$ and returns a vector of probabilistic estimates $P(C_j|x_j) := (P(C_j = 1|x_j), \ldots, P(C_j = M(j)|x_j))$ on its $M(j)$ respective cate-
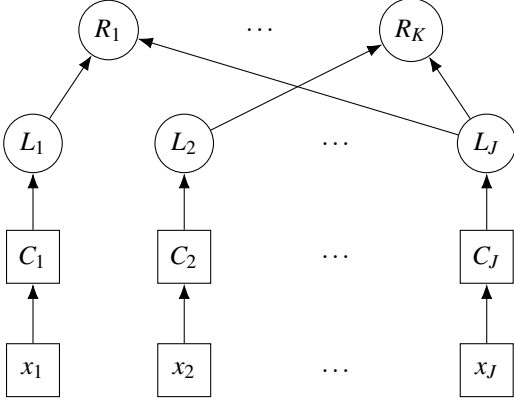
Figure 2: Graphical structure of pRSL. Squares indicate external elements plugged into the framework and circles indicate internal parts of pRSL.

gories[1] $m = 1, \ldots, M(j)$. Classifiers are not restricted to predict the same categories or use the same input data. Each classifier is connected to a label node $L_j$. The label nodes contain the same $M(j)$ categories as their connected classifier nodes, but recalibrate the predictions of their classifiers [Culakova et al., 2020]. Any function $d_j : [0,1]^{M(j)} \to [0,1]^{M(j)}$ can be used for calibration, but it can also be the identity in case the classifier already returns predictions that are not over– or underconfident. The used calibrators are detailed in the experiment sections below. Latent or zero–shot labels are connected to dummy classifiers that always output the labels' priors. As notation, $\boldsymbol{L} := (L_1, \ldots, L_J)$ denotes the vector of all labels, $\mathscr{L} := \bigotimes_{j=1}^{J} \{1, \ldots, M(j)\}$ the sample space of all possible categories $\boldsymbol{L}$ can take, and $\boldsymbol{\ell} := (\ell_1, \ldots, \ell_J) \in \mathscr{L}$ is a particular vector of categories.

The main part of pRSL are the rule nodes $R_k, k = 1, \ldots, K$. Each rule node may be connected to any selection of label nodes[2] and formulates a propositional logic formula $\varphi_k$ that holds between those labels with a probability $p_k$. The rules' conditional probability table generalizes truth tables:

$$P(R_k = 1 | \boldsymbol{L} = \boldsymbol{\ell}) = \begin{cases} p_k & , \boldsymbol{\ell} \models \varphi_k \\ 1 - p_k & , \boldsymbol{\ell} \not\models \varphi_k \end{cases}, \quad (1)$$

where $\models$ means that $\boldsymbol{\ell} \in \mathscr{L}$ logically fulfills the formula $\varphi_k$.

To perform inference, first classifiers receive their inputs $\boldsymbol{x} := (x_1, \ldots, x_J)$. Their estimations $P(C_j | x_j), j = 1, \ldots, J$, are then used as priors for the classifier nodes. Then the evidence $\boldsymbol{R} = \boldsymbol{1}$ with $\boldsymbol{R} = (R_1, \ldots, R_K)$ and $\boldsymbol{1} = (1, \ldots, 1)$ indicating that all rules are true is handed over to the network. This conditioning opens paths between the labels to communicate and update their priors using belief propagation as described in Section 3.5. Once the propagation

---

[1] We use the terms classes and labels interchangeably. They contain an arbitrary number of mutually exclusive categories.

[2] To simplify notation, rules depend on all $J$ possible label nodes in the following formulas.

has reached an equilibrium state, we obtain the updated beliefs $P(L_j | \boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x}), j = 1, \ldots, J$, with $P(L_j | \boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x}) := (P(L_j = 1 | \boldsymbol{R} = \boldsymbol{1}, x_j), \ldots, P(L_j = M(j) | \boldsymbol{R} = \boldsymbol{1}, x_j))$.

### 3.3 EXAMPLE

To better understand how pRSL combines its input beliefs with the given rules, we continue with a numerical example.

In logistics, there is a rising trend of using human activity recognition to observe workers' actions with sensors and machine learning models [Reining et al., 2019]. Suppose we have three sensors, each analyzed by a machine learning model: A camera $C_1$ that can distinguish walking movement $w$ from work at a normal height $n$ and overhead work $o$. Additionally, workers wear sensor shoes $C_2$ that can distinguish standing $s$ from gait $g$, and a wristband $C_3$ that observes high $h$, centered $c$ or low $l$ hand height. All of these classifiers have imperfect accuracy, but we assume for this example's simplicity that they are well calibrated so that we can use $C$ and $L$ interchangeably. We want to connect $C_1, C_2$, and $C_3$ in order to better detect the potentially harmful overhead work. For that, we use three rules:

$$\varphi_1 : (h \wedge s) \to o \; (p_1 = 0.8)$$

expresses that in many cases when workers stand and have their arms up, they do overhead work.

$$\varphi_2 : n \to (c \vee l) \; (p_2 = 0.9)$$

means that during normal work the workers' hands are almost always at a centered or low height. Lastly,

$$\varphi_3 : w \leftrightarrow g \; (p_3 = 1)$$

means the camera's walking category and the sensor shoes' gait category semantically mean the same activity.

Now, consider the camera is unsure whether the worker performs overhead work or normal work, the sensor shoes are quite certain the person is standing and the wristband indicates the worker's hands are probably high:

$$P(C_1 | x_1) = (w = 0.1, n = 0.4, o = 0.5),$$
$$P(C_2 | x_2) = (g = 0.05, s = 0.95),$$
$$P(C_3 | x_3) = (h = 0.5, c = 0.3, l = 0.2).$$

Table 1 shows how pRSL uses the given rules to re–weight which activities seem likely and which contradict the rules and thus are less likely. For each possible $\ell$ the priors are multiplied with the probabilities of the rules given $\ell$. The normalized result is reported in the last column. From the table, we conclude that $\boldsymbol{\ell} = (o, s, h)$ is most likely. After marginalizing over each label, we obtain the updated beliefs

$$P(L_1 | \boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x}) = (w = 0.01, n = 0.05, o = 0.94),$$
$$P(L_2 | \boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x}) = (g = 0.01, s = 0.99),$$
$$P(L_3 | \boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x}) = (h = 0.48, c = 0.31, l = 0.21).$$

Table 1: Calculation for the pRSL Example described in Section 3.3.

| $\ell_1$ | $\ell_2$ | $\ell_3$ | $P(\boldsymbol{L}=\boldsymbol{\ell}\|x)$ | $P(R_1=1\|\boldsymbol{L}=\boldsymbol{\ell})$ | $P(R_2=1\|\boldsymbol{L}=\boldsymbol{\ell})$ | $P(R_3=1\|\boldsymbol{L}=\boldsymbol{\ell})$ | $P(\boldsymbol{L}=\boldsymbol{\ell}\|\boldsymbol{R}=\boldsymbol{1},x)$ |
|---|---|---|---|---|---|---|---|
| $w$ | $s$ | $h$ | $0.1 \cdot 0.95 \cdot 0.5$ | 0.2 | 0.9 | 0 | 0 |
| $n$ | $s$ | $h$ | $0.4 \cdot 0.95 \cdot 0.5$ | 0.2 | 0.1 | 1 | 0.0078 |
| $o$ | $s$ | $h$ | $0.5 \cdot 0.95 \cdot 0.5$ | 0.8 | 0.9 | 1 | 0.3517 |
| $w$ | $g$ | $l$ | $0.1 \cdot 0.05 \cdot 0.2$ | 0.8 | 0.9 | 1 | 0.0015 |
| $n$ | $s$ | $c$ | $0.4 \cdot 0.95 \cdot 0.3$ | 0.8 | 0.9 | 1 | 0.1688 |
| $o$ | $s$ | $c$ | $0.5 \cdot 0.95 \cdot 0.3$ | 0.8 | 0.9 | 1 | 0.2110 |
| … | … | … | … | … | … | … | … |

So, after combining the information of all three sensors we can be certain that the worker is performing overhead work.

This example shows how pRSL applies principles of probabilistic logic to multi–label classification and performs stacking of $C_1, C_2$, and $C_3$ while remaining interpretable. Obviously, there are more sophisticated algorithms to obtain the updated beliefs. They are portrayed hereafter.

### 3.4 MULTICATEGORICAL NOISY–OR

While propositional logic allows for a broad variety of statements, inference and learning are correspondingly complex. However, many formulas can be written in the form of an implication $\varphi = (A \wedge B) \to (C \vee D)$, where both the body and the head may contain multiple labels. This form can be equivalently expressed as a disjunction $\varphi \equiv \neg A \vee \neg B \vee C \vee D$. The noisy–or gate [Pearl, 1988] is a parametrized distribution that extends the classical logic disjunction to the probabilistic setting. It is widely adapted in Bayesian networks [Ji et al., 2020] as it allows for faster inference and, as we will later show, efficient learning.

In Section 5.2 in the appendix, we extend the noisy–or gate to our case where the output is binary, but the labels acting as inputs may be multicategorical. For a rule $R_k$, each label $m$ in each connected node $L_j$ has an inhibition probability $q_{jm}^k$ that gives the probability that the noisy–or node is not activated even though the label is active. The rule's probability previously defined in (1) is now softened:

$$P(R_k=1|\boldsymbol{L}=\boldsymbol{\ell}) = 1 - \prod_{j=1}^{J} q_{j\ell_j}^k . \tag{2}$$

### 3.5 INFERENCE

Inference on pRSL is performed using belief propagation [Pearl, 1988]. Note that we may be interested in two kinds of queries: Either the marginal distribution of each label

$$P(L_j|\boldsymbol{R}=\boldsymbol{1},\boldsymbol{x}), j=1,\dots,J, \tag{3}$$

or the most probable explanation (MPE)

$$\operatorname*{argmax}_{\boldsymbol{\ell} \in \mathscr{L}} P(\boldsymbol{L}=\boldsymbol{\ell}|\boldsymbol{R}=\boldsymbol{1},\boldsymbol{x}), \tag{4}$$

which gives the joint setup of labels that has the highest likelihood. Dembczynski et al. [2010] have shown that the two queries are Bayes–optimal for different loss functions.

We implement pRSL as Bayesian network in R 3.6.3 using the gRain [Højsgaard, 2012] package that provides exact inference for marginal and MPE queries. This implementation allows processing both noisy–or and arbitrary propositional logic based rules, but is NP–complete even for bipartite noisy–or networks [Cooper, 1990]. To circumvent this by utilizing the noisy–or structure, we further implemented approximate algorithms for larger datasets. These loopy belief propagation [Murphy et al., 1999] algorithms are based on the exact sum–product and max–product algorithms for acyclic networks introduced by Pearl [1988] and can approximate marginal and MPE queries. They achieve a runtime complexity of $O(J \cdot J_0 \cdot K)$ for marginal and $O(J \cdot J_0 \cdot K + K \cdot 2^{J_0})$ for MPE queries, where $J_0$ is the maximum number of labels connected to a rule. An experiment on their approximation quality is found in Section 4.1.

### 3.6 LEARNING

All parts of pRSL can either be expert–given or learned from data. We assume the classifiers $C_1, \dots, C_J$ to be trained in advance and refer to Xia et al. [2020] for calibrating the classifier outputs. Thus, we focus on learning the noisy–or rules determined by $\boldsymbol{q} = (q_{11}^1, \dots, q_{JM(J)}^K)$. Our objective is to improve the predictive performance

$$\operatorname*{argmax}_{\boldsymbol{q}} \log(P(\boldsymbol{L}=\boldsymbol{\ell}^*|\boldsymbol{R}=\boldsymbol{1},\boldsymbol{x})), \tag{5}$$

where $\boldsymbol{\ell}^*$ are the true categories of a given observation with classifier input data $\boldsymbol{x}$.

To find the optimal rules, we apply batchwise ADAM optimization [Kingma and Ba, 2015]. In Section 5.1.1 in the appendix, we build partial derivatives of (5) by $q_{jm}^k$ for all rules simultaneously. The resulting gradients can be transformed into expressions that depend only on

$$P(L_j=m|R_k=0, R_v=1, v \neq k, \boldsymbol{x}), \tag{6}$$

$$P(R_k=1|R_v=1, v \neq k, \boldsymbol{x}), \tag{7}$$

and particular values of $\boldsymbol{q}$, so that for each rule's gradient we need to perform only two marginal queries on the pRSL, thus avoiding the higher costs of MPE queries.

We can also perform gradient descent when not all true labels are known such as in zero–shot problems or incomplete datasets. The set of known labels $\{\boldsymbol{L}' = \boldsymbol{\ell}'\}$ may be different for each observation. Then, the optimization goal is

$$\arg\max_{\boldsymbol{q}} \log(P(\boldsymbol{L}' = \boldsymbol{l}'|\boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x})). \qquad (8)$$

We derive corresponding gradients in Section 5.1.2 in the appendix, which is even possible for such $q_{jm}^k$ that belong to labels $L_j \notin \boldsymbol{L}'$ with missing ground–truth. The gradients consist of (6), (7), specific values of $\boldsymbol{q}$, and additionally

$$P(L_j = m|\boldsymbol{L}' = \boldsymbol{l}', R_k = 0, R_v = 1, v \neq k, \boldsymbol{x}) \text{ and} \quad (9)$$
$$P(R_k = 1|\boldsymbol{L}' = \boldsymbol{l}', R_v = 1, v \neq k, \boldsymbol{x}), \qquad (10)$$

i.e. two additional marginal queries are required per rule.

Overall, the number of parameters $\boldsymbol{q}$ to learn is linear in the number of rules and labels, and one step in the gradient descent can be performed in $O(K^2 \cdot J \cdot J_0)$ runtime when using approximate queries to compute the gradients.

## 3.7 REGULARIZATION

The maximum number of labels connected to each rule $J_0$ plays a major rule in keeping computations bearable and rules interpretable. Thus, regularization should favor $q_{jm}^k = 1$ for all inhibition probabilities related to one label $L_j$, so that it can be disconnected from rule $R_k$. In consequence, we do not regularize the $q_{jm}^k$ directly, but their soft minimum [Cook, 2011] per label $L_j$ and rule $R_k$ given by

$$s(k,j) = -\frac{1}{\alpha} \log \left( \sum_{m=1}^{M(j)} \exp(-\alpha q_{jm}^k) \right), \qquad (11)$$

where $\alpha$ controls the hardness of the soft minimum.

Regularization happens twofold: On the one hand, we apply a hard regularizer that allows only the $J_0$ label nodes $L_j$ with the lowest $s(k,j)$ to be connected to each rule $R_k$ and sets all other inhibition probabilities to 1, thus transposing $\boldsymbol{q}$ to close but computationally efficient positions during the gradient descent. On the other hand, there is a soft regularizer to push $\boldsymbol{q}$ towards $\boldsymbol{1}$. This soft regularizer is obtained via the duality between regularizers and Bayesian priors [Murphy, 2012]. Regarding $s(k,j)$ as a random variable $S$, we assume $S \sim \text{Beta}(\beta_1, \beta_2)$, where $0 < \beta_1, \beta_2 < 1$ are selected so that

$$P(S < 0.1) = \gamma_0 \text{ and} \qquad (12)$$
$$P(S > 0.9) = \gamma_1, \qquad (13)$$

with $\gamma_0 < \gamma_1$. This prior places a high mass on 1 and a smaller mass on 0, so that most labels are removed from

the rule to lower $J_0$ and the ones used in the rule are pushed towards crisp inhibition probabilities for easier interpretability. The corresponding regularizer penalty to be added to the optimization goal (5) is

$$R(\boldsymbol{q}) = \frac{\lambda^t}{\eta} ((\beta_1 - 1) \sum_{k=1}^{K} \sum_{j=1}^{J} \log(s(k,j) + \varepsilon) + \qquad (14)$$
$$(\beta_2 - 1) \sum_{k=1}^{K} \sum_{j=1}^{J} \log(1 - s(k,j) + \varepsilon)),$$

where $\eta$ is the normalization constant of the $\text{Beta}(\beta_1, \beta_2)$ distribution, $\varepsilon$ prevents dividing by zero when differentiating and $\lambda < 1$ decreases the influence of the regularization penalty with each iteration $t$ of the gradient descent. A benefit of deriving the regularizer from a prior is that the regularizer strength is given by the normalization constant $\eta$ and does not have to be tuned as a hyperparameter. In simulations not detailed here, $\alpha = 20$, $J_0 = 5$, $\varepsilon = 10^{-4}$, and $\lambda = 0.98$ proved to be good default values.

An R implementation of all above methods is found online[3].

## 4 EXPERIMENTS

### 4.1 APPROXIMATION QUALITY

Though it has no quality bound, Murphy et al. [1999] have observed that the loopy belief propagation used for pRSL's approximate queries may work well in the bipartite noisy–or setting. However, the exact conditions are still under research [Weller and Jebara, 2013]. As the scalability of pRSL relies on the approximation quality, we compare exact to approximated marginal and MPE queries in various simulated pRSL models below. The code to reproduce the simulations and further analyses can be accessed online[4].

pRSL models of three sizes were generated, each replicated 10 times. The smallest included 5 labels and 5 rules, a medium one 10 labels and 10 rules, and the largest 30 labels and 30 rules, after which exact inference became incomputable due to its exponentially rising runtime and memory usage. The sampling procedure is detailed in Section 5.3 in the appendix. In short, labels comprised a random amount of $2 - 4$ categories. Noisy–or rule related $2 - 5$ random categories with random inhibition probabilities. Finally, 100 classifier observations were simulated by drawing $\text{Dir}(1)$ distributed random variables for each label.

Marginal queries were approximated well, with correlation coefficients of $r_s = 0.995$, $r_m = 0.996$, and $r_l = 0.995$ between the exact and approximate queries averaged across all replications in the small ($r_s$), medium ($r_m$), and large ($r_l$) models. In terms of scalability, no strong reduction in quality can be seen in the rising model sizes.
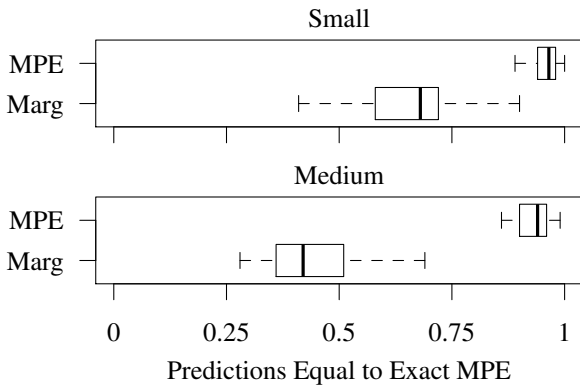
---

[3] https://github.com/mkirchhof/rsl
[4] https://github.com/mkirchhof/rslSim

Figure 3: Quality of Approximation of Exact MPEs by Approximate MPEs and Approximate Marginals.

Table 2: Overview of Benchmark Datasets. Cardinality gives the Mean Number of Positive Labels per Observation. Density is the Cardinality Normalized by the Number of Labels.

| Dataset | Domain | Labels | Cardinality | Density |
|---------|--------|--------|-------------|---------|
| Emotions | Music | 6 | 1.869 | 0.311 |
| Yeast | Genes | 14 | 4.237 | 0.303 |
| Birds | Sound | 19 | 1.014 | 0.053 |
| Medical | Diseases | 45 | 1.245 | 0.028 |
| Enron | Emails | 53 | 3.378 | 0.064 |
| Mediamill | Newscasts | 101 | 4.376 | 0.043 |

For MPE queries, the comparison could only be conducted on the small and medium datasets, as the naive implementation of exact MPE queries in the `gRain` package ruled out the large dataset. Besides approximating MPE queries by the aforementioned loopy belief propagation, they were additionally approximated in Naive–Bayes manner by combining the label–wise approximate marginal queries. Figure 3 shows how often the approximate MPE queries matched their exact counterparts. It shows that approximate MPE queries clearly outperform the label–wise approximation with a median of 96.5% versus 68% correctly resembled MPE queries on the small and 94% versus 42% on the medium datasets. This underlines that questions regarding MPEs are not well answered by combining marginal queries, but need distinguished joint–label algorithms.

In summary, loopy belief propagation showed a nearly perfect approximation quality for marginal and a good approximation quality for MPE queries. The difference between the queries might be due to the higher sensitivity to small inaccuracies of the underlying max–product algorithm for MPE queries as opposed to the sum–product algorithm for marginal queries that can marginalize them out.

## 4.2 PERFORMANCE BENCHMARK

### 4.2.1 Benchmark Datasets

In the following experiment, pRSL was compared to the state–of–the–art on six established multi–label benchmark datasets obtained from Mulan[5]. Each has a different domain and complexity as summarized in Table 2. The three smaller datasets – emotions, yeast, and birds – were split into a 10–fold cross–validation with an 8/1/1 split for train–validation–test. A 5–fold cross–validation with a 3/1/1 split was used for the three bigger datasets. pRSL used exact queries for former and approximate queries for the latter three datasets.

[5]http://mulan.sourceforge.net/datasets-mlc.html

### 4.2.2 Comparison Methods

Probabilistic random forests [Malley et al., 2012] were used as binary relevance (BR) learner to transform the raw inputs into probabilistic estimates for each label individually. We compared pRSL to three methods that further process these initial beliefs: Two benchmark models that are based on black box decisions (NN and MLWSE) and one recent approach (BOOMER) similar in nature to pRSL.

The first method is a neural network (NN) with two hidden layers that have twice the number of labels as neurons. It uses label–wise cross–entropy as loss function and utilizes the validation data for early stopping. Second, MLWSE [Xia et al., 2020] optimizes a quadratic loss by linearly combining labels based on pairwise correlations. Third, we compare with a method already mentioned in Section 2: BOOMER [Rapp et al., 2020] is a recently published soft–rule–based approach that optimizes a joint label cross–entropy loss. To maintain a fair comparison, its three hyperparameters for shrinkage, regularization strength, and number of rules were tuned via grid–search on the validation data. Last, pRSL was applied where the validation data served to find an optimal number of rules. No form of calibration was applied. Experiments could not be conducted on Ding et al. [2015], which is similar to pRSL, due to unavailable implementation.

### 4.2.3 Evaluation Measures

As shown by Dembczyński et al. [2012], measuring the percentage of misclassified observations label–wise, called Hamming loss, does not suffice in multi–label classification. Hence, following their proofs, we additionally measure the percentage of observations where all labels are correctly classified, which lays focus on the MPE estimates. To complement these crisp–decision focused metrics with a metric that judges the quality of the returned probability estimates, we measure the log–likelihood of the returned beliefs, which is the only local proper scoring rule [Parmigiani and Inoue, 2009]. Precisely, we used the median log–likelihood of the joint labels or the individual labels in case the former was not available for a learner.

Table 3: Results of the $k$–Fold Crossvalidations. Mean $\pm$ Standard Deviation Between Folds. Best Result in Bold.

| | Emotions | Yeast | Birds | Medical | Enron | Mediamill |
|---|---|---|---|---|---|---|
| | | | Joint Accuracy (higher = better) | | | |
| BR | $0.321 \pm 0.005$ | $0.174 \pm 0.018$ | $0.510 \pm 0.004$ | $0.386 \pm 0.022$ | $0.114 \pm 0.026$ | $0.148 \pm 0.001$ |
| NN | $0.309 \pm 0.075$ | $0.193 \pm 0.018$ | $0.540 \pm 0.051$ | $0.613 \pm 0.032$ | $0.120 \pm 0.020$ | $0.179 \pm 0.003$ |
| MLWSE | $0.327 \pm 0.058$ | $0.201 \pm 0.013$ | $\mathbf{0.541 \pm 0.041}$ | $0.640 \pm 0.025$ | $0.124 \pm 0.022$ | $0.151 \pm 0.001$ |
| BOOMER | $0.326 \pm 0.074$ | $0.224 \pm 0.022$ | $0.527 \pm 0.024$ | $\mathbf{0.654 \pm 0.024}$ | $0.137 \pm 0.026$ | $\mathbf{0.189 \pm 0.004}$ |
| pRSL | $\mathbf{0.348 \pm 0.067}$ | $\mathbf{0.236 \pm 0.015}$ | $0.507 \pm 0.032$ | $0.491 \pm 0.031$ | $\mathbf{0.153 \pm 0.020}$ | $0.149 \pm 0.002$ |
| | | Joint log–Likelihood (higher = better) | | | Label–wise log–Likelihood (higher = better) | |
| BR | $-2.386 \pm 0.145$ | $-5.772 \pm 0.198$ | $-2.448 \pm 0.160$ | $-1.584 \pm 0.091$ | $-6.376 \pm 0.182$ | $-6.614 \pm 0.062$ |
| NN | $-2.105 \pm 0.267$ | $-5.659 \pm 0.231$ | $-1.547 \pm 0.409$ | $\mathbf{-0.750 \pm 0.133}$ | $-6.476 \pm 0.172$ | $\mathbf{-6.016 \pm 0.038}$ |
| MLWSE | $-2.140 \pm 0.202$ | $-5.512 \pm 0.227$ | $\mathbf{-1.502 \pm 0.167}$ | $-1.125 \pm 0.110$ | $\mathbf{-6.018 \pm 0.104}$ | $-6.371 \pm 0.052$ |
| BOOMER | unavailable | unavailable | unavailable | unavailable | unavailable | unavailable |
| pRSL | $\mathbf{-1.839 \pm 0.273}$ | $\mathbf{-3.592 \pm 0.085}$ | $-2.458 \pm 0.156$ | $-1.565 \pm 0.120$ | $-6.479 \pm 0.242$ | $-6.532 \pm 0.061$ |
| | | | Label–wise Hamming Loss (lower = better) | | | |
| BR | $\mathbf{0.179 \pm 0.019}$ | $0.188 \pm 0.003$ | $0.042 \pm 0.003$ | $0.017 \pm 0.001$ | $0.045 \pm 0.001$ | $0.027 \pm 0.000$ |
| NN | $0.183 \pm 0.024$ | $0.191 \pm 0.005$ | $0.039 \pm 0.003$ | $0.012 \pm 0.001$ | $0.046 \pm 0.001$ | $\mathbf{0.025 \pm 0.000}$ |
| MLWSE | $0.183 \pm 0.024$ | $\mathbf{0.186 \pm 0.005}$ | $\mathbf{0.037 \pm 0.002}$ | $0.011 \pm 0.001$ | $\mathbf{0.044 \pm 0.001}$ | $0.026 \pm 0.000$ |
| BOOMER | $0.183 \pm 0.024$ | $0.188 \pm 0.005$ | $0.041 \pm 0.003$ | $\mathbf{0.011 \pm 0.001}$ | $0.046 \pm 0.001$ | $0.026 \pm 0.000$ |
| pRSL | $0.182 \pm 0.022$ | $0.190 \pm 0.005$ | $0.043 \pm 0.002$ | $0.015 \pm 0.001$ | $0.046 \pm 0.001$ | $0.027 \pm 0.000$ |

### 4.2.4 Results

All code required to reproduce the experiments can be found online[6]. The results are reported in Table 3. As BOOMER returns no probabilistic estimates, its log–likelihood could not be computed. In general, the examined multi–label algorithms outperformed the BR baseline across all datasets and metrics, except the hamming loss on emotions. Interestingly, on birds and medical, BR is outperformed by a larger margin even in terms of hamming loss which is a single label metric that can be optimized without modeling multi–label dependencies [Dembczyński et al., 2012].

Each of the four multi–label algorithms showed strengths on different datasets and metrics. Except on birds, no algorithm performs best across all metrics on a fixed dataset or across all datasets on a fixed metric. In particular, MLWSE, the only linear model, performed best on birds across all metrics and three datasets in terms of hamming loss, plus a nearly indistinguishable performance with BOOMER on medical. The two soft rule based algorithms BOOMER and pRSL performed on a nearly indistinguishable level on four datasets when taking the estimation uncertainty into account, except on medical and mediamill.

pRSL showed the best performance among the benchmarked state–of–the–art methods on three out of six datasets regarding accuracy and two of the three datasets where the joint log–likelihood could be estimated. It does not show im-

provements on the label–wise metrics on any dataset. This may be a consequence of its internal joint–label loss function. It also increased the accuracy against the BR input beliefs on two out of the three big datasets where it relied on approximate queries. This further reinforces the claim on approximation quality made in Section 4.1. On a side note, pRSL showed no signs of overfitting (see Table 4 in the appendix). Two hypotheses for pRSL's worse performance on birds, mediamill, and partially on medical, can be named:

The first is these datasets' low density as seen in Table 2. This issue is related to the problem of imbalanced classes, which is a common pitfall in multi–label classification, especially on these particular datasets [Zhang et al., 2020b].

The second hypothesis is that the initial beliefs of the underlying random forests that pRSL relied on were possibly inadequate. This is motivated by the fact that pRSL is the only benchmarked learner that treats the inputs as probabilities, whereas the other learners see them as scores. Additionally, on each of the three aforementioned datasets, BR was outperformed by its competitors in terms of hamming loss by a considerable margin. To investigate this, Figure 4 shows the calibration of the random forests' probability estimates, that is the relative amount of positive labels stratified by the probability estimates given on them. The dashed bisector visualizes the ideal case in which the classifier is neither under– nor overconfident. It can be seen that the random forest is underconfident for birds, medical and slightly for mediamill. To further assess the hypotheses, we calibrated the random forests' probability estimates using Platt scaling

---

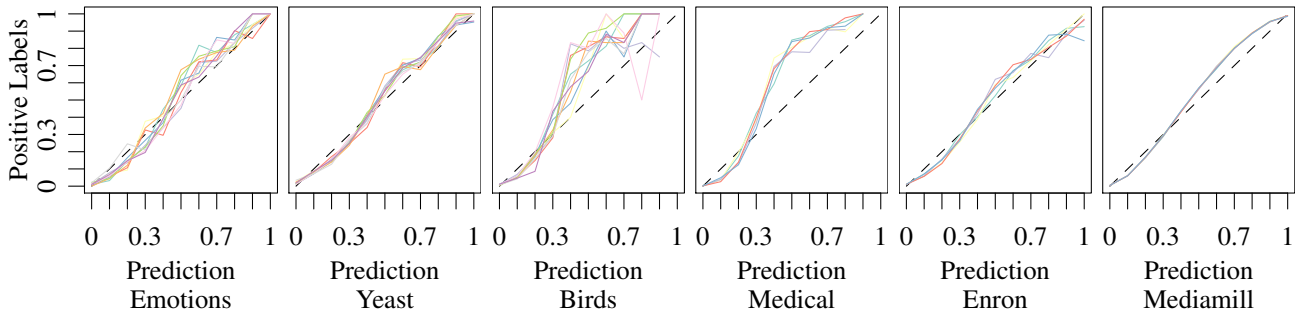[6] https://github.com/mkirchhof/rslBench

Figure 4: Calibration of the Binary Relevance Classifiers on all Benchmark Datasets. Colors Indicate Different Folds.

[Platt, 2000] and relearned pRSL on the same hyperparameter and initial values settings as before. However, the performance stayed on an equal level despite expected sampling uncertainty. This might indicate that the underconfidence is not a reason but a symptom of the performance.

As a word of caution, we want to add that the aforementioned hypotheses require more elaborate testing on simulation and real–world datasets. However, such experiments are beyond the scope of this paper.

# 5 CONCLUSION

We introduced pRSL, a learner for multi–label classification that relies on probabilistically extended propositional logic rules. pRSL takes predictions of arbitrary underlying classifiers as input and models the structure between labels by weighting up and down combinations of labels with respect to their contradiction or fulfillment of rules. These rules can be given as any propositional logic formula or learned in a noisy–or form. In the process, we extended the noisy–or gate to multicategorical input. In comparing pRSL to state–of–the–art black–box and interpretable methods on several common benchmark datasets, we found that pRSL is on par with them, albeit it does not outperform them by a larger margin. We ascribe this to pRSL's interpretable nature, and to the fact that it is not yet implemented on GPU, which could allow for a better gradient descent. Moreover, we report promising results of loopy belief propagation algorithms for approximate inference in the marginal and most–probable–explanation case even on larger datasets.

Focusing on the methodological part in this paper, further advanced classification problems have not been explored. These include learning with missing data, with latent labels, or when dynamically adding and removing classifiers and labels from the ensemble. We will test pRSL's robustness under such conditions in upcoming work, with a focus on the motivating logistics example. Further, we aim to gain more theoretical insight into pRSL's location in the field of similar approaches such as knowledge graphs and probabilistic logic. This might allow to combine pRSL with approaches that model the between–class structures in orthogonal ways.

## Author Contributions

M. Kirchhof conceived and implemented the idea and wrote the paper. L. Schmid conceived the idea and edited the paper. C. Reining advised on required conditions from a practical perspective. M. t. Hompel supervised and acquired funding for the project. M. Pauly supervised and edited the paper.

## Acknowledgements

## References

Yuval Atzmon and Gal Chechik. Probabilistic AND-OR Attribute Grouping for Zero-Shot Learning. In *Proceedings of the Thirty-Forth Conference on Uncertainty in Artificial Intelligence*, 2018.

Jun Chen, Xiaoya Dai, Quan Yuan, Chao Lu, and Haifeng Huang. Towards Interpretable Clinical Diagnosis with Bayesian Network Ensembles Stacked on Entity-Aware CNNs. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3143–3153, 2020.

John D Cook. Basic Properties of the Soft Maximum. *UT*

*MD Anderson Cancer Center Department of Biostatistics Working Paper Series, Working Paper 70*, 2011.

Gregory F Cooper. The Computational Complexity of Probabilistic Inference using Bayesian Belief Networks. *Artificial Intelligence*, 42(2-3):393–405, 1990.

Natalia Culakova, Dan Murphy, Joao Gante, Carlos Ledezma, Vahan Hovhannisyan, and Alan Mosca. How to Calibrate your Neural Network Classifier: Getting True Probabilities from a Classification Model. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 3499–3500, 2020.

Krzysztof Dembczynski, Weiwei Cheng, and Eyke Hüllermeier. Bayes Optimal Multilabel Classification via Probabilistic Classifier Chains. In *Proceedings of the 27th International Conference on International Conference on Machine Learning*, pages 279–286, 2010.

Krzysztof Dembczyński, Willem Waegeman, Weiwei Cheng, and Eyke Hüllermeier. On label dependence and loss minimization in multi-label classification. *Machine Learning*, 88(1-2):5–45, 2012.

Nan Ding, Jia Deng, Kevin P Murphy, and Hartmut Neven. Probabilistic Label Relation Graphs with Ising Models. In *Proceedings of the 2015 IEEE International Conference on Computer Vision*, pages 1161–1169, 2015.

Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On Calibration of Modern Neural Networks. In *Proceedings of the 34th International Conference on Machine Learning*, pages 1321–1330. PMLR, 2017.

Søren Højsgaard. Graphical Independence Networks with the gRain Package for R. *Journal of Statistical Software*, 46(10):1–26, 2012.

Geng Ji, Dehua Cheng, Huazhong Ning, Changhe Yuan, Hanning Zhou, Liang Xiong, and Erik B. Sudderth. Variational Training for Large-Scale Noisy-OR Bayesian Networks. In *Proceedings of The 35th Uncertainty in Artificial Intelligence Conference*, pages 873–882, 2020.

Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations*, 2015.

Eyal Krupka, Kfir Karmon, Noam Bloom, Daniel Freedman, Ilya Gurvich, Aviv Hurvitz, Ido Leichter, Yoni Smolin, Yuval Tzairi, Alon Vinnikov, et al. Toward Realistic Hands Gesture Interface: Keeping it Simple for Developers and Machines. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, pages 1887–1898, 2017.

Meelis Kull, Miquel Perello Nieto, Markus Kängsepp, Telmo Silva Filho, Hao Song, and Peter Flach. Beyond Temperature Scaling: Obtaining Well–calibrated Multi–class Probabilities with Dirichlet Calibration. In *Advances in Neural Information Processing Systems*, volume 32, pages 12316–12326, 2019.

Christoph H Lampert, Hannes Nickisch, and Stefan Harmeling. Attribute–Based Classification for Zero–Shot Visual Object Categorization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(3):453–465, 2013.

Chung-Wei Lee, Wei Fang, Chih-Kuan Yeh, and Yu-Chiang Frank Wang. Multi–label Zero–Shot Learning with Structured Knowledge Graphs. In *Proceedings of the 2018 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1576–1585, 2018.

Lu Liu, Tianyi Zhou, Guodong Long, Jing Jiang, and Chengqi Zhang. Attribute Propagation Network for Graph Zero–Shot Learning. 34(04):4868–4875, 2020.

James D Malley, Jochen Kruppa, Abhijit Dasgupta, Karen G Malley, and Andreas Ziegler. Probability Machines: Consistent Probability Estimation using Nonparametric Learning Machines. *Methods of Information in Medicine*, 51(1):74–81, 2012.

Pauline Maurice, Adrien Malaisé, Clélie Amiot, Nicolas Paris, Guy-Junior Richard, Olivier Rochel, and Serena Ivaldi. Human Movement and Ergonomics: An Industry–Oriented Dataset for Collaborative Robotics. *The International Journal of Robotics Research*, 38(14):1529–1537, 2019.

Eneldo Loza Mencía and Frederik Janssen. Learning Rules for Multi–label Classification: A Stacking and a Separate–And–Conquer Approach. *Machine Learning*, 105(1):77–126, 2016.

Kevin P Murphy. *Machine Learning: A Probabilistic Perspective*. MIT press, 2012.

Kevin P. Murphy, Yair Weiss, and Michael I. Jordan. Loopy Belief Propagation for Approximate Inference: An Empirical Study. In *Proceedings of the Fifteenth Conference on Uncertainty in Artificial Intelligence*, pages 467—-475, 1999.

Friedrich Niemann, Christopher Reining, Fernando Moya Rueda, Nilah Ravi Nair, Janine Anika Steffens, Gernot A Fink, and Michael ten Hompel. Lara: Creating a dataset for human activity recognition in logistics using semantic attributes. *Sensors*, 20(15):4083, 2020.

Arjun Pakrashi, Derek Greene, and Brian MacNamee. Benchmarking Multi–label Classification Algorithms. In *24th Irish Conference on Artificial Intelligence and Cognitive Science*, 2016.

Giovanni Parmigiani and Lurdes Inoue. *Decision Theory: Principles and Approaches*. John Wiley & Sons, 2009.

Judea Pearl. *Probabilistic Reasoning in Intelligent Systems*. Morgan Kaufmann, 1988.

DT Pham and MS Aksoy. RULES: A Simple Rule Extraction System. *Expert Systems with Applications*, 8(1): 59–65, 1995.

John Platt. Probabilistic Outputs for Support Vector Machines and Comparisons to Regularized Likelihood Methods. *Adv. Large Margin Classif.*, 10, 06 2000.

Arwa B Raies and Vladimir B Bajic. In Silico Toxicology: Comprehensive Benchmarking of Multi–label Classification Methods Applied to Chemical Toxicity Data. *Wiley Interdisciplinary Reviews: Computational Molecular Science*, 8(3):e1352, 2018.

Michael Rapp, Eneldo Loza Mencía, Johannes Fürnkranz, Vu-Linh Nguyen, and Eyke Hüllermeier. Learning Gradient Boosted Multi–label Classification Rules. In *Machine Learning and Knowledge Discovery in Databases, ECML PKDD 2020*, 2020.

Christopher Reining, Friedrich Niemann, Fernando Moya Rueda, Gernot A Fink, and Michael ten Hompel. Human Activity Recognition for Production and Logistics – A Systematic Literature Review. *Information*, 10 (8):245, 2019.

Pedro Manuel Santos Ribeiro, Ana Clara Matos, Pedro Henrique Santos, and Jaime S Cardoso. Machine Learning Improvements to Human Motion Tracking with IMUs. *Sensors*, 20(21):6383, 2020.

Fernando Moya Rueda and Gernot Fink. From Human Pose to On–Body Devices for Human–Activity Recognition. In *26th International Conference on Pattern Recognition (ICPR)*, pages 10066–10073, 2021.

Rongbo Shen, Fuhao Zou, Jingkuan Song, Kezhou Yan, and Ke Zhou. EFUI: An Ensemble Framework using Uncertain Inference for Pornographic Image Recognition. *Neurocomputing*, 322:166–176, 2018.

Qing Wang and Ping Li. Randomized Bayesian Network Classifiers. In *International Workshop on Multiple Classifier Systems*, pages 319–330, 2013.

Adrian Weller and Tony Jebara. On MAP inference by MWSS on Perfect Graphs. In *Proceedings of the Twenty–Ninth Conference on Uncertainty in Artificial Intelligence*, pages 684–693, 2013.

Baoyuan Wu, Fan Jia, Wei Liu, Bernard Ghanem, and Siwei Lyu. Multi–label Learning With Missing Labels using Mixed Dependency Graphs. *International Journal of Computer Vision*, 126(8):875–896, 2018.

Yuelong Xia, Ke Chen, and Yun Yang. Multi–label Classification with Weighted Classifier Selection and Stacked Ensemble. *Information Sciences*, 2020.

Yongqin Xian, Christoph H Lampert, Bernt Schiele, and Zeynep Akata. Zero–Shot Learning – A Comprehensive Evaluation of the Good, the Bad and the Ugly. *IEEE transactions on pattern analysis and machine intelligence*, 41(9):2251–2265, 2018.

Aston Zhang, Zachary C. Lipton, Mu Li, and Alexander J. Smola. *Dive into Deep Learning*. 2020a.

Min-Ling Zhang, Yu-Kun Li, Hao Yang, and Xu-Ying Liu. Towards Class–Imbalance Aware Multi–label Learning. *IEEE Transactions on Cybernetics*, 2020b.

# SUPPLEMENTARY MATERIAL

## 5.1 GRADIENTS

To perform gradient descent, we differentiate the joint log likelihood of the true categories $\boldsymbol{\ell}^* = (\ell_1^*, \ldots, \ell_J^*)$ by the noisy–or's inhibition probabilities $q_{jm}^k$ for all rules $k = 1, \ldots, K$, and all $m = 1, \ldots, M(j)$ categories of labels $j = 1, \ldots, J$. To simplify notation, we define

$$\boldsymbol{R}_{-k} := (R_1, \ldots, R_{k-1}, R_{k+1}, \ldots, R_K),$$
$$\boldsymbol{L}_{-j} := (L_1, \ldots, L_{j-1}, L_{j+1} \ldots, L_J),$$
$$\boldsymbol{\ell}_{-j} := (\ell_1, \ldots, \ell_{j-1}, \ell_{j+1} \ldots, \ell_J),$$
$$\boldsymbol{\ell}_{-j}^* := (\ell_1^*, \ldots, \ell_{j-1}^*, \ell_{j+1}^* \ldots, \ell_J^*),$$
$$\mathscr{L}_{j=m} := \bigotimes_{v=1}^{j-1} \{1, \ldots, M(v)\} \otimes \{m\} \otimes \bigotimes_{v=j+1}^{J} \{1, \ldots, M(v)\},$$
$$\mathscr{L}_{j\neq m} := \bigotimes_{v=1}^{j-1} \{1, \ldots, M(v)\} \otimes$$
$$\{1, \ldots, m-1, m+1, \ldots, M(j)\} \otimes$$
$$\bigotimes_{v=j+1}^{J} \{1, \ldots, M(v)\} \text{ and}$$
$$q_{-j\ell}^k := \prod_{v=1}^{j-1} q_{v\ell_v}^k \cdot \prod_{v=j+1}^{J} q_{v\ell_v}^k \text{ and}$$
$$q_{-j\ell^*}^k := \prod_{v=1}^{j-1} q_{v\ell_v^*}^k \cdot \prod_{v=j+1}^{J} q_{v\ell_v^*}^k.$$

### 5.1.1 All Labels Known

As discussed in the paper, we first differentiate in the case where all true categories $\ell_1^*, \ldots, \ell_J^*$ are known.

**Case 1: Category is incorrect.** We start for those $q_{jm}^k$ where $m$ is not the true category, that is $\ell_j^* \neq m$:

$$D_{jm}^k := \frac{\partial}{\partial q_{jm}^k} P(\boldsymbol{L} = \boldsymbol{\ell}^* | \boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x})$$
$$= \frac{\partial}{\partial q_{jm}^k} \frac{P(R_k = 1, \boldsymbol{L} = \boldsymbol{\ell}^* | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x})}{P(R_k = 1 | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x})}.$$

From here on, all probabilities stay conditioned on $\boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$, so we will write this condition behind the expression in the following formulas.

$$D_{jm}^k = \frac{\partial}{\partial q_{jm}^k} \frac{P(R_k = 1, \boldsymbol{L} = \boldsymbol{\ell}^*)}{\sum\limits_{\boldsymbol{\ell} \in \mathscr{L}} (1 - q_{j\ell_j}^k q_{-j\ell}^k) P(\boldsymbol{L} = \boldsymbol{\ell})} \, | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$$

$$= \frac{\partial}{\partial q_{jm}^k} (P(R_k = 1, \boldsymbol{L} = \boldsymbol{\ell}^*)) \cdot$$

$$\left( \sum_{\boldsymbol{\ell} \in \mathscr{L}} P(\boldsymbol{L} = \boldsymbol{\ell}) - \sum_{\boldsymbol{\ell} \in \mathscr{L}_{j \neq m}} q_{j\ell_j}^k q_{-j\ell}^k P(\boldsymbol{L} = \boldsymbol{\ell}) - q_{jm}^k \sum_{\boldsymbol{\ell} \in \mathscr{L}_{j=m}} q_{-j\ell}^k P(\boldsymbol{L} = \boldsymbol{\ell}) \right)^{-1} | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}.$$

Simplifying the notation, the above expression can be written as:

$$D_{jm}^k = \frac{\partial}{\partial q_{jm}^k} \frac{a_1}{a_2 - q_{jm}^k a_3} \, | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$$
$$= \frac{a_1 a_3}{(a_2 - q_{jm}^k a_3)^2} \, | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$$

for suitable choices of $a_1$, $a_2$, and $a_3$. Substituting back and multiplying with $q_{jm}^k (q_{jm}^k)^{-1}$ we get:

$$D_{jm}^k = \frac{P(R_k = 1, \boldsymbol{L} = \boldsymbol{\ell}^*) \sum\limits_{\boldsymbol{\ell} \in \mathscr{L}_{j=m}} q_{jm}^k q_{-j\ell}^k P(\boldsymbol{L} = \boldsymbol{\ell})}{q_{jm}^k P(R_k = 1)^2} \, | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$$

$$= \frac{P(R_k = 1, \boldsymbol{L} = \boldsymbol{\ell}^*) P(L_j = m, R_k = 0)}{q_{jm}^k P(R_k = 1)^2} \, | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}.$$

The first term in the nominator is computationally complex and might get numerically instable with an increasing number of labels $J$, but since we optimize for the log–likelihood, it vanishes via the chain rule:

$$\frac{\partial}{\partial q_{jm}^k} \log(P(\boldsymbol{L} = \boldsymbol{\ell}^* | \boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x}))$$
$$= \frac{1}{P(\boldsymbol{L} = \boldsymbol{\ell}^* | \boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x})} D_{jm}^k$$
$$= \frac{P(R_k = 1) P(R_k = 1, \boldsymbol{L} = \boldsymbol{\ell}^*) P(L_j = m, R_k = 0)}{P(\boldsymbol{L} = \boldsymbol{\ell}^*, R_k = 1) q_{jm}^k P(R_k = 1)^2} \, | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$$
$$= \frac{P(L_j = m, R_k = 0)}{q_{jm}^k P(R_k = 1)} \, | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$$
$$= \frac{P(L_j = m | R_k = 0)(1 - P(R_k = 1))}{q_{jm}^k P(R_k = 1)} \, | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}.$$

Note that $P(L_j = m | R_k = 0)$ is returned for all categories $m = 1, \ldots, M(j)$ of all labels $L_j, j = 1, \ldots, J$, by a single marginal query. Hence, two marginal queries to the network are sufficient to compute the gradient of all inhibition probabilities related to a rule.

**Case 2: Category is correct.** Let us now differentiate for $q_{jm}^k$ where $m$ is the true category, that is $\ell_j^* = m$:

$$D_{jm}^k = \frac{\partial}{\partial q_{jm}^k} P(\boldsymbol{L} = \boldsymbol{\ell}^* | \boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x})$$

$$= \frac{\partial}{\partial q_{jm}^k} (P(\boldsymbol{L} = \boldsymbol{\ell}^*) - q_{j\ell_j^*}^k q_{-j\ell^*}^k P(\boldsymbol{L} = \boldsymbol{\ell}^*)) \cdot$$

$$\left( \sum_{\boldsymbol{\ell} \in \mathscr{L}} P(\boldsymbol{L} = \boldsymbol{\ell}) - \sum_{\boldsymbol{\ell} \in \mathscr{L}_{j \neq m}} q_{j\ell_j}^k q_{-j\ell}^k P(\boldsymbol{L} = \boldsymbol{\ell}) - \right.$$

$$\left. q_{jm}^k \sum_{\boldsymbol{\ell} \in \mathscr{L}_{j=m}} q_{-j\ell}^k P(\boldsymbol{L} = \boldsymbol{\ell}) \right)^{-1} | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x} .$$

Again, we simplify the notation to make differentiation easier to see:

$$= \frac{\partial}{\partial q_{jm}^k} \frac{b_1 - q_{jm}^k b_2}{b_3 - q_{jm}^k b_4} | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$$

$$= \frac{b_1 b_4 - b_2 b_3}{(b_3 - q_{jm}^k b_4)^2} | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x} .$$

Substituting back and multiplying by $q_{jm}^k (q_{jm}^k)^{-1}$, we get:

$$= \frac{P(\boldsymbol{L} = \boldsymbol{\ell}^*)(q_{jm}^k b_4 - q_{jm}^k q_{-j\ell}^k b_3)}{q_{jm}^k P(R_k = 1)^2} | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$$

$$= (q_{jm}^k P(R_k = 1)^2)^{-1} (P(\boldsymbol{L} = \boldsymbol{l}^*)(P(R_k = 0, L_j = m) -$$

$$q_{jm}^k q_{-j\ell}^k (P(L_j = m) + P(L_j \neq m) -$$

$$P(R_k = 0, L_j \neq m)))) | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$$

$$= (q_{jm}^k P(R_k = 1)^2)^{-1} (P(\boldsymbol{L} = \boldsymbol{\ell}^*)(P(R_k = 0, L_j = m) -$$

$$q_{jm}^k q_{-j\ell}^k (P(L_j = m) + P(R_k = 1, L_j \neq m)))) | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$$

$$= (q_{jm}^k P(R_k = 1)^2)^{-1} (P(\boldsymbol{L} = \boldsymbol{\ell}^*)(P(R_k = 0, L_j = m) -$$

$$q_{jm}^k q_{-j\ell}^k (P(R_k = 0, L_j = m) + P(R_k = 1, L_j = m) +$$

$$P(R_k = 1, L_j \neq m)))) | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$$

$$= (q_{jm}^k P(R_k = 1)^2)^{-1} (P(\boldsymbol{L} = \boldsymbol{\ell}^*)(P(R_k = 0, L_j = m) -$$

$$q_{jm}^k q_{-j\ell}^k (P(R_k = 0, L_j = m) + P(R_k = 1)))) | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$$

$$= (q_{jm}^k P(R_k = 1)^2)^{-1} (P(\boldsymbol{L} = \boldsymbol{l}^*)((1 - q_{jm}^k q_{-j\ell}^k) \cdot$$

$$P(R_k = 0, L_j = m) - q_{jm}^k q_{-j\ell}^k P(R_k = 1))) | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$$

$$= (q_{jm}^k P(R_k = 1)^2)^{-1} (P(\boldsymbol{L} = \boldsymbol{\ell}^*, R_k = 1) P(R_k = 0, L_j = m) -$$

$$P(\boldsymbol{L} = \boldsymbol{\ell}^*) q_{jm}^k q_{-j\ell}^k P(R_k = 1))) | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x} .$$

Just as before, the computational heavy terms cancel out when we optimize for log–likelihood:

$$\frac{\partial}{\partial q_{jm}^k} \log(P(\boldsymbol{L} = \boldsymbol{l}^* | \boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x}))$$

$$= \frac{P(L_j = m | R_k = 0)(1 - P(R_k = 1))}{q_{jm}^k P(R_k = 1)} -$$

$$\frac{q_{-j\ell^*}^k}{1 - q_{jm}^k q_{-j\ell^*}^k} | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x} .$$

This term is similar to the previous gradient and requires the same marginal queries.

### 5.1.2 Missing Labels

In this case only a subset of labels $\boldsymbol{L}' \subset \boldsymbol{L}$ is known and takes the categories $\boldsymbol{L}' = \boldsymbol{\ell}'$, while the ground truth for all other labels $\boldsymbol{L}^0 = \boldsymbol{L} \backslash \boldsymbol{L}'$ is unknown. In consequence, the optimization goal changes slightly. We define $\mathscr{L}^0, \mathscr{L}_{j=m}^0, \mathscr{L}_{j \neq m}^0$ and $\boldsymbol{\ell}^0$ in analogy to before.

**Case 1: Category is known and incorrect.** As in the previous section, we will first consider the case where $L_j \in \boldsymbol{L}'$ is known and $m$ is not the true category, that is $\ell_j \neq m$:

$$\frac{\partial}{\partial q_{jm}^k} P(\boldsymbol{L}' = \boldsymbol{l}' | \boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x})$$

$$= \frac{\partial}{\partial q_{jm}^k} \frac{\sum_{\boldsymbol{\ell}^0 \in \mathscr{L}^0} P(R_k = 1, \boldsymbol{L}' = \boldsymbol{\ell}', \boldsymbol{L}^0 = \boldsymbol{\ell}^0 | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x})}{P(R_k = 1 | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x})}$$

As the nominator is a factor independent of $q_{jm}^k$, the sum can be placed before the expression and the computations follow those in Section 5.1.1. When taking the logarithm, the gradient is divided by $P(\boldsymbol{L}' = \boldsymbol{\ell}' | \boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x})$, so that the resulting gradient remains the same as in Section 5.1.1:

$$\frac{\partial}{\partial q_{jm}^k} \log(P(\boldsymbol{L}' = \boldsymbol{\ell}' | \boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x}))$$

$$= \frac{P(L_j = m | R_k = 0)(1 - P(R_k = 1))}{q_{jm}^k P(R_k = 1)} | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x} .$$

**Case 2: Category is known and correct.** In the case that $m$ is known and the true category, that is $L_j \in \boldsymbol{L}'$ and $\ell_j = m$, the sum can again be factored out. Thus, we can again make use of the gradients from Section 5.1.1:

$$\frac{\partial}{\partial q_{jm}^k} P(\boldsymbol{L}' = \boldsymbol{\ell}' | \boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x})$$

$$= \sum_{\boldsymbol{\ell}^0 \in \mathscr{L}^0} \frac{\partial}{\partial q_{jm}^k} P(\boldsymbol{L}' = \boldsymbol{\ell}', \boldsymbol{L}^0 = \boldsymbol{\ell}^0 | \boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x})$$

$$= (q_{jm}^k P(R_k = 1)^2)^{-1} (P(\boldsymbol{L}' = \boldsymbol{\ell}', R_k = 1) P(R_k = 0, L_j = m) -$$

$$P(\boldsymbol{L}' = \boldsymbol{\ell}', R_k = 0) P(R_k = 1))) | \boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x} .$$

When taking the logarithm, the expression gets slightly more complicated than in Section 5.1.1:

$$\frac{\partial}{\partial q_{jm}^k} \log(P(\boldsymbol{L}' = \boldsymbol{\ell}'|\boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x}))$$

$$= (P(\boldsymbol{L}' = \boldsymbol{\ell}'|\boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x}))^{-1} \frac{\partial}{\partial q_{jm}^k} P(\boldsymbol{L}' = \boldsymbol{\ell}'|\boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x})$$

$$= \frac{P(L_j = m|R_k = 0)(1 - P(R_k = 1))}{q_{jm}^k P(R_k = 1)} -$$

$$\frac{P(\boldsymbol{L}' = \boldsymbol{\ell}', R_k = 0)}{P(\boldsymbol{L}' = \boldsymbol{\ell}', R_k = 1)} |\boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$$

$$= \frac{P(L_j = m|R_k = 0)(1 - P(R_k = 1))}{q_{jm}^k P(R_k = 1)} -$$

$$\frac{P(R_k = 0|\boldsymbol{L}' = \boldsymbol{\ell}')}{P(R_k = 1|\boldsymbol{L}' = \boldsymbol{\ell}')} |\boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x} .$$

So, one additional marginal query has to be calculated per rule. Note that this expression simplifies to that of Section 5.1.1 if all labels are known, that is if $\boldsymbol{L}' = \boldsymbol{L}$, because the conditional probabilities in the second fraction are then the rules conditional probability tables, given by $\boldsymbol{q}$ alone.

**Case 3: Label is unknown.** Interestingly, the gradients can also be computed for labels that have no ground truth, that is $L_j \in \boldsymbol{L}^0$. Obviously, we do not need to distinguish whether $m$ is correct or not. Let $\boldsymbol{L}_{-j}^0 := \boldsymbol{L}^0 \backslash L_j$ and let $q_{-j\ell}^k := \prod_{v:L_v \in \boldsymbol{L}'} q_{v\ell_v}^k \cdot \prod_{v:L_v \in \boldsymbol{L}_{-j}^0} q_{v\ell_v}^k$ for simpler notation.

$$\frac{\partial}{\partial q_{jm}^k} P(\boldsymbol{L}' = \boldsymbol{\ell}'|\boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x})$$

$$= \frac{\partial}{\partial q_{jm}^k} \frac{\sum\limits_{\boldsymbol{\ell}^0 \in \mathscr{L}^0} (1 - q_{j\ell_j^*}^k q_{-j\ell^*}^k) P(\boldsymbol{L}' = \boldsymbol{\ell}', \boldsymbol{L}^0 = \boldsymbol{\ell}^0)}{P(R_k = 1)} |\boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$$

$$= \frac{\partial}{\partial q_{jm}^k} \left( \sum_{\boldsymbol{\ell}^0 \in \mathscr{L}^0} P(\boldsymbol{L}' = \boldsymbol{\ell}', \boldsymbol{L}^0 = \boldsymbol{\ell}^0) - \right.$$

$$\sum_{\boldsymbol{\ell}^0 \in \mathscr{L}_{-j}^0} q_{j\ell_j^*}^k q_{-j\ell^*}^k P(\boldsymbol{L}' = \boldsymbol{\ell}', \boldsymbol{L}^0 = \boldsymbol{\ell}^0) -$$

$$\left. q_{jm}^k \sum_{\boldsymbol{\ell}^0 \in \mathscr{L}_{-j}^0} q_{-j\ell^*}^k P(\boldsymbol{L}' = \boldsymbol{\ell}', \boldsymbol{L}^0 = \boldsymbol{\ell}^0) \right) \cdot$$

$$\left( \sum_{\boldsymbol{\ell} \in \mathscr{L}} P(\boldsymbol{L} = \boldsymbol{\ell}) - \sum_{\boldsymbol{\ell} \in \mathscr{L}_{j \neq m}} q_{j\ell_j}^k q_{-j\ell}^k P(\boldsymbol{L} = \boldsymbol{\ell}) - \right.$$

$$\left. q_{jm}^k \sum_{\boldsymbol{\ell} \in \mathscr{L}_{j=m}} q_{-j\ell}^k P(\boldsymbol{L} = \boldsymbol{\ell}) \right)^{-1} |\boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x} .$$

Again, we substitute to make differentiation easier to see:

$$= \frac{\partial}{\partial q_{jm}^k} \frac{c_1 - q_{jm}^k c_2}{c_3 - q_{jm}^k c_4} |\boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$$

$$= \frac{c_1 c_4 - c_2 c_3}{(c_3 - q_{jm}^k c_4)^2} |\boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x} .$$

By substituting back and replacing the sum expressions with corresponding probability terms, we receive:

$$= (P(R_k = 1)^2)^{-1} (P(\boldsymbol{L}' = \boldsymbol{\ell}') - P(L_j \neq m, \boldsymbol{L}' = \boldsymbol{\ell}', R_k = 0)) \cdot$$

$$(q_{jm}^k)^{-1} P(L_j = m, R_k = 0) - (q_{jm}^k)^{-1} \cdot$$

$$P(L_j = m, \boldsymbol{L}' = \boldsymbol{\ell}', R_k = 0) \cdot (P(L_j = m) +$$

$$P(L_j \neq m) - P(L_j \neq m, R_k = 0)) |\boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$$

$$= (q_{jm}^k P(R_k = 1)^2)^{-1} (P(L_j = m, \boldsymbol{L}' = \boldsymbol{\ell}') +$$

$$P(L_j \neq m, \boldsymbol{L}' = \boldsymbol{\ell}') - P(L_j \neq m, \boldsymbol{L}' = \boldsymbol{\ell}', R_k = 0)) \cdot$$

$$P(L_j = m, R_k = 0) - P(L_j = m, \boldsymbol{L}' = \boldsymbol{\ell}', R_k = 0) \cdot$$

$$(P(L_j = m) + P(L_j \neq m, R_k = 1))) |\boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$$

$$= (q_{jm}^k P(R_k = 1)^2)^{-1} (P(L_j = m, \boldsymbol{L}' = \boldsymbol{\ell}', R_k = 0) +$$

$$P(L_j = m, \boldsymbol{L}' = \boldsymbol{\ell}', R_k = 1) + P(L_j \neq m, \boldsymbol{L}' = \boldsymbol{\ell}', R_k = 1)) \cdot$$

$$P(L_j = m, R_k = 0) - P(L_j = m, \boldsymbol{L}' = \boldsymbol{\ell}', R_k = 0) \cdot$$

$$(P(L_j = m, R_k = 0) + P(L_j = m, R_k = 1) +$$

$$P(L_j \neq m, R_k = 1))) |\boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$$

$$= (q_{jm}^k P(R_k = 1)^2)^{-1} \cdot$$

$$(P(L_j = m, \boldsymbol{L}' = \boldsymbol{\ell}', R_k = 1) P(L_j = m, R_k = 0) +$$

$$P(L_j = m, \boldsymbol{L}' = \boldsymbol{\ell}', R_k = 0) P(L_j = m, R_k = 0) +$$

$$P(L_j \neq m, \boldsymbol{L}' = \boldsymbol{\ell}', R_k = 1) P(L_j = m, R_k = 0) -$$

$$P(L_j = m, \boldsymbol{L}' = \boldsymbol{\ell}', R_k = 0) P(L_j = m, R_k = 0) -$$

$$P(L_j = m, \boldsymbol{L}' = \boldsymbol{\ell}', R_k = 0) P(L_j = m, R_k = 1) -$$

$$P(L_j = m, \boldsymbol{L}' = \boldsymbol{\ell}', R_k = 0) P(L_j \neq m, R_k = 1)) |\boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x}$$

$$= (q_{jm}^k P(R_k = 1)^2)^{-1} (P(\boldsymbol{L}' = \boldsymbol{\ell}', R_k = 1) P(L_j = m, R_k = 0) -$$

$$P(L_j = m, \boldsymbol{L}' = \boldsymbol{\ell}', R_k = 0) P(R_k = 1)) |\boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x} .$$

As before, the log likelihood removes the computationally complicated terms:

$$\frac{\partial}{\partial q_{jm}^k} \log(P(\boldsymbol{L}' = \boldsymbol{\ell}'|\boldsymbol{R} = \boldsymbol{1}, \boldsymbol{x}))$$

$$= \frac{P(L_j = m, R_k = 0)}{q_{jm}^k P(R_k = 1)} -$$

$$\frac{P(L_j = m, R_k = 0|\boldsymbol{L}' = \boldsymbol{\ell}')}{q_{jm}^k P(R_k = 1|\boldsymbol{L}' = \boldsymbol{\ell}')} |\boldsymbol{R}_{-k} = \boldsymbol{1}, \boldsymbol{x} .$$

At this point it can be seen that the gradient reduces to 0 if no label has a ground truth, that is $\boldsymbol{L}' = \emptyset$, which makes intuitive sense. Applying one last transformation gives a
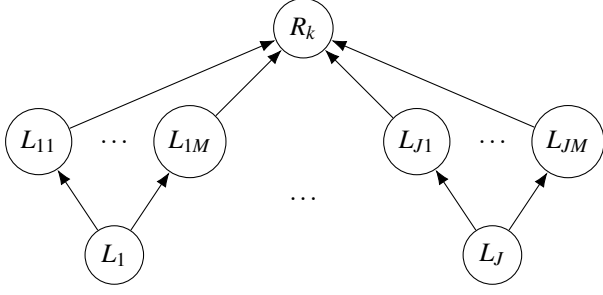
Figure 5: Decomposition of Multicategorical Inputs for a Binary–Input Noisy–Or Gate.

form that is easier to compute:

$$
\begin{aligned}
=&\frac{P(L_j = m|R_k = 0)(1 - P(R_k = 1))}{q^k_{jm}P(R_k = 1)} - \\
&\frac{P(L_j = m|R_k = 0, \boldsymbol{L}' = \boldsymbol{\ell}')(1 - P(R_k = 1|\boldsymbol{L}' = \boldsymbol{\ell}'))}{q^k_{jm}P(R_k = 1|\boldsymbol{L}' = \boldsymbol{\ell}')} \\
&|\boldsymbol{R}_{-k} = \mathbf{1}, \boldsymbol{x} .
\end{aligned}
$$

So, overall four marginal queries are required to compute this gradient, or two more than in the case of known labels.

## 5.2 EXTENSION OF NOISY–OR

The ordinary noisy–or gate as defined in Pearl [1988] is connected to a set of binary input variables $L_1, \ldots, L_J$. Each variable can only set the gate $R_k$ to $R_k = 1$ with a probability $1 - q^k_{j1}$ if the variable itself is $L_j = 1$. Thus the conditional probability distribution that defines $R_k$ is:

$$
P(R_k = 0|L_1 = \ell_1, \ldots, L_J = \ell_J) = \prod_{j=1}^{J}(q^k_{j1})^{\ell_j} .
$$

In our case, the input variables may have multiple categories $m = 1, \ldots, M(j)$ and each of these categories can trigger the gate to be $R_k = 1$ with a probability $1 - q^k_{jm}$. To find the conditional probability distribution of $R_k$ in this case, we start by splitting each input variables $L_j$ up into several binary auxiliary variables $L_{jm}$ where

$$
P(L_{jm} = 1|L_j = a) = \begin{cases} 1, a = m \\ 0, a \neq m \end{cases} = \mathbb{1}_{L_j = m}(L_j) .
$$

These variables can then be connected to the ordinary binary–input noisy–or gate as visualized in Figure 5 with their corresponding inhibition probability $q^k_{jm}$. We will now show that this process extends the binary–input noisy–or gate to multicategorical input naturally. We start with the conditional probability via the above described structure with auxiliary variables.

**Algorithm 1** Simulation Dataset Generation
**Require:** nLabels, nRules, nData
 1: pRSL ← empty model
 2: **for** $i$ in $1, \ldots,$ nLabels **do**
 3: nCategories $\sim U\{2, \ldots, 4\}$
 4: Add label node with nCategories to pRSL
 5: **for** $i$ in $1, \ldots,$ nRules **do**
 6: nCategories $\sim U\{2, \ldots, 5\}$
 7: categories ← Draw nCategories from all
          $\{1, \ldots, M(1), \ldots, 1, \ldots, M(J)\}$
 8: **for** each categories **do**
 9: inhProb $\sim$ Distribution with density
          $f(x) = 2 \cdot (1-x) \cdot \mathbb{1}_{[0,1]}(x)$
10: Add rule with categories and inhProbs to pRSL
11: **for** $i$ in $1, \ldots,$ nData **do**
12: **for** each label node in pRSL **do**
13: classifierOutput[label][i] $\sim$ Dir(1)
14: **return** pRSL, classifierOutput

$$
\begin{aligned}
&P(R_k = 0|\boldsymbol{L} = \boldsymbol{\ell}) \\
=& \sum_{\ell_{11}, \ldots, \ell_{JM(J)} \, in\{0,1\}} P(R_k = 0|L_{11} = \ell_{11}, \ldots, L_{JM(J)} = \ell_{JM(J)}) \cdot \\
& P(L_{11} = \ell_{11}|L_1 = \ell_1) \cdot \ldots \cdot P(L_{JM(J)} = \ell_{JM(J)}|L_J = \ell_J) \\
=& \sum_{\ell_{11}, \ldots, \ell_{JM(J)} \in \{0,1\}} \prod_{j=1}^{J}(q^k_{jm})^{L_{jm}} (\mathbb{1}_{L_1=1}(L_1))^{\ell_{11}} \cdot \ldots \cdot \\
& (\mathbb{1}_{L_1=M(1)}(L_1))^{\ell_{1M(1)}} \cdot \ldots \cdot (\mathbb{1}_{L_J=1}(L_J))^{\ell_{J1}} \cdot \ldots \cdot \\
& (\mathbb{1}_{L_J=M(J)}(L_J))^{\ell_{JM(J)}} .
\end{aligned}
$$

With $0^0 := 1$, we can see that the only time the term inside the sum is not 0 is when $L_{jm} = 1$ iff $L_j = m$ for all $j = 1, \ldots, J$. This leaves open only one possible allocation of the binary auxiliary variables due to their XOR relation within $j$, so that the sum and the auxiliary variables vanish. We finally get a familiar expression that naturally extends the binary noisy–or to the multicategorical case:

$$
P(R_k = 0|\boldsymbol{L} = \boldsymbol{\ell}) = \prod_{j=1}^{J} q^k_{j\ell_j} .
$$

## 5.3 SIMULATION DATASET GENERATION

Algorithm 1 describes the sampling procedure used to generate the simulation datasets used in Section 4.1. In the first ten lines of code, a pRSL model containing nLabels label nodes and nRules rule nodes is randomly generated. In lines 11 to 13, the classifier outputs for each classifier node are simulated by dirichlet noise for nData observations. Both the simulated data and the data–generating pRSL model are returned to allow performing approximate marginal and MPE queries on the correct model in Section 4.1.

14

Table 4: Performance of pRSL on Train, Validation, and Test Data. Mean ± Standard Deviation Between Folds.

| | Emotions | Yeast | Birds | Medical | Enron | Mediamill |
|---|---|---|---|---|---|---|
| | Joint Accuracy (higher = better) | | | | | |
| Train | $0.351 \pm 0.015$ | $0.227 \pm 0.010$ | $0.516 \pm 0.010$ | $0.500 \pm 0.019$ | $0.153 \pm 0.010$ | $0.146 \pm 0.001$ |
| Validation | $0.339 \pm 0.031$ | $0.251 \pm 0.026$ | $0.516 \pm 0.042$ | $0.497 \pm 0.027$ | $0.154 \pm 0.010$ | $0.149 \pm 0.006$ |
| Test | $0.348 \pm 0.067$ | $0.236 \pm 0.015$ | $0.507 \pm 0.032$ | $0.491 \pm 0.031$ | $0.153 \pm 0.020$ | $0.149 \pm 0.002$ |
| | Joint log–Likelihood (higher = better) | | | Label–wise log–Likelihood (higher = better) | | |
| Train | $-1.802 \pm 0.077$ | $-3.589 \pm 0.054$ | $-2.534 \pm 0.060$ | $-1.587 \pm 0.022$ | $-6.598 \pm 0.079$ | $-6.585 \pm 0.006$ |
| Validation | $-1.921 \pm 0.241$ | $-3.538 \pm 0.177$ | $-2.532 \pm 0.269$ | $-1.625 \pm 0.051$ | $-6.564 \pm 0.178$ | $-6.563 \pm 0.005$ |
| Test | $-1.839 \pm 0.273$ | $-3.592 \pm 0.085$ | $-2.458 \pm 0.156$ | $-1.565 \pm 0.120$ | $-6.479 \pm 0.242$ | $-6.532 \pm 0.061$ |
| | Label–wise Hamming Loss (lower = better) | | | | | |
| Train | $0.182 \pm 0.005$ | $0.191 \pm 0.003$ | $0.043 \pm 0.001$ | $0.015 \pm 0.001$ | $0.046 \pm 0.001$ | $0.027 \pm 0.000$ |
| Validation | $0.181 \pm 0.018$ | $0.188 \pm 0.004$ | $0.042 \pm 0.004$ | $0.015 \pm 0.001$ | $0.046 \pm 0.001$ | $0.027 \pm 0.000$ |
| Test | $0.182 \pm 0.022$ | $0.190 \pm 0.005$ | $0.043 \pm 0.002$ | $0.015 \pm 0.001$ | $0.046 \pm 0.001$ | $0.027 \pm 0.000$ |