

pRSL: Interpretable Multi-label Stacking by Learning Probabilistic Rules

Supplementary Material

Michael Kirchhof¹ Lena Schmid¹ Christopher Reining² Michael ten Hompel² Markus Pauly¹

¹Department of Statistics, TU Dortmund University, Dortmund, Germany

²Chair of Material Handling and Warehousing, TU Dortmund University, Dortmund, Germany

1 GRADIENTS

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

$$\mathbf{R}_{-k} := (R_1, \dots, R_{k-1}, R_{k+1}, \dots, R_K),$$

$$\mathbf{L}_{-j} := (L_1, \dots, L_{j-1}, L_{j+1}, \dots, L_J),$$

$$\ell_{-j} := (\ell_1, \dots, \ell_{j-1}, \ell_{j+1}, \dots, \ell_J),$$

$$\ell_{-j}^* := (\ell_1^*, \dots, \ell_{j-1}^*, \ell_{j+1}^*, \dots, \ell_J^*),$$

$$\mathcal{L}_{j=m} := \bigotimes_{v=1}^{j-1} \{1, \dots, M(v)\} \otimes$$

$$\{m\} \otimes \bigotimes_{v=j+1}^J \{1, \dots, M(v)\},$$

$$\mathcal{L}_{j \neq m} := \bigotimes_{v=1}^{j-1} \{1, \dots, M(v)\} \otimes$$

$$\{1, \dots, m-1, m+1, \dots, M(j)\} \otimes$$

$$\bigotimes_{v=j+1}^J \{1, \dots, M(v)\} \text{ and}$$

$$q_{-j\ell}^k := \prod_{v=1}^{j-1} q_{v\ell_v}^k \cdot \prod_{v=j+1}^J q_{v\ell_v}^k \text{ and}$$

$$q_{-j\ell^*}^k := \prod_{v=1}^{j-1} q_{v\ell_v^*}^k \cdot \prod_{v=j+1}^J q_{v\ell_v^*}^k.$$

1.1 ALL LABELS KNOWN

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

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

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

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

$$\begin{aligned} D_{jm}^k &= \frac{\partial}{\partial q_{jm}^k} \frac{P(R_k = 1, \mathbf{L} = \ell^*)}{\sum_{\ell \in \mathcal{L}} (1 - q_{j\ell_j}^k q_{-j\ell}^k) P(\mathbf{L} = \ell)} \Big|_{\mathbf{R}_{-k} = \mathbf{1}, \mathbf{x}} \\ &= \frac{\partial}{\partial q_{jm}^k} (P(R_k = 1, \mathbf{L} = \ell^*) \cdot \\ &\quad \left(\sum_{\ell \in \mathcal{L}} P(\mathbf{L} = \ell) - \sum_{\ell \in \mathcal{L}_{j \neq m}} q_{j\ell_j}^k q_{-j\ell}^k P(\mathbf{L} = \ell) - \right. \\ &\quad \left. q_{jm}^k \sum_{\ell \in \mathcal{L}_{j=m}} q_{-j\ell}^k P(\mathbf{L} = \ell) \right)^{-1} \Big|_{\mathbf{R}_{-k} = \mathbf{1}, \mathbf{x}}. \end{aligned}$$

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

$$\begin{aligned} D_{jm}^k &= \frac{\partial}{\partial q_{jm}^k} \frac{a_1}{a_2 - q_{jm}^k a_3} \Big|_{\mathbf{R}_{-k} = \mathbf{1}, \mathbf{x}} \\ &= \frac{a_1 a_3}{(a_2 - q_{jm}^k a_3)^2} \Big|_{\mathbf{R}_{-k} = \mathbf{1}, \mathbf{x}} \end{aligned}$$

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:

$$\begin{aligned} D_{jm}^k &= P(R_k = 1, \mathbf{L} = \ell^*) \cdot \\ &\quad \frac{\sum_{\ell \in \mathcal{L}_{j=m}} q_{jm}^k q_{-j\ell}^k P(\mathbf{L} = \ell)}{q_{jm}^k P(R_k = 1)^2} \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} \\ &= P(R_k = 1, \mathbf{L} = \ell^*) \cdot \\ &\quad \frac{P(L_j = m, R_k = 0)}{q_{jm}^k P(R_k = 1)^2} \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} . \end{aligned}$$

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:

$$\begin{aligned} &\frac{\partial}{\partial q_{jm}^k} \log(P(\mathbf{L} = \ell^* \mid \mathbf{R} = \mathbf{1}, \mathbf{x})) \\ &= \frac{1}{P(\mathbf{L} = \ell^* \mid \mathbf{R} = \mathbf{1}, \mathbf{x})} D_{jm}^k \\ &= \frac{P(R_k = 1)}{P(\mathbf{L} = \ell^*, R_k = 1)} \cdot \\ &\quad P(R_k = 1, \mathbf{L} = \ell^*) \frac{P(L_j = m, R_k = 0)}{q_{jm}^k P(R_k = 1)^2} \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} \\ &= \frac{P(L_j = m, R_k = 0)}{q_{jm}^k P(R_k = 1)} \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} \\ &= \frac{P(L_j = m \mid R_k = 0)(1 - P(R_k = 1))}{q_{jm}^k P(R_k = 1)} \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} . \end{aligned}$$

Note that $P(L_j = m \mid R_k = 0)$ is returned for all categories $m = 1, \dots, M(j)$ of all labels $L_j, j = 1, \dots, J$, 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 $\ell_j^* = m$:

$$\begin{aligned} D_{jm}^k &= \frac{\partial}{\partial q_{jm}^k} P(\mathbf{L} = \ell^* \mid \mathbf{R} = \mathbf{1}, \mathbf{x}) \\ &= \frac{\partial}{\partial q_{jm}^k} (P(\mathbf{L} = \ell^*) - q_{j\ell_j^*}^k q_{-j\ell^*}^k P(\mathbf{L} = \ell^*)) \cdot \\ &\quad \left(\sum_{\ell \in \mathcal{L}} P(\mathbf{L} = \ell) - \sum_{\ell \in \mathcal{L}_{j \neq m}} q_{j\ell_j}^k q_{-j\ell}^k P(\mathbf{L} = \ell) - \right. \\ &\quad \left. q_{jm}^k \sum_{\ell \in \mathcal{L}_{j=m}} q_{-j\ell}^k P(\mathbf{L} = \ell) \right)^{-1} \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} . \end{aligned}$$

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

$$\begin{aligned} &= \frac{\partial}{\partial q_{jm}^k} \frac{b_1 - q_{jm}^k b_2}{b_3 - q_{jm}^k b_4} \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} \\ &= \frac{b_1 b_4 - b_2 b_3}{(b_3 - q_{jm}^k b_4)^2} \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} . \end{aligned}$$

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

$$\begin{aligned} &= \frac{P(\mathbf{L} = \ell^*)(q_{jm}^k b_4 - q_{jm}^k q_{-j\ell}^k b_3)}{q_{jm}^k P(R_k = 1)^2} \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} \\ &= (q_{jm}^k P(R_k = 1)^2)^{-1} \cdot \\ &\quad (P(\mathbf{L} = \ell^*)(P(R_k = 0, L_j = m) - \\ &\quad q_{jm}^k q_{-j\ell}^k (P(L_j = m) + P(L_j \neq m) - \\ &\quad P(R_k = 0, L_j \neq m)))) \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} \\ &= (q_{jm}^k P(R_k = 1)^2)^{-1} \cdot \\ &\quad (P(\mathbf{L} = \ell^*)(P(R_k = 0, L_j = m) - \\ &\quad q_{jm}^k q_{-j\ell}^k (P(L_j = m) + \\ &\quad P(R_k = 1, L_j \neq m)))) \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} \\ &= (q_{jm}^k P(R_k = 1)^2)^{-1} \cdot \\ &\quad (P(\mathbf{L} = \ell^*)(P(R_k = 0, L_j = m) - \\ &\quad q_{jm}^k q_{-j\ell}^k (P(R_k = 0, L_j = m) + P(R_k = 1, L_j = m) + \\ &\quad P(R_k = 1, L_j \neq m)))) \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} \\ &= (q_{jm}^k P(R_k = 1)^2)^{-1} \cdot \\ &\quad (P(\mathbf{L} = \ell^*)(P(R_k = 0, L_j = m) - \\ &\quad q_{jm}^k q_{-j\ell}^k (P(R_k = 0, L_j = m) \\ &\quad + P(R_k = 1)))) \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} \\ &= (q_{jm}^k P(R_k = 1)^2)^{-1} \cdot \\ &\quad (P(\mathbf{L} = \ell^*)((1 - q_{jm}^k q_{-j\ell}^k)P(R_k = 0, L_j = m) - \\ &\quad q_{jm}^k q_{-j\ell}^k P(R_k = 1))) \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} \\ &= (q_{jm}^k P(R_k = 1)^2)^{-1} \cdot \\ &\quad (P(\mathbf{L} = \ell^*, R_k = 1)P(R_k = 0, L_j = m) - \\ &\quad P(\mathbf{L} = \ell^*)q_{jm}^k q_{-j\ell}^k P(R_k = 1))) \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} . \end{aligned}$$

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

$$\begin{aligned} &\frac{\partial}{\partial q_{jm}^k} \log(P(\mathbf{L} = \ell^* \mid \mathbf{R} = \mathbf{1}, \mathbf{x})) \\ &= \frac{P(L_j = m \mid R_k = 0)(1 - P(R_k = 1))}{q_{jm}^k P(R_k = 1)} \cdot \\ &\quad \frac{q_{-j\ell^*}^k}{1 - q_{jm}^k q_{-j\ell^*}^k} \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} . \end{aligned}$$

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

1.2 MISSING LABELS

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

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

$$\begin{aligned} & \frac{\partial}{\partial q_{jm}^k} P(\mathbf{L}' = \ell' | \mathbf{R} = \mathbf{1}, \mathbf{x}) \\ &= \frac{\partial}{\partial q_{jm}^k} \frac{\sum_{\ell^0 \in \mathcal{L}^0} P(R_k = 1, \mathbf{L}' = \ell', \mathbf{L}^0 = \ell^0 | \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x})}{P(R_k = 1 | \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x})} \end{aligned}$$

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 1.1. When taking the logarithm, the gradient is divided by $P(\mathbf{L}' = \ell' | \mathbf{R} = \mathbf{1}, \mathbf{x})$, so that the resulting gradient remains the same as in Section 1.1:

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

Case 2: Category is known and correct. In the case that m is known and the true category, that is $L_j \in \mathbf{L}'$ and $\ell_j = m$, the sum can again be factored out. Thus, we can again make use of the gradients from Section 1.1:

$$\begin{aligned} & \frac{\partial}{\partial q_{jm}^k} P(\mathbf{L}' = \ell' | \mathbf{R} = \mathbf{1}, \mathbf{x}) \\ &= \sum_{\ell^0 \in \mathcal{L}^0} \frac{\partial}{\partial q_{jm}^k} P(\mathbf{L}' = \ell', \mathbf{L}^0 = \ell^0 | \mathbf{R} = \mathbf{1}, \mathbf{x}) \\ &= (q_{jm}^k P(R_k = 1)^2)^{-1} \cdot \\ & \quad (P(\mathbf{L}' = \ell', R_k = 1)P(R_k = 0, L_j = m) - \\ & \quad P(\mathbf{L}' = \ell', R_k = 0)P(R_k = 1)) \Big|_{\mathbf{R}_{-k} = \mathbf{1}, \mathbf{x}} . \end{aligned}$$

When taking the logarithm, the expression gets slightly more complicated than in Section 1.1:

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

So, one additional marginal query has to be calculated per rule. Note that this expression simplifies to that of Section 1.1 if all labels are known, that is if $\mathbf{L}' = \mathbf{L}$, because the conditional probabilities in the second fraction are then the rules conditional probability tables, given by \mathbf{q} alone.

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

$$\begin{aligned} & \frac{\partial}{\partial q_{jm}^k} P(\mathbf{L}' = \ell' | \mathbf{R} = \mathbf{1}, \mathbf{x}) \\ &= \frac{\partial}{\partial q_{jm}^k} \sum_{\ell^0 \in \mathcal{L}^0} \left((1 - q_{j\ell'_j}^k q_{-j\ell'^*}^k) \cdot \frac{P(\mathbf{L}' = \ell', \mathbf{L}^0 = \ell^0)}{P(R_k = 1)} \right) \Big|_{\mathbf{R}_{-k} = \mathbf{1}, \mathbf{x}} \end{aligned}$$

$$\begin{aligned} &= \frac{\partial}{\partial q_{jm}^k} \left(\sum_{\ell^0 \in \mathcal{L}^0} P(\mathbf{L}' = \ell', \mathbf{L}^0 = \ell^0) - \sum_{\ell^0 \in \mathcal{L}_{-j}^0} q_{j\ell'_j}^k q_{-j\ell'^*}^k P(\mathbf{L}' = \ell', \mathbf{L}^0 = \ell^0) - q_{jm}^k \sum_{\ell^0 \in \mathcal{L}_{-j}^0} q_{-j\ell'^*}^k P(\mathbf{L}' = \ell', \mathbf{L}^0 = \ell^0) \right) \cdot \\ & \quad \left(\sum_{\ell \in \mathcal{L}} P(\mathbf{L} = \ell) - \sum_{\ell \in \mathcal{L}_{j \neq m}} q_{j\ell'_j}^k q_{-j\ell}^k P(\mathbf{L} = \ell) - q_{jm}^k \sum_{\ell \in \mathcal{L}_{j=m}} q_{-j\ell}^k P(\mathbf{L} = \ell) \right)^{-1} \Big|_{\mathbf{R}_{-k} = \mathbf{1}, \mathbf{x}} . \end{aligned}$$

Again, we substitute to make differentiation easier to see:

$$\begin{aligned} &= \frac{\partial}{\partial q_{jm}^k} \frac{c_1 - q_{jm}^k c_2}{c_3 - q_{jm}^k c_4} \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} \\ &= \frac{c_1 c_4 - c_2 c_3}{(c_3 - q_{jm}^k c_4)^2} \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} . \end{aligned}$$

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

$$\begin{aligned} &= (P(R_k = 1)^2)^{-1} \cdot \\ & \quad (P(\mathbf{L}' = \ell') - P(L_j \neq m, \mathbf{L}' = \ell', R_k = 0)) \cdot \\ & \quad (q_{jm}^k)^{-1} P(L_j = m, R_k = 0) - (q_{jm}^k)^{-1} \cdot \\ & \quad P(L_j = m, \mathbf{L}' = \ell', R_k = 0) \cdot (P(L_j = m) + \\ & \quad P(L_j \neq m) - P(L_j \neq m, R_k = 0)) \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} \\ &= (q_{jm}^k P(R_k = 1)^2)^{-1} \cdot \\ & \quad (P(L_j = m, \mathbf{L}' = \ell') + \\ & \quad P(L_j \neq m, \mathbf{L}' = \ell') - P(L_j \neq m, \mathbf{L}' = \ell', R_k = 0)) \cdot \\ & \quad (P(L_j = m, R_k = 0) - P(L_j = m, \mathbf{L}' = \ell', R_k = 0) \cdot \\ & \quad (P(L_j = m) + P(L_j \neq m, R_k = 1))) \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} \\ &= (q_{jm}^k P(R_k = 1)^2)^{-1} \cdot \\ & \quad (P(L_j = m, \mathbf{L}' = \ell', R_k = 0) + \\ & \quad P(L_j = m, \mathbf{L}' = \ell', R_k = 1) + \\ & \quad P(L_j \neq m, \mathbf{L}' = \ell', R_k = 1)) \cdot \\ & \quad (P(L_j = m, R_k = 0) - P(L_j = m, \mathbf{L}' = \ell', R_k = 0) \cdot \\ & \quad (P(L_j = m, R_k = 0) + P(L_j = m, R_k = 1) + \\ & \quad P(L_j \neq m, R_k = 1))) \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} \\ &= (q_{jm}^k P(R_k = 1)^2)^{-1} \cdot \\ & \quad (P(L_j = m, \mathbf{L}' = \ell', R_k = 1)P(L_j = m, R_k = 0) + \\ & \quad P(L_j = m, \mathbf{L}' = \ell', R_k = 0)P(L_j = m, R_k = 0) + \\ & \quad P(L_j \neq m, \mathbf{L}' = \ell', R_k = 1)P(L_j = m, R_k = 0) - \\ & \quad P(L_j = m, \mathbf{L}' = \ell', R_k = 0)P(L_j = m, R_k = 0) - \\ & \quad P(L_j = m, \mathbf{L}' = \ell', R_k = 0)P(L_j = m, R_k = 1) - \\ & \quad P(L_j = m, \mathbf{L}' = \ell', R_k = 0) \cdot \\ & \quad P(L_j \neq m, R_k = 1)) \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} \\ &= (q_{jm}^k P(R_k = 1)^2)^{-1} \cdot \\ & \quad (P(\mathbf{L}' = \ell', R_k = 1)P(L_j = m, R_k = 0) - \\ & \quad P(L_j = m, \mathbf{L}' = \ell', R_k = 0) \cdot \\ & \quad P(R_k = 1)) \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} . \end{aligned}$$

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

$$\begin{aligned} &= \frac{\partial}{\partial q_{jm}^k} \log(P(\mathbf{L}' = \ell' \mid \mathbf{R} = \mathbf{1}, \mathbf{x})) \\ &= \frac{P(L_j = m, R_k = 0)}{q_{jm}^k P(R_k = 1)} - \\ & \quad \frac{P(L_j = m, R_k = 0 \mid \mathbf{L}' = \ell')}{q_{jm}^k P(R_k = 1 \mid \mathbf{L}' = \ell')} \mid \mathbf{R}_{-k} = \mathbf{1}, \mathbf{x} . \end{aligned}$$

At this point it can be seen that the gradient reduces to 0 if no label has a ground truth, that is $L' = \emptyset$, which makes intuitive sense. Applying one last transformation gives a form that is easier to compute:

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

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

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, \dots, L_J . Each variable can only set the gate R_k to $R_k = 1$ with a probability $1 - q_{j1}^k$ if the variable itself is $L_j = 1$. Thus the conditional probability distribution that defines R_k is:

$$P(R_k = 0 \mid L_1 = \ell_1, \dots, L_J = \ell_J) = \prod_{j=1}^J (q_{j1}^k)^{\ell_j} .$$

In our case, the input variables may have multiple categories $m = 1, \dots, M(j)$ and each of these categories can trigger the gate to be $R_k = 1$ with a probability $1 - q_{jm}^k$. 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

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

These variables can then be connected to the ordinary binary-input noisy-or gate as visualized in Figure 1 with their corresponding inhibition probability q_{jm}^k . 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.

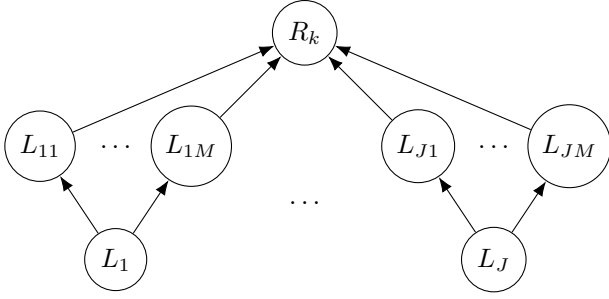


Figure 1: Decomposition of Multicategorical Inputs for a Binary-Input Noisy-Or Gate.

Algorithm 1 Simulation Dataset Generation

Require: nLabels, nRules, nData

- 1: pRSL \leftarrow empty model
- 2: **for** i in $1, \dots, \text{nLabels}$ **do**
- 3: nCategories $\sim U\{2, \dots, 4\}$
- 4: Add label node with nCategories to pRSL
- 5: **for** i in $1, \dots, \text{nRules}$ **do**
- 6: nCategories $\sim U\{2, \dots, 5\}$
- 7: categories \leftarrow Draw nCategories from all
 $\{1, \dots, M(1), \dots, 1, \dots, M(J)\}$
- 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 $1, \dots, \text{nData}$ **do**
- 12: **for** each label node in pRSL **do**
- 13: classifierOutput[label][i] $\sim \text{Dir}(1)$
- 14: **return** pRSL, classifierOutput

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

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 $j = 1, \dots, J$. 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:

$$P(R_k = 0 | \mathbf{L} = \ell) = \prod_{j=1}^J q_j^k \ell_j .$$

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.

References

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

Table 1: Performance of pRSL on Train, Validation, and Test Data. Mean \pm Standard Deviation Between Folds.

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