

A family of admissible heuristics for A* to perform inference in probabilistic classifier chains

Deiner Mena^{1,2} · Elena Montañés¹ ·
José Ramón Quevedo¹ · Juan José del Coz¹

Received: 14 September 2015 / Accepted: 15 September 2016 / Published online: 26 September 2016
© The Author(s) 2016

Abstract Probabilistic classifier chains have recently gained interest in multi-label classification, due to their ability to optimally estimate the joint probability of a set of labels. The main hindrance is the excessive computational cost of performing inference in the prediction stage. This pitfall has opened the door to propose efficient inference alternatives that avoid exploring all the possible solutions. The ϵ -approximate algorithm, beam search and Monte Carlo sampling are appropriate techniques, but only ϵ -approximate algorithm with $\epsilon = 0$ theoretically guarantees reaching an optimal solution in terms of subset 0/1 loss. This paper offers another alternative based on heuristic search that keeps such optimality. It consists of applying the A* algorithm providing an admissible heuristic able to explore fewer nodes than the ϵ -approximate algorithm with $\epsilon = 0$. A preliminary study has already coped with this goal, but at the expense of the high computational time of evaluating the heuristic and only for linear models. In this paper, we propose a family of heuristics defined by a parameter that controls the trade-off between the number of nodes explored and the cost of computing the heuristic. Besides, a certain value of the parameter provides a method that is also suitable for the non-linear case. The experiments reported over several benchmark datasets show that the number of nodes explored remains quite steady for different values of the parameter, although the time considerably increases for high values. Hence, low values of the parameter

Editor: Eyke Hüllermeier.

✉ Elena Montañés
elena@aic.uniovi.es

Deiner Mena
deiner@aic.uniovi.es; deiner.mena@utch.edu.co

José Ramón Quevedo
quevedo@aic.uniovi.es

Juan José del Coz
juanjo@aic.uniovi.es

¹ Artificial Intelligence Center, University of Oviedo, Campus de Viesques, s/n,
33204 Gijón, Asturias, Spain

² Cra. 22 No 18B-10 B/ Nicolás Medrano - Ciudadela Universitaria, Quibdó - Chocó, Colombia

give heuristics that theoretically guarantee exploring fewer nodes than the ϵ -approximate algorithm with $\epsilon = 0$ and show competitive computational time. Finally, the results exhibit the good behavior of the A* algorithm using these heuristics in complex situations such as the presence of noise.

Keywords Multi-label classification · Probabilistic classifier chains · Inference · A* · Admissible heuristics

1 Introduction

Multi-label classification (MLC) is a machine learning problem in which models are sought that assign a subset of (classes) labels to each instance, unlike conventional (single-class) classification that involves predicting only a single class. Multi-label classification problems are ubiquitous and naturally occur, for instance, in assigning keywords to a paper, tags to resources in a social network, objects to images or emotional expressions to human faces.

In general, the problem of multi-label learning comes with two fundamental challenges. The first refers to the computational complexity of the algorithms. If the number of labels is large, then a complex approach might not be applicable in practice. Therefore, the scalability of algorithms is a key issue in this field. The second problem is related to the 'own nature' of multi-label data. Not only is the number of labels typically large, but each instance also belongs to a variable-sized subset of labels simultaneously. Moreover, and perhaps even more importantly, the labels will normally not occur independently of each other; instead, there are statistical dependencies between them. From a learning and prediction point of view, these relationships constitute a promising source of information, in addition to that coming from the mere description of the instances. Thus, it is hardly surprising that research on MLC has very much focused on the design of new methods that are able to detect—and benefit from—interdependencies among labels.

Several approaches have been proposed in the literature to cope with MLC. Firstly, researchers tried to adapt and extend different state-of-the-art binary or multi-class classification algorithms (Elisseeff and Weston 2005; McCallum 1999; Zhang and Zhou 2007). Secondly, they further analyzed in depth the label dependence and attempt to design new approaches that exploit label correlations (Dembczyński et al. 2012). In this regard, two kinds of label dependence have been formally distinguished, namely, conditional dependence (Dembczyński et al. 2010; Montañés 2011, 2014; Read et al. 2011) and marginal (unconditional) dependence (Cheng and Hüllermeier 2009). Also, pairwise relations (Elisseeff and Weston 2005), relations in sets of different sizes (Read et al. 2011; Tsoumakas and Vlahavas 2007), or relations in the whole set of labels (Cheng and Hüllermeier 2009; Montañés 2011) have also been exploited.

Regarding conditional label dependence, the approach called probabilistic classifier chains (PCC) has aroused great interest among the multi-label community, since it offers the excellent property of being able to estimate the conditional joint distribution of the labels. However, the original PCC algorithm (Dembczyński et al. 2010) suffers from high computational cost, since it performs an exhaustive search as inference strategy to obtain optimal solutions in terms of a given loss function. Several efforts that use different searching and sampling strategies in order to overcome this drawback are being made currently. This includes uniform-cost search (Dembczynski et al. 2012), beam search (Kumar et al. 2012, 2013) and Monte Carlo sampling (Dembczynski et al. 2012; Read et al. 2014). All of these algorithms successfully

estimate an optimal solution reached by the original PCC (Dembczyński et al. 2010), at the same time that they reduce the computational cost in terms of both the number of candidate solutions explored and execution time. The main contribution of this paper is to propose an alternative based on an heuristic search strategy. In particular, the proposal consists of obtaining admissible heuristics for the well-known A* algorithm (Hart et al. 1968). In this respect, we have already started to fill this gap in the literature with a recent published preliminary work (Mena et al. 2015), concluding that the proposal guarantees, not only optimal predictions in terms of subset 0/1 loss, but also that it explores fewer nodes than all previous methods that also provide optimal predictions. Unfortunately, and after studying in depth this heuristic, two main drawbacks can be stated: (i) it could only be defined for linear models and (ii) its computation is moderately high. In this direction, the goal of this paper is twofold. On the one hand, this work goes further by defining a family of heuristics through a parameter that controls the trade-off between the number of nodes explored and the cost of computing the heuristic. Besides, a special value of the parameter leads to an heuristic suitable for non-linear models. On the other hand, this work also studies and analyzes situations in which the computation of the heuristic compensates the whole computational cost, showing a steady behavior with regard to other algorithms. All these methods are analyzed and experimentally compared over a wide range of multi-label datasets.

The rest of the paper is organized as follows. Section 2 formally describes the multi-label framework and the principles of PCC. Section 3 describes and discusses the properties and behavior of the different state-of-the-art approaches for inference in PCCs. Section 4 details the heuristic search framework and defines admissible heuristics for the A* algorithm. Exhaustive experiments are shown and discussed in Sect. 5. Finally, Sect. 6 offers some conclusions and includes new directions for future work.

2 Probabilistic classifier chains in multi-label classification

This section formally describes the MLC task and the PCC methods.

2.1 Formal settings of multi-label classification and loss functions

Let be $\mathcal{L} = \{\ell_1, \ell_2, \dots, \ell_m\}$ a finite and non-empty set of m labels and $S = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$ a training set independently and randomly drawn according to an unknown probability distribution $\mathbf{P}(\mathbf{X}, \mathbf{Y})$ on $\mathcal{X} \times \mathcal{Y}$, where \mathcal{X} and \mathcal{Y} are the input and the output space, respectively. The former is the space of the instance description, whereas the latter is given by the power set $\mathcal{P}(\mathcal{L})$ of \mathcal{L} . To ease notation, we define \mathbf{y}_i as a binary vector $\mathbf{y}_i = (y_{i,1}, y_{i,2}, \dots, y_{i,m})$ in which $y_{i,j} = 1$ indicates the presence (relevance) and $y_{i,j} = 0$ the absence (irrelevance) of ℓ_j in the labeling of \mathbf{x}_i . Hence, \mathbf{y}_i is the realization of a corresponding random vector $\mathbf{Y} = (\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_m)$. Using this convention, the output space can also be defined as $\mathcal{Y} = \{0, 1\}^m$. The goal in MLC is to induce from S a hypothesis $f : \mathcal{X} \rightarrow \mathcal{Y}$ that minimizes the risk in terms of certain loss function $L(\cdot)$ when it provides a vector of relevant labels $\mathbf{y} = f(\mathbf{x}) = (f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_m(\mathbf{x}))$ for unlabeled query instances \mathbf{x} . This risk can be defined as the expected loss over the joint distribution $\mathbf{P}(\mathbf{X}, \mathbf{Y})$, that is,

$$R_L(f) = \mathbb{E}_{\mathbf{X}, \mathbf{Y}} L(\mathbf{Y}, f(\mathbf{X})). \quad (1)$$

therefore, denoted by $\mathbf{P}(\mathbf{y} | \mathbf{x})$ the conditional distribution $\mathbf{Y} = \mathbf{y}$ given $\mathbf{X} = \mathbf{x}$, then the so-called risk minimizer f^* can be expressed by

$$f^*(\mathbf{x}) = \arg \min_f \sum_{y \in \mathcal{Y}} \mathbf{P}(y | \mathbf{x}) L(y, f(\mathbf{x})). \tag{2}$$

Let us comment that the conditional distribution $\mathbf{P}(y | \mathbf{x})$ presents different properties which are crucial for optimizing different loss functions. In this respect, the strategy followed by a certain MLC algorithm for modeling label dependence determines the loss function that is optimized. But, unfortunately, for most of the algorithms it is quite complex and confusing to discover the loss function they attempt to optimize.

With regard to the loss functions, several performance measures have been taken for evaluating MLC. The most specific ones are the subset 0/1 loss and the Hamming loss, but there exist other measures that have been taken from other research fields, as such F_1 or the *Jaccard* index. Here, we will focus on just the subset 0/1 loss, since it is the measure PCCs are able to optimize. The subset 0/1 loss checks if the predicted and relevant label subsets are equal or not and is defined by¹

$$L_{S_{0/1}}(y, f(\mathbf{x})) = \mathbb{I}[y \neq f(\mathbf{x})]. \tag{3}$$

In the case of this evaluation measure, it is sufficient to take the mode of the entire joint conditional distribution for optimizing this loss. Formally, the risk minimizer adopts the following simplified form

$$f^*(\mathbf{x}) = \arg \max_{y \in \mathcal{Y}} \mathbf{P}(y | \mathbf{x}). \tag{4}$$

2.2 Probabilistic classifier chains

PCCs (Dembczyński et al. 2010) [such as CC (Read et al. 2009, 2011)] are based on learning a chain of classifiers. These methods take an order of the label set and train a probabilistic binary classifier for estimating $\mathbf{P}(y_j | \mathbf{x}, y_1, \dots, y_{j-1})$ for each label ℓ_j following this order. Hence, the probabilistic model obtained for predicting label ℓ_j , denoted by f_j is of the form

$$f_j : \mathcal{X} \times \{0, 1\}^{j-1} \longrightarrow [0, 1]. \tag{5}$$

The training data for this classifier is the set $S_j = \{(\bar{\mathbf{x}}_1, y_{1,j}), \dots, (\bar{\mathbf{x}}_n, y_{n,j})\}$ where $\bar{\mathbf{x}}_i = (\mathbf{x}_i, y_{i,1}, \dots, y_{i,j-1})$, that is, the features are \mathbf{x}_i supplemented by the relevance of the labels $\ell_1, \dots, \ell_{j-1}$ preceding ℓ_j in the chain and the category is the relevance of the label ℓ_j .

In the testing stage of the methods based on learning a chain of classifiers, the goal is to perform inference for each instance, which consists of estimating the risk minimizer for a given loss function over the estimated entire joint conditional distribution. The idea revolves around repeatedly applying the general product rule of probability to the joint distribution of the labels $\mathbf{Y} = (\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_m)$, that is, computing

$$\mathbf{P}(\mathbf{y} | \mathbf{x}) = \prod_{j=1}^m \mathbf{P}(y_j | \mathbf{x}, y_1, \dots, y_{j-1}). \tag{6}$$

Before analyzing this issue in the next section, note that from a theoretical point of view, this expression holds for any order considered for the labels. But, in practice, these methods are label order dependent for several reasons. On the one hand, it is not possible to assure that the models obtained in the training stage perfectly estimate the joint conditional probability $P(\mathbf{y} | \mathbf{x})$. On the other hand, predicted values instead of true values are successively taken in

¹ For a predicate p , the expression $\mathbb{I}[p]$ evaluates to 1 if p is true and 0 otherwise.

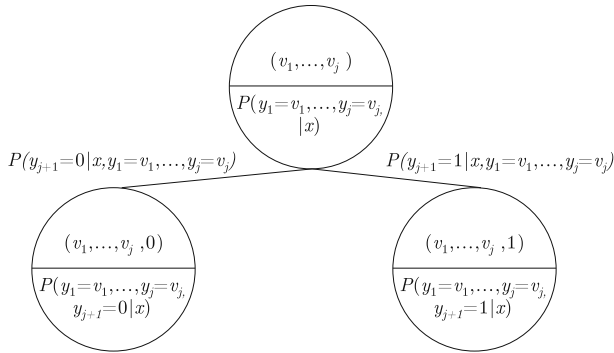


Fig. 1 A generic node and its children of the probability binary tree. The *top part* of each node contains the combination of labels and the *bottom part* includes the joint probability of such a combination. The *edges* are labeled with the conditional probability

the testing stage. This is more serious if the highest errors occur at the beginning of the chain, since error predictions are successively propagated (Montañés 2014; Senge et al. 2012a, b). In any case, in this paper we assume the order of the labels in the chain to be given, since the goal is just to analyze the performance of the methods, without taking into account the effect of different orders. Hence, we do not include any study about which order can be the best.

Before going into the detailed description of the state-of-the-art of the inference approaches and of our proposal, we note that the training phase is common to all of them, thus, the models f_j induced by the binary classifiers will be the same. So, in what follows we will focus just on the testing stage.

3 Inference in probabilistic classifier chains

First of all, let us consider the task of performing different inference procedures as different manners of exploring a probability binary tree in order to facilitate the explanation and analysis of the inference approaches in the next section. In such a tree, the k -th node of level $j < m$ with $k \leq 2^j$ is labeled by $y_j^k = (v_1, v_2, \dots, v_j)$ with $v_i \in \{0, 1\}$ for $i = 1, \dots, j$. This node has two children respectively labeled as $y_{j+1}^{2k-1} = (v_1, v_2, \dots, v_j, 0)$ and $y_{j+1}^{2k} = (v_1, v_2, \dots, v_j, 1)$ and with marginal joint conditional probability $\mathbf{P}(y_1 = v_1, \dots, y_j = v_j, y_{j+1} = 0 | \mathbf{x})$ and $\mathbf{P}(y_1 = v_1, \dots, y_j = v_j, y_{j+1} = 1 | \mathbf{x})$. The weights of the edges between both parent and children are respectively $\mathbf{P}(y_{j+1} = 0 | \mathbf{x}, y_1 = v_1, \dots, y_j = v_j)$ and $\mathbf{P}(y_{j+1} = 1 | \mathbf{x}, y_1 = v_1, \dots, y_j = v_j)$, which are respectively estimated by $1 - f_{j+1}(\mathbf{x}, v_1, \dots, v_j)$ and $f_{j+1}(\mathbf{x}, v_1, \dots, v_j)$. The marginal joint conditional probability of the children is computed by the product rule of probability. Then, $\mathbf{P}(y_1 = v_1, \dots, y_j = v_j, y_{j+1} = 0 | \mathbf{x}) = \mathbf{P}(y_{j+1} = 0 | \mathbf{x}, y_1 = v_1, \dots, y_j = v_j) \cdot \mathbf{P}(y_1 = v_1, \dots, y_j = v_j | \mathbf{x})$ and $\mathbf{P}(y_1 = v_1, \dots, y_j = v_j, y_{j+1} = 1 | \mathbf{x}) = \mathbf{P}(y_{j+1} = 1 | \mathbf{x}, y_1 = v_1, \dots, y_j = v_j) \cdot \mathbf{P}(y_1 = v_1, \dots, y_j = v_j | \mathbf{x})$. The root node is labeled by the empty set. Figure 1 illustrates this.

Several approaches have been proposed for inference in PCCs. The method the first proposed is the one based on greedy search (GS), being the integral part of the original CC method (Read et al. 2009). Its successor is the exhaustive search (ES), called the PCC method (Dembczyński et al. 2010). The ϵ -approximate (ϵ -A) algorithm (Dembczynski et al. 2012) is a

uniform-cost (UC) search algorithm that can output optimal predictions in terms of subset 0/1 loss and also reduces significantly the computational cost of ES. A more recent approach based on beam search (Kumar et al. 2012, 2013) (BS) presents good behavior both in terms of performance and computational cost. Finally, Monte Carlo sampling (Dembczynski et al. 2012) is an appealing and simpler alternative (Dembczynski et al. 2012; Read et al. 2014) to overcome the high computational cost of ES.

This section copes with the particularities of these inference methods, except the Monte Carlo sampling approaches. We have already studied Monte Carlo approaches (Dembczynski et al. 2012; Read et al. 2014) and compared them with the ϵ -A algorithm and BS techniques (Mena et al. 2015). The conclusions of that work were that, although they are well suited for minimization of other losses, e.g., Hamming loss or example-based F-measure, (i) they need to explore many more nodes to be closer to the optimal, despite the fact that their performance could sometimes be better, (ii) they enlarge as the size of the sample drawn grows and (iii) they are quite slow even for low values of the sample. Hence, we do not consider them as competitive methods in the present work, although they could sometimes be appealing.

3.1 Greedy search

At the testing stage, the GS strategy, originally called CC (Read et al. 2009, 2011), provides an output $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_m)$ for a new unlabeled instance \mathbf{x} by successively querying each classifier f_j that estimates the conditional probability $\mathbf{P}(y_j | \mathbf{x}, y_1, \dots, y_{j-1})$. This means exploring just one node in each level j . Given that only the two children of the explored node in level j are taken, their marginal joint conditional probability only differs in the marginal conditional probability $\mathbf{P}(y_j | \mathbf{x}, y_1, \dots, y_{j-1})$, since both children have the same parent. Thus, the path selected is that of the child with the highest marginal conditional probability $\mathbf{P}(y_j | \mathbf{x}, y_1, \dots, y_{j-1})$ and the prediction for an instance \mathbf{x} is of the form

$$\hat{\mathbf{y}} = (f_1(\mathbf{x}), f_2(\mathbf{x}, f_1(\mathbf{x})), f_3(\mathbf{x}, f_1(\mathbf{x}), f_2(\mathbf{x}, f_1(\mathbf{x}))), \dots). \quad (7)$$

Figure 2a shows the path followed by an instance using this strategy. In this example, only the right node is explored in each level. The optimal solution is not reached, since the optimal solution is that which ends in the sixth leaf, whereas the method falls in the last leaf.

Concerning the optimization of subset 0/1 loss, a rigorous analysis (Dembczynski et al. 2012) establishes bounds for the performance of the GS, showing the poor performance of the method for this loss, although it tends to optimize it.

3.2 Exhaustive search

Unlike GS that explores only a single label combination, ES estimates the entire joint conditional distribution $\mathbf{P}(\cdot | \mathbf{x})$ for a new unlabeled instance \mathbf{x} , since it provides a Bayes optimal inference. Hence, it explores all possible paths in the tree. Then, for each $h(\mathbf{x})$, it computes $\mathbf{P}(\mathbf{y} | \mathbf{x})$ and $L(\mathbf{y}, \mathbf{f}(\mathbf{x}))$ for all combination of labels $\mathbf{y} = (y_1, y_2, \dots, y_m)$ and outputs the combination $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_m) = \mathbf{f}^*(\mathbf{x})$ with minimum risk for the given loss $L(\cdot, \cdot)$. By doing so, it generally improves in terms of performance, since it perfectly estimates the risk minimizer, albeit at the cost of a higher computational cost, as it comes down to summing over an exponential (2^m) number of label combinations for each $h(\mathbf{x})$.

Figure 2b illustrates this approach where, by exploring all paths, the optimal solution is always reached.

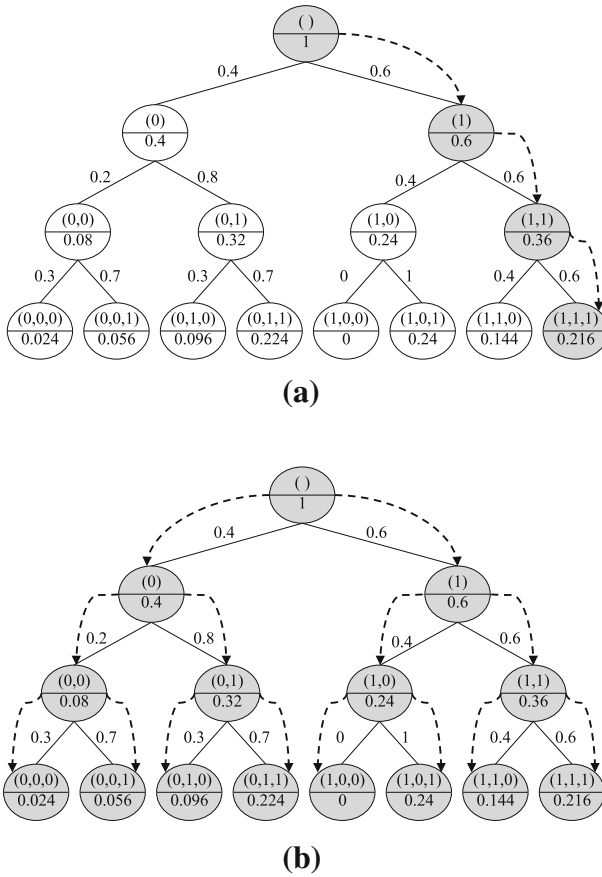


Fig. 2 An example of paths followed by an instance using **a** Greedy Search (CC) and **b** exhaustive Search (PCC). The *dotted arrows* show the path followed by the algorithm

3.3 ϵ -approximate algorithm

The ϵ -Approximate (ϵ -A) algorithm (Dembczynski et al. 2012) arises as an alternative to the high computational cost of ES and to the poor performance of CC. In terms of the probability tree defined above, it expands only the nodes whose marginal joint conditional probability exceeds the threshold $\epsilon = 2^{-k}$ with $1 \leq k \leq m$ (notice that $\epsilon = 0$ and all values of ϵ between 0 and 2^{-m} are in fact the same case of $\epsilon = 2^{-m}$). This marginal joint conditional probability for a node in level j , which deals with the label ℓ_j and for an unlabeled instance \mathbf{x} is

$$\mathbf{P}(y_1, \dots, y_j | \mathbf{x}) = \prod_{i=1}^j \mathbf{P}(y_i | \mathbf{x}, y_1, \dots, y_{i-1}), \tag{8}$$

where $\mathbf{P}(y_i | \mathbf{x}, y_1, \dots, y_{i-1})$ is estimated by $f_i(\mathbf{x}, y_1, \dots, y_{i-1})$.

The nodes are expanded in the order established by this probability, calculating the marginal joint conditional probability for their children. So, the algorithm does not follow a specific path, otherwise it changes from one path to another depending on the marginal joint

conditional probabilities. In the end, two situations can be found: (i) the node expanded is a leaf or (ii) there are no more nodes that exceed the threshold. If the former situation happens, the prediction for the unlabeled instance \mathbf{x} will be $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_m)$ corresponding to the combination of the leaf reached (see Fig. 3a). Conversely, if situation (ii) takes place, then GS is applied to the nodes whose children do not exceed the threshold, and the prediction $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_m)$ for the unlabeled instance \mathbf{x} in this case will be that with the highest entire joint conditional probability $\mathbf{P}(y_1, \dots, y_m | \mathbf{x})$ (see Fig. 3b, c).

The parameter ϵ plays an important role in the algorithm. The particular case of $\epsilon = 0$ (or any value in the interval $[0, 2^{-m}]$, that is, $k = m$) is of special interest, since the algorithm performs a UC that always finds the optimal solution. Figure 3a illustrates this situation.

Conversely, the method is looking to GS as ϵ grows, being the GS in the case of $\epsilon = 0.5$ (or equivalently $\epsilon = 2^{-1}$, that is, $k = 1$). This is so because in this case two situations are possible: (i) only one node has a marginal joint conditional probability greater than ϵ , in which case the algorithm follows one path, or (ii) no nodes have a marginal joint conditional probability greater than ϵ , in which case a GS is applied from here to the bottom of the tree. Figure 3b shows an example of this particular case.

Notice that this method provides an optimal prediction if the entire joint conditional probability of the corresponding label combination is greater than ϵ . The interpretation of the method for a generic value of $\epsilon = 2^{-k}$ is that the method guarantees reaching a partial optimal solution at least until the $k - th$ level on the tree. Figure 3c shows the particular case of $\epsilon = 0.25$.

Consequently, this algorithm estimates the risk minimizer for the subset 0/1 loss to a greater or lesser extend, depending on the value of ϵ . Moreover, a theoretical analysis of this estimation (Dembczynski et al. 2012) allows bounding its goodness as a function of the number of iterations, which in turn depends on ϵ .

3.4 Beam search

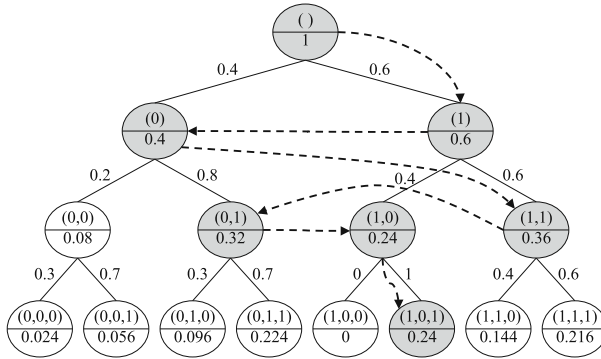
Beam Search (BS) (Kumar et al. 2012, 2013) also explores more than one path in the probabilistic tree. This method includes a parameter b called the beam width that limits the number of combinations of labels explored. The idea is to explore b possible candidate sets of labels at each level of the tree. Hence, depending on such a value, a certain number of the top levels are exhaustively explored, particularly a total of $k^* - 1$ levels, k^* being the lowest integer such that $b < 2^{k^*}$. Then, only b number of possibilities are explored for each of the remainder levels. The combinations explored from the level k^* to the bottom are those with the highest marginal joint conditional probability seen thus far. This marginal joint conditional probability for a node of level j for an unlabeled instance \mathbf{x} is the same as for the ϵ -A algorithm. Hence, such a probability is

$$\mathbf{P}(y_1, \dots, y_j | \mathbf{x}) = \prod_{i=1}^j \mathbf{P}(y_i | \mathbf{x}, y_1, \dots, y_{i-1}), \tag{9}$$

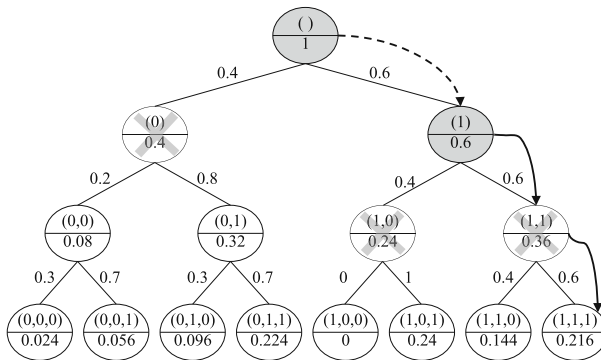
where $\mathbf{P}(y_i | \mathbf{x}, y_1, \dots, y_{i-1})$ is estimated by $f_i(\mathbf{x}, y_1, \dots, y_{i-1})$.

In the end, the algorithm outputs $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_m)$ with the highest entire joint conditional probability $\mathbf{P}(y_1, \dots, y_m | \mathbf{x})$.

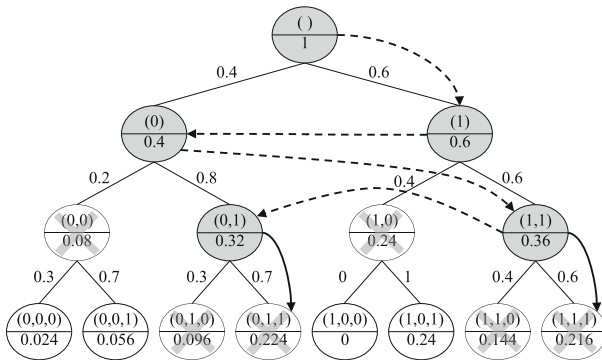
BS differs from GS in that (i) BS explores more than one combination whereas GS just explores one and also in (ii) the probability taken for deciding a path to follow in the tree is different from one method to the other. Concerning (ii), both take the marginal joint conditional probability $\mathbf{P}(y_1, \dots, y_j | \mathbf{x})$, but in the case of GS, that is equivalent to taking



(a)



(b)



(c)

Fig. 3 Several examples of the paths followed by the ϵ -A algorithm for different values of ϵ . The nodes with a cross are those that have a marginal joint conditional probability lower than ϵ and, hence, they are not explored any more. The dotted arrows show the path followed by the algorithm. The solid arrows indicate the path followed by the algorithm when the marginal joint conditional probability does not exceed the value of ϵ , and, hence a GS is applied to this node from here to the bottom of the tree. **a** ϵ -A algorithm with $\epsilon = 0$ ($k = m$), **b** ϵ -A algorithm with $\epsilon = 0.5$ ($k = 1$) and **c** ϵ -A algorithm with $\epsilon = 0.25$ ($k = 2$)

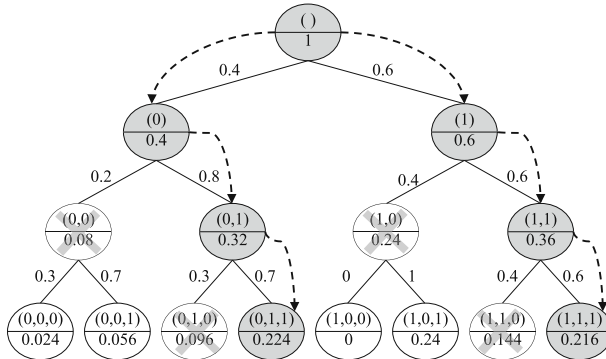


Fig. 4 An example of paths followed by an instance using Beam Search (BS) with $b = 2$. The nodes with a cross mean that this node has none of the highest marginal joint conditional probabilities and, hence, it is not explored any more. The dotted arrows show the path followed by the algorithm

the marginal conditional probability $\mathbf{P}(y_j, | \mathbf{x}, y_1, \dots, y_{j-1})$ since the two nodes explored in each level of the tree have the same parent, as explained before. However, in the case of BS, the b nodes explored do not have to have the same parent, even in the case of $b = 2$. Of course, if $b > 2$, this is impossible to happen in a binary tree.

Let consider the case of $b = 1$. In this case, BS expands just one node in each level, the one with the highest marginal joint conditional probability that coincides with the highest marginal conditional probability, so BS when $b = 1$ follows just one path that coincides with the one followed by GS. Also, if $b = 2^m$, BS performs an ES. Hence, BS encapsulates both GS and ES respectively considering $b = 1$ and $b = 2^m$. This makes it possible to control the trade-off between computational cost and performance of the method by tuning b between 1 and 2^m .

As a final remark, the fact that BS considers marginal joint conditional probabilities makes the method tend to estimate the risk minimizer for the subset 0/1 loss. Kumar et al. (2012) and Kumar et al. (2013) do not include any theoretical analysis about the goodness of the estimation of the risk minimizer, but they empirically show that by taking certain values of b ($b < 15$), the risk minimizer provided by the method converges to the one obtained using ES.

Figure 4 shows an example of the paths explored by the BS algorithm when $b = 2$.

4 A* algorithm for inference in PCC

The algorithm A* is the most widely-known form of best-first search (Pearl 1984), in which the best node at each iteration, according to an evaluation function e , is expanded. The particularity of the A* algorithm is that e can be seen as a function E of the other two functions g and h , $e(k) = E(g(k), h(k))$, where $g(k)$ evaluates the cost of reaching node k from the root and $h(k)$ evaluates the cost of reaching a solution (a leaf) from k . Hence, $e(k)$ evaluates the total cost of reaching a solution from the root through k . In general, it is possible to obtain the exact value of the known information (g), but the unknown information must be estimated through an heuristic (h). To obtain an optimal solution, h must not overestimate the actual cost of reaching a solution, that is, it must be an admissible heuristic. These kinds of heuristics are optimistic, because they estimate that the cost of obtaining a solution is less

than it actually is. Also, the A^* algorithm is optimally efficient for any heuristic, because no other optimal algorithm using the same heuristic guarantees to expand fewer nodes than A^* .

4.1 Building an admissible heuristic

In order to adapt A^* for inference in PCC, we must take into account that we have probabilities instead of costs. So, (1) A^* must select the node with the highest estimated probability, (2) h must not underestimate the probability from the node to a leaf, that is, h must be an admissible heuristic² and (3) E must be the product function, $e = g \cdot h$. Considering all these aspects, e will provide an estimation of the entire joint conditional probability $\mathbf{P}(y_1, \dots, y_m | \mathbf{x})$ for optimizing subset 0/1 loss. In order to derive g and h , let us say that the product rule of probability to the joint distribution of the labels $\mathbf{Y} = (\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_m)$ can be rewritten as

$$\mathbf{P}(y_1, \dots, y_m | \mathbf{x}) = \mathbf{P}(y_1, \dots, y_j | \mathbf{x}) \times \mathbf{P}(y_{j+1}, \dots, y_m | \mathbf{x}, y_1, \dots, y_j). \tag{10}$$

Hence, for a node at level j , let us consider g to be the marginal joint conditional probability $\mathbf{P}(y_1, \dots, y_j | \mathbf{x})$ and h an heuristic that does not underestimate the marginal joint conditional probability $\mathbf{P}(y_{j+1}, \dots, y_m | \mathbf{x}, y_1, \dots, y_j)$. Let us remember that the values of y_1, \dots, y_j are known at level j .

Before discussing the heuristic h proposed here, let us notice that g is the same marginal joint conditional probability that both the ϵ -A algorithm and the BS calculate to select the nodes to be expanded. Even more, ϵ -A with $\epsilon = 0$ is not only equivalent to the UC search, but also to A^* with the constant heuristic $h = 1$. This heuristic is admissible too, since no probability is greater than 1. But, on the other hand, it is also the worst admissible heuristic, since any other admissible heuristic will dominate it³ and consequently the A^* algorithm using such an heuristic will never expand more nodes than the A^* algorithm using $h = 1$.

Let us go now to discuss the heuristic proposed in this paper. Since our heuristic must not underestimate the marginal joint conditional probability $\mathbf{P}(y_{j+1}, \dots, y_m | \mathbf{x}, y_1, \dots, y_j)$ in order to be admissible, it is quite straightforward to pick the maximum value of such probability for obtaining an optimal heuristic h^* :

$$\begin{aligned} h^* &= \max_{(y_{j+1}, \dots, y_m) \in \{0,1\}^{m-j}} \mathbf{P}(y_{j+1}, \dots, y_m | \mathbf{x}, y_1, \dots, y_j) \\ &= \max_{(y_{j+1}, \dots, y_m) \in \{0,1\}^{m-j}} \prod_{i=j+1}^m \mathbf{P}(y_i | \mathbf{x}, y_1, \dots, y_j, y_{j+1}, \dots, y_{i-1}). \end{aligned} \tag{11}$$

However, obtaining such a maximum is, in fact, applying an ES over the set of labels $\mathcal{L} = \{\ell_{j+1}, \dots, \ell_m\}$. Hence, this optimal heuristic is not computationally applicable. So, we need to obtain a tractable heuristic in exchange for renouncing such optimality. This leads to design an heuristic \hat{h} also admissible, but less dominant than h^* ($h^* < \hat{h}$). For this purpose, let us consider $(\bar{y}_{j+1}, \dots, \bar{y}_m)$ the values that define h^* , that is

$$h^* = \prod_{i=j+1}^m \mathbf{P}(y_i | \mathbf{x}, y_1, \dots, y_j, \bar{y}_{j+1}, \dots, \bar{y}_{i-1}), \tag{12}$$

² An heuristic h is admissible for our search problem if and only if it satisfies that $h^*(k) \leq h(k)$ for any node k , where $h^*(k)$ is the highest (and unknown) probability from the node k to a leaf.

³ In our case, an heuristic h_1 dominates another heuristic h_2 (denoted by $h_1 < h_2$) if and only if it satisfies that $h_1(k) \leq h_2(k)$ for any node k .

and let us notice that values $(\bar{y}_{j+1}, \dots, \bar{y}_m)$ which maximize the product do not have to maximize each individual term, then

$$\begin{aligned} & \mathbf{P}(y_i | \mathbf{x}, y_1, \dots, y_j, \bar{y}_{j+1}, \dots, \bar{y}_{i-1}) \\ & \leq \max_{\substack{(y_{j+1}, \dots, y_{i-1}) \\ \in \{0,1\}^{i-1-j}}} \mathbf{P}(y_i | \mathbf{x}, y_1, \dots, y_j, y_{j+1}, \dots, y_{i-1}). \end{aligned} \tag{13}$$

Hence, by defining \hat{h} as

$$\hat{h} = \prod_{i=j+1}^m \max_{\substack{(y_{j+1}, \dots, y_{i-1}) \\ \in \{0,1\}^{i-1-j}}} \mathbf{P}(y_i | \mathbf{x}, y_1, \dots, y_j, y_{j+1}, \dots, y_{i-1}). \tag{14}$$

it is easy to deduce that \hat{h} is admissible and less dominant than h^* ($h^* < \hat{h}$). But again, \hat{h} is not computationally applicable in general. However, this is not the case if we restrict ourselves to the case of linear models, for instance using *logistic regression*. Remember that $\mathbf{P}(y_i | \mathbf{x}, y_1, \dots, y_{i-1})$ is estimated through the model $f_i(\mathbf{x}, y_1, \dots, y_{i-1})$ using a sigmoid function to transform the output of $f_i(\mathbf{x}, y_1, \dots, y_{i-1})$ into a probability. Particularly,

$$\mathbf{P}(y_i = 1 | \mathbf{x}, y_1, \dots, y_{i-1}) = \frac{1}{1 + \exp^{-f_i(\mathbf{x}, y_1, \dots, y_{i-1})}} \tag{15}$$

$$\mathbf{P}(y_i = 0 | \mathbf{x}, y_1, \dots, y_{i-1}) = 1 - \frac{1}{1 + \exp^{-f_i(\mathbf{x}, y_1, \dots, y_{i-1})}}. \tag{16}$$

In order to obtain the maximum value between $\mathbf{P}(y_i = 1 | \mathbf{x}, y_1, \dots, y_{i-1})$ and $\mathbf{P}(y_i = 0 | \mathbf{x}, y_1, \dots, y_{i-1})$, we need first to obtain the maximum value of both terms on their own. In this direction and according to the above expressions of both probabilities, $\mathbf{P}(y_i = 1 | \mathbf{x}, y_1, \dots, y_{i-1})$ will be maximum if $f_i(\mathbf{x}, y_1, \dots, y_{i-1})$ is the maximum and analogously $\mathbf{P}(y_i = 0 | \mathbf{x}, y_1, \dots, y_{i-1})$ will be maximum when $f_i(\mathbf{x}, y_1, \dots, y_{i-1})$ is the minimum. Hence, let us now focus on obtaining the maximum and the minimum of $f_i(\mathbf{x}, y_1, \dots, y_{i-1})$.

From now on, let us consider that f_i is a linear model, then f_i adopts the following form

$$f_i(\mathbf{x}, y_1, \dots, y_{i-1}) = \langle \mathbf{w}_x^i, \mathbf{x} \rangle + \langle \mathbf{w}_y^i, (y_1, \dots, y_{i-1}) \rangle + \beta^i, \tag{17}$$

that splitting the second term in the known part (from ℓ_1 to ℓ_j) and unknown part (from ℓ_{j+1} to ℓ_i) it leads to

$$f_i(\mathbf{x}, y_1, \dots, y_j, y_{j+1}, \dots, y_{i-1}) = \langle \mathbf{w}_x^i, \mathbf{x} \rangle + \sum_{k=1}^j w_{y,k}^i y_k + \sum_{k=j+1}^{i-1} w_{y,k}^i y_k + \beta^i. \tag{18}$$

Since \mathbf{x} is given and y_1, \dots, y_j are fixed, the second summation contains the variables for which the maximum and the minimum must be obtained. Let C_i be the constant part of f_i with regard obtaining the maximum and the minimum, that is,

$$C_i(\mathbf{x}, y_1, \dots, y_j) = \langle \mathbf{w}_x^i, \mathbf{x} \rangle + \sum_{k=1}^j w_{y,k}^i y_k + \beta^i. \tag{19}$$

Then f_i can be rewritten as

$$f_i(\mathbf{x}, y_1, \dots, y_j, y_{j+1}, \dots, y_{i-1}) = C_i(\mathbf{x}, y_1, \dots, y_j) + \sum_{k=j+1}^{i-1} w_{y,k}^i y_k. \tag{20}$$

Consequently, the function whose maximum must be obtained is

$$\sum_{k=j+1}^{i-1} w_{y,k}^i y_k = w_{y,j+1}^i y_{j+1} + \dots + w_{y,i-1}^i y_{i-1}. \tag{21}$$

Let us now denote by $K_{i,j}^+$ and $K_{i,j}^-$ the positive and negative indexes of the coefficients $w_{y,k}^i$ with $j + 1 \leq k \leq i - 1$, that is,

$$\begin{aligned} K_{i,j}^+ &= \{k \mid j + 1 \leq k \leq i - 1, w_{y,k}^i \geq 0\} \\ K_{i,j}^- &= \{k \mid j + 1 \leq k \leq i - 1, w_{y,k}^i < 0\}. \end{aligned} \tag{22}$$

Hence, (21) is maximum when y_k for $j + 1 \leq k \leq i - 1$ are

$$y_k = \begin{cases} 1 & \text{if } k \in K_{i,j}^+ \\ 0 & \text{if } k \in K_{i,j}^- \end{cases} \tag{23}$$

and (21) is minimum when y_k for $j + 1 \leq k \leq i - 1$ are

$$y_k = \begin{cases} 1 & \text{if } k \in K_{i,j}^- \\ 0 & \text{if } k \in K_{i,j}^+ \end{cases} \tag{24}$$

Hence,

- (i) Let $y_{j+1}^{i,1}, \dots, y_{i-1}^{i,1} \in \{0, 1\}$ be the values which maximize (21), that is, the values that maximize $f_i(\mathbf{x}, y_1, \dots, y_j, y_{j+1}, \dots, y_{i-1})$ and hence, the values which makes $\mathbf{P}(y_i = 1 \mid \mathbf{x}, y_1, \dots, y_j, y_{j+1}, \dots, y_{i-1})$ be maximum and,
- (ii) Let $y_{j+1}^{i,0}, \dots, y_{i-1}^{i,0} \in \{0, 1\}$ be the values that minimize (21), that is, the values that minimize $f_i(\mathbf{x}, y_1, \dots, y_j, y_{j+1}, \dots, y_{i-1})$ and hence, the values which makes $\mathbf{P}(y_i = 0 \mid \mathbf{x}, y_1, \dots, y_j, y_{j+1}, \dots, y_{i-1})$ be maximum.

Hence, $\max_{v \in \{0,1\}} \mathbf{P}(y_i = v \mid \mathbf{x}, y_1, \dots, y_j, y_{j+1}, \dots, y_{i-1})$ will be

$$\max_{v \in \{0,1\}} \{\mathbf{P}(y_i = v \mid \mathbf{x}, y_1, \dots, y_j, y_{j+1}^{i,v}, \dots, y_{i-1}^{i,v})\}. \tag{25}$$

Notice that $y_k^{i,1} = 1 - y_k^{i,0}$ for $j + 1 \leq k \leq i - 1$ and that according to the above definition of the sigmoid function, if $f_i(\mathbf{x}, y_1, \dots, y_j, y_{j+1}^{i,1}, \dots, y_{i-1}^{i,1}) \geq -f_i(\mathbf{x}, y_1, \dots, y_j, y_{j+1}^{i,0}, \dots, y_{i-1}^{i,0})$, then the maximum of $\mathbf{P}(y_i = v \mid \mathbf{x}, y_1, \dots, y_j, y_{j+1}^{i,v}, \dots, y_{i-1}^{i,v})$ is reached when $v = 1$ and otherwise when $v = 0$.

Therefore, the final expression for the heuristic will be

$$\hat{h} = \prod_{i=j+1}^m \max_{v \in \{0,1\}} \{\mathbf{P}(y_i = v \mid \mathbf{x}, y_1, \dots, y_j, y_{j+1}^{i,v}, \dots, y_{i-1}^{i,v})\}. \tag{26}$$

Remember that h^* and \hat{h} only differ on the values of y_{j+1}, \dots, y_{i-1} . In the former, the same values are common for all the factors of the product, whereas in the latter these values depend on each term i of the product, and hence, they can be different, which makes \hat{h} not be an optimal heuristic. On the other hand, the cost of computing \hat{h} is of polynomial order, unlike h^* which is of exponential order. Figure 5 exemplifies \hat{h} . Let us focus on the root node of the subtree in Fig. 5. The marginal joint conditional probability until this node is $g = 0.6$. For computing \hat{h} , first we evaluate the maximum marginal conditional probabilities

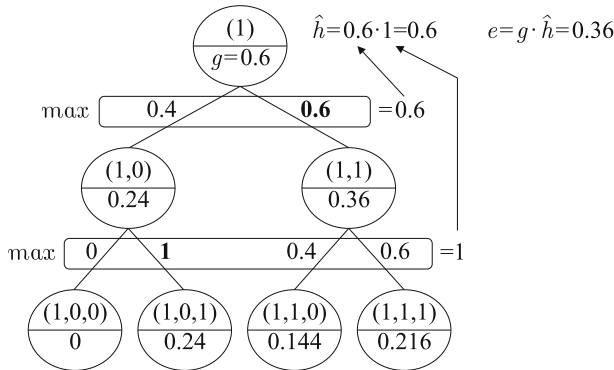


Fig. 5 An example of the computation of heuristic \hat{h}

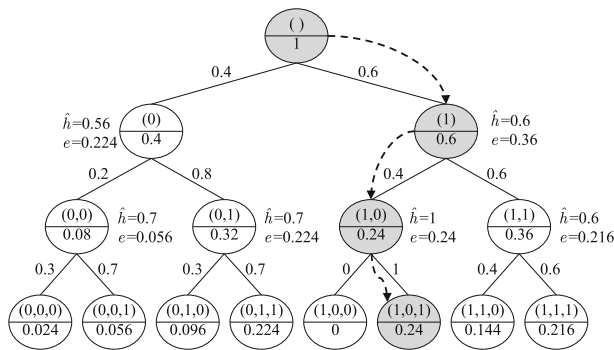


Fig. 6 An example of A* using the heuristic \hat{h} . The dotted arrows show the path followed by the algorithm. The values of g are provided inside each node

$P(y_2 | \mathbf{x}, y_1 = 1)$ provided by $f_2(\mathbf{x}, y_1)$ and $P(y_3 | \mathbf{x}, y_1 = 1, y_2)$ provided by $f_3(\mathbf{x}, y_1, y_2)$ which respectively are 0.6 and 1.0. Then we carry out their product by applying the product rule of the probability to estimate the marginal joint conditional probability from that node to a leaf. Notice that the maximum at each level does not correspond to the same branch of the tree. In other words, the maximum in the first level corresponds to $y_2 = 1$, whereas the maximum in the second level is obtained by $P(y_3 | \mathbf{x}, y_1 = 1, y_2 = 0)$ when $y_3 = 1$ fixing $y_2 = 0$. This is what makes the heuristic \hat{h} not to be the optimal h^* .

Figure 6 shows the path followed by A* using \hat{h} . It reaches the optimal leaf as is theoretically expected. Comparing this graph with the one in Fig. 3a that illustrates ϵ -A with $\epsilon = 0$, and taking into account the properties of the heuristics related to the dominance, ϵ -A, $\epsilon = 0$ (equivalent to A* using $h = 1$ or UC search) explores more nodes than A* using \hat{h} .

As a final remark, the A* algorithm using \hat{h} perfectly estimates the risk minimizer for the subset 0/1 loss, as both ϵ -A with $\epsilon = 0$ and ES do it. Even more, A* with \hat{h} expands equal or fewer nodes, since \hat{h} is more dominant than the heuristic $h = 1$ ($\hat{h} < h = 1$). Obviously, computing \hat{h} is more costly than computing $h = 1$ or applying just a UC search. The question is then if this additional computing time compensates the theoretical guarantee it has of expanding fewer nodes.

4.2 A general admissible heuristic for trading off the number of nodes explored and its computing time

The previous section pointed out that both \hat{h} and $h = 1$ are admissible heuristics. Besides, using the former theoretically guarantees that the A* algorithm explores fewer nodes than using the latter. But the cost of computing \hat{h} could be high in comparison to just considering a constant heuristic $h = 1$. This trade-off sheds light on including a parameter d for limiting the depth of the heuristic. Hence, for a node of level j and a value of d with $0 \leq d \leq m - j$, only the terms $\mathbf{P}(y_i | \mathbf{x}, y_1, \dots, y_j, y_{j+1}^i, \dots, y_{i-1}^i)$ of \hat{h} are evaluated using f_i for nodes from level $j + 1$ to level $j + k$ whereas these terms are estimated by the constant 1 for nodes from level $j + k + 1$ to level m , that is,

$$\hat{h}^d = \prod_{i=j+1}^{j+d} \mathbf{P}(y_i | \mathbf{x}, y_1, \dots, y_j, y_{j+1}^i, \dots, y_{i-1}^i) \cdot \prod_{i=j+d+1}^m 1, \tag{27}$$

or equivalently

$$\hat{h}^d = \prod_{i=j+1}^{j+d} \mathbf{P}(y_i | \mathbf{x}, y_1, \dots, y_j, y_{j+1}^i, \dots, y_{i-1}^i). \tag{28}$$

It is clear that \hat{h} is more dominant than \hat{h}^d ($\hat{h} < \hat{h}^d$) and it continues being admissible. In turn, \hat{h}^d is more dominant than $h = 1$ ($\hat{h}^d < h = 1$). Hence, it is expected that using the heuristic \hat{h}^d with $0 \leq d \leq m - j$ and $1 \leq j \leq m - 1$, the A* algorithm explores more number or equal number of nodes than \hat{h} , but less number or equal number of nodes than $h = 1$. In fact, \hat{h}^d encapsulates both heuristics, since taking the extreme values of d leads to them. On the one hand, taking the maximum value of d leads to \hat{h}^{m-j} with $1 \leq j \leq m - 1$ (we will denote this heuristic as \hat{h}^∞), which is in fact the heuristic \hat{h} detailed in Sect. 4.1. On the other hand, taking the minimum value of d leads to \hat{h}^0 which is indeed the heuristic $h = 1$. In general, \hat{h}^{d_1} is more dominant than \hat{h}^{d_2} if $d_1 > d_2$. However, the computational time of obtaining the values of \hat{h}^d increases as d increases. Hence, tuning d adequately one can obtain a balance between the number of nodes explored and the computational time employed to evaluate the heuristic.

The case of $d = 1$ has especial interest, since the heuristic is also valid for non-linear models f_j . The form of the heuristic is

$$\hat{h}^1 = \prod_{i=j+1}^{j+1} \mathbf{P}(y_i | \mathbf{x}, y_1, \dots, y_j, y_{j+1}^i, \dots, y_{i-1}^i) \cdot \prod_{i=j+2}^m 1, \tag{29}$$

that simplifying leads to

$$\hat{h}^1 = \mathbf{P}(y_{j+1} | \mathbf{x}, y_1, \dots, y_j). \tag{30}$$

As seen, this heuristic \hat{h}^1 means just evaluating the model of the node $f_j(\mathbf{x}, y_1, \dots, y_j)$, hence without making any restriction on the linearity of f_j .

Figure 7 shows an example of the path followed by the A* algorithm when \hat{h}^1 is taken as an heuristic. As seen, the A* algorithm explores more nodes using this heuristic than using \hat{h} , but the computational time of evaluating the heuristic is considerably reduced, as we will show later on in the experiments.

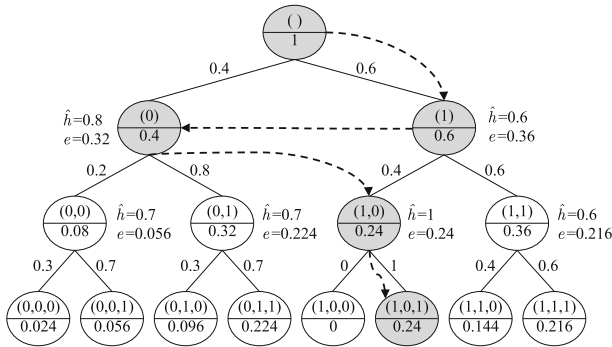


Fig. 7 An example of A* using the heuristic \hat{h}^1 . The dotted arrows show the path followed by the algorithm. The values of g are provided inside each node

4.3 Implementation details

Here we significantly extend the results of A* using \hat{h} presented in the preliminary work (Mena et al. 2015). However, they present differences because we have carried out an improvement in the implementation of the algorithms in order to become more efficient. In fact, the time results reported there in comparison with the number of nodes explored was the key to make us reconsider the implementation of the algorithms. Algorithm 1 shows the pseudocode of A* using \hat{h}^d .

Particularly, the differences between this version and that of the preliminary work (Mena et al. 2015) are that now:

1. The models f_i (parameters w_x^i , w_y^i and β^i) are ordered according to the label order of the chain.
2. All repeated computations are calculated just once and stored:
 - (a) The evaluation of $\langle w_x^i, x \rangle$ for each linear model f_i , which is common for all nodes corresponding to f_i , is computed once and saved in variable WX (line 2).
 - (b) Sets K^+ and K^- , that is the best values of y_k for the unknown part of the heuristic, are computed only once before starting the main loop of A* (lines 3–8).
3. Open set, Q , stores tuples of four elements: label combination, level, e and g . Q is stored in a resizable vector whose positions are reused. The left child is stored in the position of its parent (line 21). Notice that parent information is not required to obtain the final solution, since Q contains all the required information for each node.
4. The design of Q allows us to optimize the function Max , the operation to obtain the best node to be expanded. The function Max only needs to evaluate the first elements ($1 \dots last$) of Q (line 13), because the rest of the vector is unused.
5. Algorithm 1 contains several auxiliary functions to make the pseudocode shorter and more easily understood. However, in the actual program all these functions were coded inline, making the code faster because all the calls to functions with large parameters, like the function *Heuristic*, are avoided.

4.4 Complexity analysis

Despite the theoretical optimality properties of the A* algorithm, it might not be useful in some problems because it may explore an exponential number of nodes with regard to

Algorithm 1 Pseudocode of the implementation of A* algorithm using \hat{h}^d

```

1: function A*
   Input:  $x, [W, \beta]$  a CC Linear Model,  $m, d, BlockSize$ 
   Output: Label combination with highest probability for  $x$ 
2:  $WX \leftarrow W_x * x$  // Computes  $\langle w_x, x \rangle$  for all labels
3:  $K^+ \leftarrow AllocMemory(m)$ 
4:  $K^- \leftarrow AllocMemory(m)$ 
5: for  $label = [2 : m]$  do // Starts at 2nd level, no label attributes in the 1st model
6:    $K^+[label] \leftarrow \llbracket W_y[label] \geq 0 \rrbracket$ 
7:    $K^-[label] \leftarrow \llbracket W_y[label] < 0 \rrbracket$ 
8: end for
9:  $Q \leftarrow AllocMemory(BlockSize)$  // tuples  $\{Labels, Level, e, g\}$ 
10:  $Q[1] \leftarrow \{\{ \}, 0, 1, 1\}$  // root node, empty label set, level 0,  $e$  and  $g = 1$ 
11:  $last \leftarrow 1$  // last element used in  $Q$ 
12: while true do
13:    $[Best, Position] \leftarrow Max(Q, last)$ 
14:   if  $Best.Level = m$  then
15:     return  $Best.Labels$  // Leaf node
16:   end if
17:    $level \leftarrow Best.Level + 1$ 
18:    $P \leftarrow 1/(1 + exp(-(WX[level] + \beta[level] + \langle w_y[level], Best.Labels \rangle)))$ 
19:   // Left child
20:    $\hat{h}^d \leftarrow Heuristic(d, level, [Best.Labels 0], K^+, K^-, WX, W_y, x)$  // Eq(28)
21:    $Q[Position] \leftarrow \{\{Best.Labels 0\}, level, Best.g * (1 - P) * \hat{h}^d, Best.g * (1 - P)\}$ 
22:   // Right child
23:    $last \leftarrow last + 1$ 
24:   if  $last > Q.size$  then
25:      $Q \leftarrow resize(Q, BlockSize)$ 
26:   end if
27:    $\hat{h}^d \leftarrow Heuristic(d, level, [Best.Labels 1], K^+, K^-, WX, W_y, x)$  // Eq(28)
28:    $Q[last] \leftarrow \{\{Best.Labels 1\}, level, Best.g * P * \hat{h}^d, Best.g * P\}$ 
29: end while
30: end function

```

the depth of the tree [see Pearl (1984) and Russell and Norvig (2003)]. Only under certain conditions over the heuristic taken, is it possible to obtain an algorithm with a theoretical complexity less than exponential. This is the case if one is able to prove that the heuristic h satisfies the following condition for any level j with regard to the optimal heuristic h^* (Russell and Norvig 2003)

$$|h^*(j) - h(j)| \leq \mathcal{O}(\log h^*(j)). \tag{31}$$

In general, it occurs that that error is at least linear with regard to the path cost, hence leading to exponential growth. However, by taking care of designing good heuristics, one can obtain huge profit in relation to not providing any kind of heuristic information. That is the case with our heuristic \hat{h}^d with regard to the heuristic $h = 1$, which does not provide any kind of heuristic information, in spite of being an improvement of the ES.

On what follows, let us focus on the heuristic \hat{h}^∞ , since it is the most dominant heuristic among the heuristics proposed. According to the results where the noise is increasing (see Figs. 9, 10 of Sect. 5.3), it is not clear that \hat{h}^∞ is exponential at all as the rest of the algorithms clearly show, at least for the datasets taken. This behavior of \hat{h}^∞ sheds light on performing a deeper analysis over other situations. In this sense, let us consider a theoretical case when the difference between \hat{h}^∞ and h^* is maximum. For instance, this situation occurs (see Fig. 8) when the first level has probability 0 in the left branch and 1 in the right branch, all the

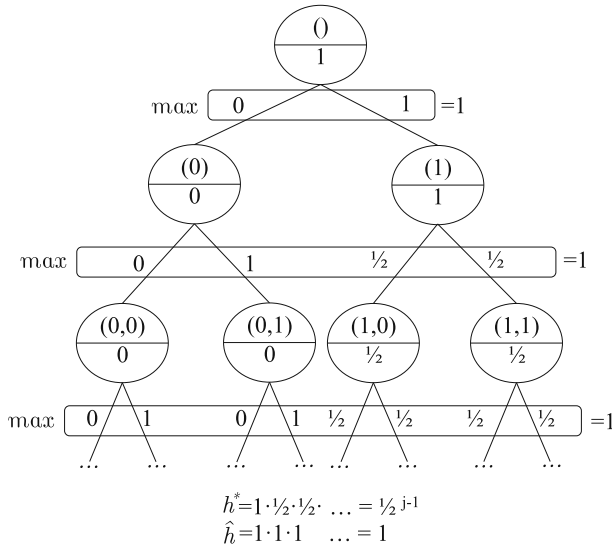


Fig. 8 An example of worst cases for heuristic \hat{h}

probabilities of the left subtree are 0 (for the left branches) and 1 (for the right branches) and all the probabilities of the right subtree are $1/2$ (for all the branches). In this case, the error can be bounded as follows for any level j

$$|h^*(j) - \hat{h}^\infty(j)| \leq |1/2^{j-1} - 1|. \tag{32}$$

This bound tends to 1 as the level j grows, hence, in this case, it is not possible to guarantee that the theoretical complexity is less than exponential. However, these extreme cases are hardly likely to occur. Let us remember that the probabilities are evaluated from the models, where the description x of the examples plays an important role. This description is equal for obtaining the probabilities of all the nodes of the same level, hence, such probabilities will probably not differ so much among them as in the extreme case considered.

5 Experiments

This section deals with the experiments carried out with all the approaches. Before discussing the experimental results in Sects. 5.2 and 5.3, let us describe in Sect. 5.1 the common settings of all experiments: datasets, learning algorithms and parameter selection procedures. Finally, Sect. 5.3 reports the computational results of the methods for more complex problems, such as those that include noise.

5.1 Settings

The experiments were performed over several benchmark multi-label datasets whose main properties are shown in Table 1. As can be seen, there are significant differences in the number of attributes, instances and labels. The cardinality—number of labels per instance—varies from 1.07 to 4.27. Concerning the number of labels, there are some datasets with just 5, 6 or 7 labels, whereas others have more than 100, one of them even has almost 400 labels.

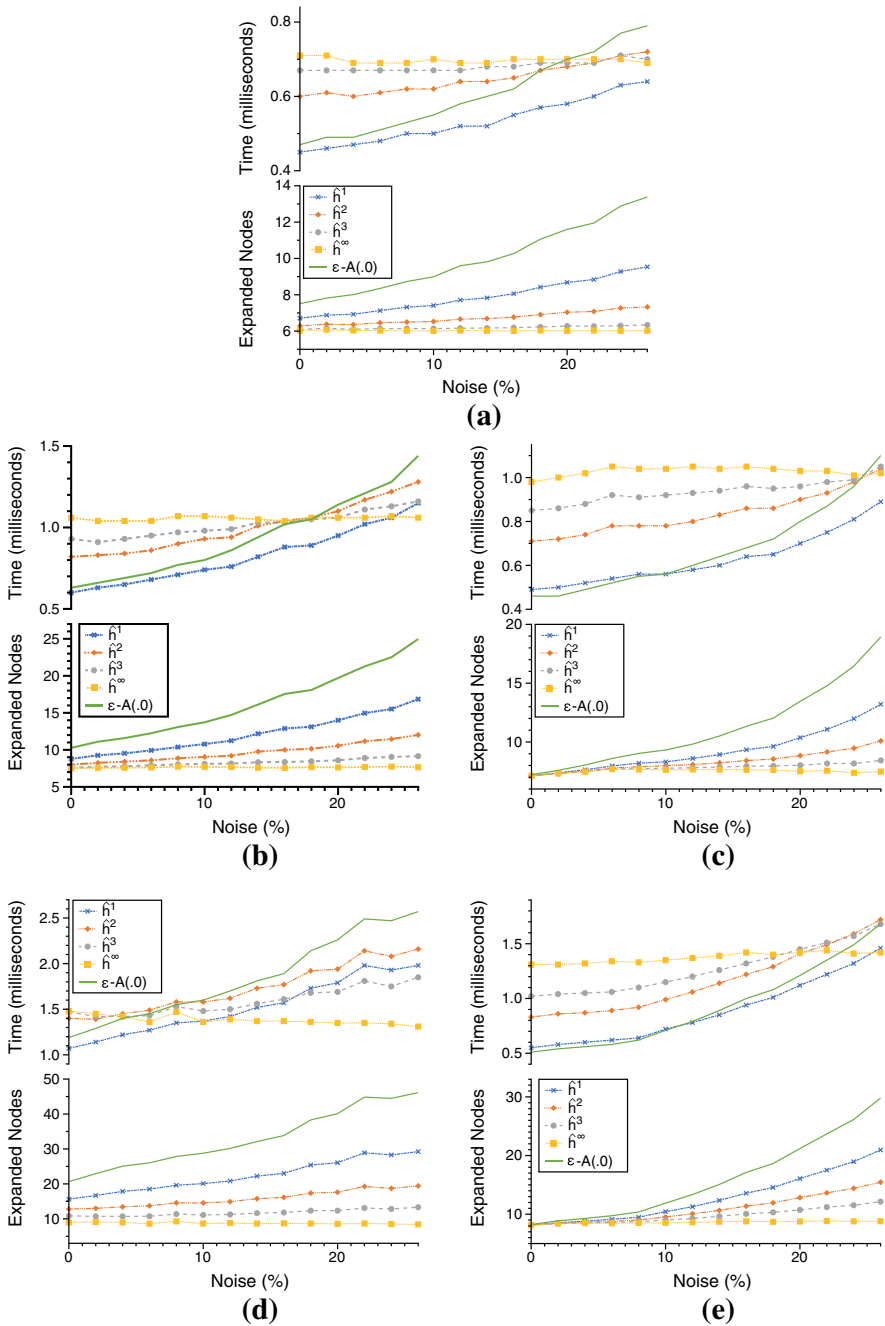


Fig. 9 Number of nodes expanded and computational time employed (ms) for the algorithm A* with the heuristic \hat{h} for different values of the parameter d (1, 2, 3, ∞) and for the ϵ -A with $\epsilon = 0$, for different percentages of noise and for datasets with few labels (from 5 to 7). **a** Image (5), **b** emotions (6), **c** scene (6), **d** flags (7) and **e** reuters (7)

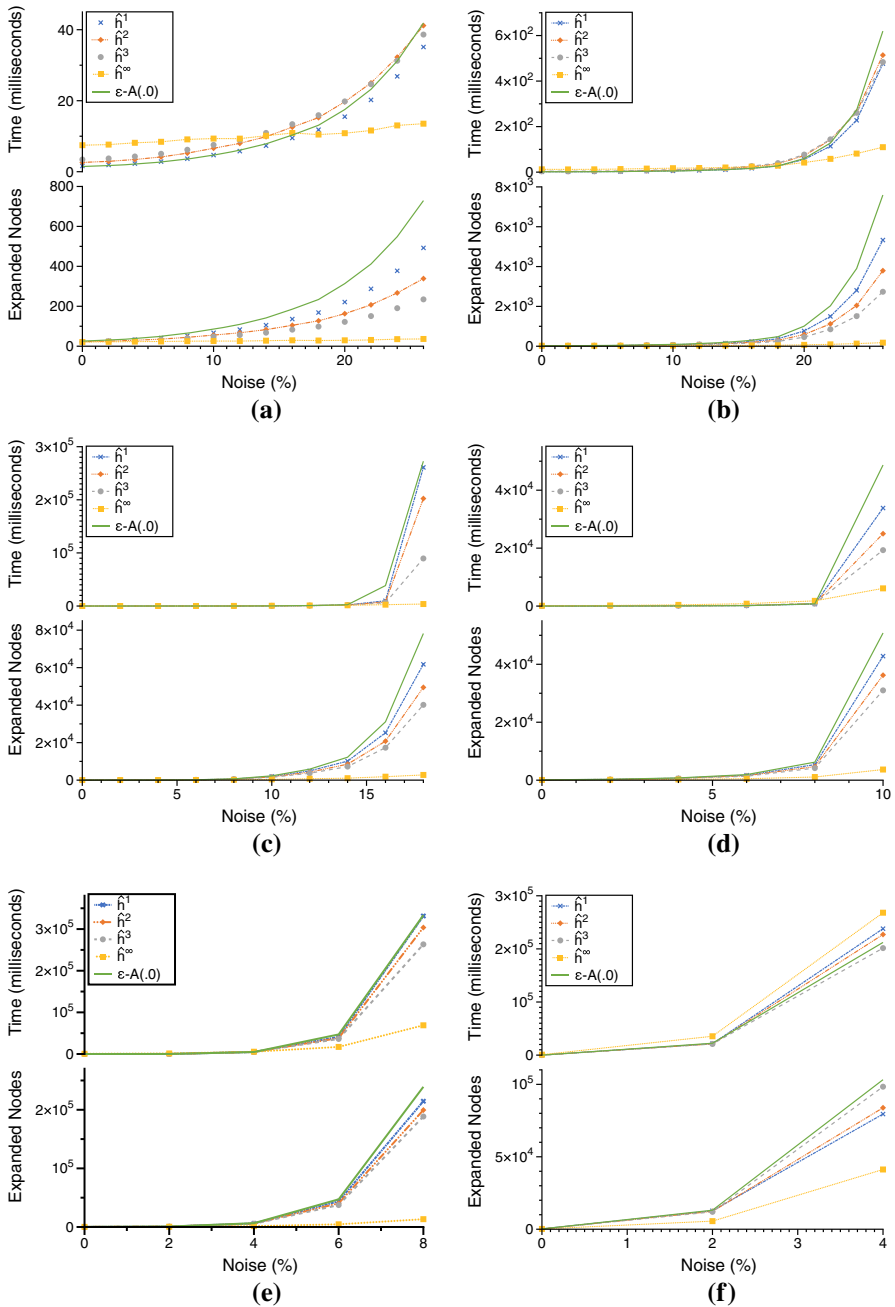


Fig. 10 Number of nodes expanded and computational time employed (ms) for the algorithm A* with the heuristic \hat{h} for different values of the parameter d (1, 2, 3, ∞) and for the ϵ -A with $\epsilon = 0$, for different percentages of noise and for datasets with more labels (from 14 to 159). **a** Yeast (14), **b** slashdot (22), **c** medical (45), **d** enron (53), **e** mediamill (101) and **f** bibtex (159)

Table 1 Properties of the datasets

| Datasets | Instances | Attributes | Labels | Cardinality |
|------------|-----------|------------|--------|-------------|
| Bibtex | 7395 | 1836 | 159 | 2.40 |
| Corel5k | 5000 | 499 | 374 | 3.52 |
| Emotions | 593 | 72 | 6 | 1.87 |
| Enron | 1702 | 1001 | 53 | 3.38 |
| Flags | 194 | 19 | 7 | 3.39 |
| Image | 2000 | 135 | 5 | 1.24 |
| Mediamill* | 5000 | 120 | 101 | 4.27 |
| Medical | 978 | 1449 | 45 | 1.25 |
| Reuters | 7119 | 243 | 7 | 1.24 |
| Scene | 2407 | 294 | 6 | 1.07 |
| Slashdot | 3782 | 1079 | 22 | 1.18 |
| Yeast | 2417 | 103 | 14 | 4.24 |

The approaches for inference in PCC compared with our proposal were those discussed throughout the paper, except for the ES. No experiment was carried out with the ES method due to its computational cost. Hence, the methods compared with the A* algorithm with the heuristic \hat{h} for different values of the parameter d (1, 2, 3, ∞) were GS, ϵ -A algorithm for different values of ϵ (.0, .25, .5) and BS for different values of beam width b (1, 2, 3, 10). Let us remember that the ϵ -A algorithm with $\epsilon = 0.5$ is equivalent to GS and to BS with $b = 1$.

The results will be presented in terms of the example-based subset 0/1 loss estimated by means of a 10-fold cross-validation.

The base learner employed to obtain the binary classifiers that compose all these multi-label models was *logistic regression* (Lin et al. 2008) with probabilistic output. The regularization parameter C was established for each *individual* binary classifier performing a grid search over the values $C \in \{10^{-3}, 10^{-2}, \dots, 10^3\}$ optimizing the brier loss estimated by means of a balanced 2-fold cross validation repeated 5 times. The brier loss (Brier 1950) is a proper score that measures the accuracy of probabilistic predictions, as *logistic regression* does. The expression is as follows

$$\frac{1}{n} \sum_{i=1}^n (\hat{p}_i - a_i)^2, \quad (33)$$

where for an instance i , p_i is the predicted probability of a certain label and a_i is the actual value of the label (0 or 1).

5.2 Results over benchmark datasets

Tables 2, 3 and 4 respectively show the subset 0/1 loss, the number of nodes explored and the computational time (in seconds) averaged per test instances for the different methods compared.

Before discussing the results of the tables, let us remember that only the A* algorithm using \hat{h}^d (and consequently also the ϵ -A with $\epsilon = 0$, since it is the A* algorithm with $h = 1$ or \hat{h}^0) provides a Bayes optimal inference, as the ES does. This means that they always predict

Table 2 Subset 0/1 loss for the different methods

| Datasets | \hat{h}^* (ϵ -A(.0)) | ϵ -A(.25) | GS/BS(1) ϵ -A(.5) | BS(2) | BS(3) | BS(10) |
|------------|----------------------------------|--------------------|----------------------------|--------------|--------------|--------------|
| Bibtex | 81.92 | 81.95 | 82.19 | 81.88 | 81.92 | 81.92 |
| Corel5k | 97.48 | 98.62 | 98.90 | 98.30 | 98.04 | 97.48 |
| Emotions | 71.16 | 71.82 | 72.83 | 72.16 | 71.32 | 71.16 |
| Enron | 83.14 | 84.26 | 85.43 | 83.43 | 83.37 | 83.14 |
| Flags | 87.13 | 87.16 | 86.13 | 88.21 | 87.13 | 87.13 |
| Image | 68.35 | 68.35 | 69.75 | 68.35 | 68.35 | 68.35 |
| Mediamill* | 83.86 | 84.58 | 85.80 | 84.10 | 83.86 | 83.86 |
| Medical | 30.37 | 30.37 | 30.67 | 30.37 | 30.37 | 30.37 |
| Reuters | 22.73 | 22.70 | 23.60 | 22.69 | 22.73 | 22.73 |
| Scene | 31.86 | 31.86 | 33.28 | 31.90 | 31.86 | 31.86 |
| Slashdot | 51.80 | 52.22 | 54.49 | 51.77 | 51.80 | 51.80 |
| Yeast | 76.95 | 77.62 | 79.77 | 76.83 | 77.08 | 76.95 |

Those scores that are equal to or better than optimal predictions reached by ϵ -A and ES are shown in bold

the label combination with the highest joint conditional probability. Despite the fact that other methods may predict other label combinations with lower joint conditional probability for some examples, in a few cases they obtain better subset 0/1 scores. This fact is due to several reasons, mainly (a) the relatively small size of testing sets, and (b) that the models f_j obtained to estimate the joint conditional probability $P(\mathbf{y} | \mathbf{x})$ do not usually return true estimations. Theoretically, under perfect conditions (large test sets and perfect models), the A* algorithm using \hat{h}^d would obtain the best scores. In general, the performance of the ϵ -A algorithm decreases as the value of ϵ increases and the performance of the BS method increases as b increases. The BS method reaches stability for low values of the beam b , it even converges to the performance of the A* algorithm but at the cost of exploring many more nodes.

With regard to the number of nodes explored (see Table 3), GS (equivalent to ϵ -A algorithm with $\epsilon = 0.5$ and to BS(1)) is the method which explores the least number of nodes, since it only goes over one path in the tree. In fact, such a number corresponds to the number of labels (plus one if the root is considered as an explored node). It follows the A* algorithm using our heuristic \hat{h}^d for any value of d , although it is exceeded by the ϵ -A algorithm with $\epsilon > 0$ and by BS with $b > 1$ in some cases. However, let us remember that neither GS nor the ϵ -A algorithm with $\epsilon > 0$ nor BS with any value of the beam b guarantee reaching an optimal solution as the A* algorithm using heuristic \hat{h}^d does. So, it is not surprising that they explore fewer nodes. What it is actually surprising is that for most of the cases they explore more nodes even without reaching an optimal solution.

Let us now focus on the methods that theoretically reach the optimum (the A* algorithm with the heuristic \hat{h}^d or the ϵ -A algorithm with $\epsilon = 0$). As it has theoretically been shown, the A* algorithm with the heuristic \hat{h}^∞ explores the least number of nodes, followed by the same algorithm but with the heuristic \hat{h}^3 , then with the heuristic \hat{h}^2 , after that with the heuristic \hat{h}^1 and finally the ϵ -A algorithm with $\epsilon = 0$.

The computational time is expected to be higher as more nodes are explored, but looking at Table 4 one can see that this does not happen at all. This is true for GS, the ϵ -A algorithm and the BS technique, but the A* algorithm obtains higher computational time in spite of exploring considerably fewer nodes. The reason for that is the time spent in computing the

Table 3 Number of explored nodes for the different methods

| Datasets | \hat{h}^1 | \hat{h}^2 | \hat{h}^3 | \hat{h}^∞ | $\epsilon - A(.0)$ | $\epsilon - A(.25)$ | GS/BS(1) $\epsilon - A(.5)$ | BS(2) | BS(3) | BS(10) |
|------------|-------------|-------------|-------------|------------------|--------------------|---------------------|--------------------------------|--------|---------|---------|
| Bibtex | 283.31 | 277.97 | 273.23 | 215.91 | 289.27 | 184.00 | 160.00 | 319.00 | 477.00 | 1575.00 |
| Core15k | 1456.92 | 1443.08 | 1431.57 | 1338.62 | 1474.17 | 517.11 | 375.00 | 749.00 | 1122.00 | 3725.00 |
| Emotions | 9.12 | 8.28 | 7.80 | 7.68 | 10.67 | 10.78 | 7.00 | 13.00 | 18.00 | 45.00 |
| Enron | 110.32 | 106.50 | 103.28 | 75.51 | 114.81 | 77.29 | 54.00 | 107.00 | 159.00 | 515.00 |
| Flags | 16.45 | 12.69 | 10.43 | 8.79 | 22.56 | 16.33 | 8.00 | 15.00 | 21.00 | 55.00 |
| Image | 6.64 | 6.28 | 6.14 | 6.11 | 7.33 | 7.66 | 6.00 | 11.00 | 15.00 | 35.00 |
| Mediamill* | 188.18 | 185.64 | 183.68 | 178.89 | 191.76 | 142.37 | 102.00 | 203.00 | 303.00 | 995.00 |
| Medical | 46.61 | 46.58 | 46.55 | 46.40 | 46.64 | 46.65 | 46.00 | 91.00 | 135.00 | 435.00 |
| Reuters | 8.15 | 8.14 | 8.13 | 8.13 | 8.24 | 8.26 | 8.00 | 15.00 | 21.00 | 55.00 |
| Scene | 7.16 | 7.15 | 7.15 | 7.15 | 7.25 | 7.25 | 7.00 | 13.00 | 18.00 | 45.00 |
| Slashdot | 24.98 | 24.88 | 24.86 | 24.84 | 25.29 | 24.89 | 23.00 | 45.00 | 66.00 | 205.00 |
| Yeast | 23.60 | 22.80 | 22.38 | 21.01 | 26.09 | 26.02 | 15.00 | 29.00 | 42.00 | 125.00 |

Table 4 Average prediction time (in seconds) per example for all methods

| Datasets | \hat{h}^1 | \hat{h}^2 | \hat{h}^3 | \hat{h}^∞ | $\epsilon - A(.0)$ | $\epsilon - A(.25)$ | GS/BS(1) $\epsilon - A(.5)$ | BS(2) | BS(3) | BS(10) |
|------------|-------------|-------------|-------------|------------------|--------------------|---------------------|--------------------------------|--------|--------|--------|
| Bibtex | 0.0195 | 0.0314 | 0.0411 | 0.8407 | 0.0162 | 0.0099 | 0.0079 | 0.0110 | 0.0156 | 0.0507 |
| Core15k | 0.1684 | 0.2578 | 0.3333 | 23.8344 | 0.1360 | 0.0161 | 0.0110 | 0.0278 | 0.0404 | 0.1361 |
| Emotions | 0.0006 | 0.0008 | 0.0009 | 0.0011 | 0.0006 | 0.0006 | 0.0004 | 0.0005 | 0.0006 | 0.0013 |
| Enron | 0.0085 | 0.0141 | 0.0185 | 0.1084 | 0.0072 | 0.0030 | 0.0019 | 0.0039 | 0.0055 | 0.0172 |
| Flags | 0.0011 | 0.0014 | 0.0015 | 0.0015 | 0.0012 | 0.0007 | 0.0004 | 0.0005 | 0.0007 | 0.0016 |
| Image | 0.0005 | 0.0006 | 0.0007 | 0.0007 | 0.0005 | 0.0005 | 0.0004 | 0.0004 | 0.0005 | 0.0009 |
| Mediamill* | 0.0144 | 0.0242 | 0.0326 | 0.4847 | 0.0115 | 0.0055 | 0.0037 | 0.0072 | 0.0105 | 0.0333 |
| Medical | 0.0033 | 0.0057 | 0.0077 | 0.0495 | 0.0027 | 0.0027 | 0.0027 | 0.0033 | 0.0046 | 0.0144 |
| Reuters | 0.0006 | 0.0008 | 0.0010 | 0.0013 | 0.0005 | 0.0005 | 0.0005 | 0.0005 | 0.0007 | 0.0016 |
| Scene | 0.0005 | 0.0007 | 0.0009 | 0.0010 | 0.0005 | 0.0005 | 0.0004 | 0.0005 | 0.0006 | 0.0013 |
| Slashdot | 0.0018 | 0.0029 | 0.0039 | 0.0124 | 0.0015 | 0.0014 | 0.0012 | 0.0016 | 0.0022 | 0.0063 |
| Yeast | 0.0016 | 0.0026 | 0.0034 | 0.0075 | 0.0015 | 0.0010 | 0.0006 | 0.0010 | 0.0014 | 0.0040 |

heuristic. In this respect, the A^* algorithm is faster as the parameter d diminishes, although this implies exploring more nodes. Hence, considering the methods that reach the optimum, the ϵ -A approximation algorithm with $\epsilon = 0$ is the fastest method, followed by A^* using the heuristic \hat{h}^1 , then \hat{h}^2 and so on until ending with the heuristic \hat{h}^∞ .

As a conclusion, the A^* algorithm with the heuristic h^1 can be considered a good alternative for guaranteeing optimal performance in terms of subset 0/1 and balancing the trade-off between the number of nodes explored and the computational time. Besides, it is also applicable for non-linear models.

5.3 Results with noise

We have also performed experiments including noise in the datasets in order to analyze more in depth the power of the A^* algorithm for inference in PCCs. This study arises from the fact that the number of nodes explored by A^* is quite low in comparison with the number of labels (depth of the tree). This is so even using the heuristic $h = 1$ (ϵ -A with $\epsilon = 0$) which does not provide information as other heuristics do, for instance, \hat{h} . This means that the algorithm A^* performs few backtracking in the tree, hence, the probabilities provided by the models may be not so close from 0.5 or, even more, they may be close to 0 or 1. In this way, by adding noise we expect the probabilities to be closer to 0.5, making the problem more complex where more informative heuristics can show their strength.

Certain grades of noise in terms of percentage were added to the datasets. No kinds of noise were included in the description of the examples, that is, in the x part, otherwise, the noise was only introduced in the labels of the examples, that is, in the y part. In this sense, a percentage of the values of the labels was swapped from being relevant to become irrelevant and vice versa for the whole dataset. This means that for a given example the relevance of either all their labels or only some of them or even none of them was changed. The percentage of noise included ranges from 0 to 25% for datasets with fewer than 22 labels. However, the experiments did not finish in a prudential time by using all the percentages of this range for datasets with more than 45 labels. In particular, the maximum percentage considered in this respect for medical (45 labels), enron (53 labels), mediamill (101 labels) and bibtex (159 labels) respectively was 18, 10, 8 and 4%.

Figures 9 and 10, respectively show for datasets with few labels (from 5 to 7) and for datasets with more labels (from 14 to 159), the number of nodes expanded and the computational time employed by the A^* algorithm with the heuristic \hat{h} for different values of the parameter d (1, 2, 3, ∞) and by the ϵ -A with $\epsilon = 0$, all of them for different percentages of noise.

Obviously, both the number of nodes explored and the computational time increases as the percentage of noise increases. Regarding the number of nodes explored, it decreases as d increases as theoretically expected. In this respect, the number of nodes expanded using the heuristic \hat{h}^∞ keeps quite steady as the noise increases in comparison to using the rest of the heuristics. In fact, by using these other heuristics, the number of nodes rockets from a certain percentage of noise. This is so in general, but it especially happens for datasets with more than 14 labels.

Concerning the computational time and for datasets with few labels, A^* using \hat{h}^d from $d > 1$ is steadier than using \hat{h}^1 or ϵ -A with $\epsilon = 0$ as the noise increases. In fact, these two last approaches are the best options, \hat{h}^1 being clearly better than ϵ -A with $\epsilon = 0$. However, for high percentages of noise the other options begin to show their strength. Only flags among the datasets with few labels presents a different behavior. In this case, the behavior of the computational time is similar to that of the datasets with more labels, which coincides with

the behavior of the number of nodes. So, the computational time increases as d decreases and \hat{h}^∞ becomes clearly the best option. In any case, we are measuring the computational time in milliseconds, so the worth of the heuristics becomes crucial for datasets with more than 22 labels.

6 Conclusions

This paper proposes a family of admissible heuristics for the A* algorithm for performing inference in PCCs. In this way, providing admissible heuristics leads obtaining optimal predictions in terms of subset 0/1 loss and exploring fewer nodes than previous approaches that also produce optimal predictions. The family of heuristics (\hat{h}^d) is defined by a parameter d which determines the depth in the tree for evaluating the heuristic and, hence, controls the trade-off between the number of nodes explored and the computational time of the algorithm due to the evaluation of the heuristic. In this manner, the algorithm A*(\hat{h}^d) explores fewer nodes but it expends more time in reaching a solution as d increases. So far, only uniform-cost search (or ϵ -A($\epsilon = .0$), which in turn is the particular case of the A*(\hat{h}^0)) has been shown to provide optimal subset 0/1 loss, unlike other approaches such as GS, ϵ -A($\epsilon > 0$) or BS that only estimate it. The algorithm A*(\hat{h}^d) with $d > 1$ is limited to linear models, but this is not a major drawback because it is quite common to employ linear models. This is usually the case for MLC problems with many labels in which the examples of some classes may be scarce, since using non-linear models in such problems can lead to overfitting. Only the particular case of $d = 1$ is also suitable for non-linear models, since it just involves evaluating the models from all the known values of the labels.

In the experiments performed over the benchmark datasets, the algorithm A*(\hat{h}^1) or uniform-cost search (or ϵ -A($\epsilon = .0$) or A*(\hat{h}^0)) are better choices than A*(\hat{h}^d) for $d > 1$, since the computational cost of evaluating the heuristic does not seem to compensate for the reduction in nodes explored. In this respect, adding noise to the benchmark datasets, and, hence, obtaining more complex problems, allows us to show the strength of A*(\hat{h}^d) for $d > 1$. In this sense, the behavior of A*(\hat{h}^∞) is quite steady as the percentage of noise increases, unlike A*(\hat{h}^d) for the rest of the values of d , including $d = 0$ (ϵ -A($\epsilon = .0$)). Furthermore, only for small datasets in a number of labels, \hat{h}^1 becomes the best alternative.

In our opinion, the heuristic search is a promising line of research for inference in PCCs. As future work, new admissible heuristics can be devised, keeping in mind that we should look for a trade-off between their computational complexity and the quality of the estimations.

Acknowledgements This research has been supported by the Spanish Ministerio de Economía y Competitividad (Grants TIN2011-23558, TIN2015-65069).

References

- Brier, G. W. (1950). Verification of forecasts expressed in terms of probability. *Monthly Weather Review*, 78(1), 1–3.
- Cheng, W., & Hüllermeier, E. (2009). Combining instance-based learning and logistic regression for multilabel classification. *Machine Learning*, 76(2–3), 211–225.
- Dembczyński, K., Cheng, W., & Hüllermeier, E. (2010). Bayes optimal multilabel classification via probabilistic classifier chains. *ICML, 2010*, 279–286.
- Dembczyński, K., Waegeman, W., Cheng, W., & Hüllermeier, E. (2012). On label dependence and loss minimization in multi-label classification. *Machine Learning*, 88(1–2), 5–45.

- Dembczynski, K., Waegeman, W., & Hüllermeier, E. (2012). An analysis of chaining in multi-label classification. In L. D. Raedt, C. Bessire, D. Dubois, P. Doherty, P. Frasconi, F. Heintz, & P. J. F. Lucas (Eds.), *ECAI, frontiers in artificial intelligence and applications* (Vol. 242, pp. 294–299). Amsterdam: IOS Press.
- Elisseeff, A., & Weston, J. (2005). A kernel method for multi-labelled classification. In: *ACM conference on research and development in information retrieval* (pp. 274–281).
- Hart, P. E., Nilsson, N. J., & Raphael, B. (1968). A formal basis for the heuristic determination of minimum cost paths. *IEEE Transactions on Systems Science and Cybernetics*, 4(2), 100–107.
- Kumar, A., Vembu, S., Menon, A. K., & Elkan, C. (2012). Learning and inference in probabilistic classifier chains with beam search. *ECML/PKDD, 2012*, 665–680.
- Kumar, A., Vembu, S., Menon, A. K., & Elkan, C. (2013). Beam search algorithms for multilabel learning. *Machine Learning*, 92(1), 65–89.
- Lin, C. J., Weng, R. C., & Keerthi, S. S. (2008). Trust region Newton method for logistic regression. *Journal of Machine Learning Research*, 9(Apr), 627–650.
- McCallum, A. K. (1999). Multi-label text classification with a mixture model trained by EM. In: *AAAI 99 workshop on text learning*.
- Mena, D., Montañés, E., Quevedo, J. R., & del Coz, J. J. (2015). *An overview of inference methods in probabilistic classifier chains for multi-label classification*. Technical Report: Artificial Intelligence Center, University of Oviedo, Spain, Campus de Viesques, Gijón.
- Mena, D., Montañés, E., Quevedo, J. R., & del Coz, J. J. (2015). Using a* for inference in probabilistic classifier chains. In: *IJCAI 2015, Proceedings of the 20th international joint conference on artificial intelligence* (pp. 3707–3713), Buenos Aires, Argentina, July 25–31.
- Montañés, E., Quevedo, J. R., & del Coz, J. J. (2011). Aggregating independent and dependent models to learn multi-label classifiers. In: *ECML/PKDD'11—Volume Part II* (pp. 484–500). Springer-Verlag.
- Montañés, E., Senge, R., Barranquero, J., Quevedo, J. R., del Coz, J. J., & Hüllermeier, E. (2014). Dependent binary relevance models for multi-label classification. *Pattern Recognition*, 47(3), 1494–1508. Handwriting Recognition and other PR Applications
- Pearl, J. (1984). *Heuristics: Intelligent search strategies for computer problem solving*. Boston, MA: Addison-Wesley Longman Publishing Co. Inc.
- Read, J., Martino, L., & Luengo, D. (2014). Efficient monte carlo methods for multi-dimensional learning with classifier chains. *Pattern Recognition*, 47(3), 1535–1546.
- Read, J., Pfahringer, B., Holmes, G., & Frank, E. (2009). Classifier chains for multi-label classification. In: *ECML/PKDD'09, LNCS* (pp. 254–269). Springer
- Read, J., Pfahringer, B., Holmes, G., & Frank, E. (2011). Classifier chains for multi-label classification. *Machine Learning*, 85(3), 333–359.
- Russell, S. J., & Norvig, P. (2003). *Artificial intelligence: A modern approach* (2nd ed.). New York: Pearson Education.
- Senge, R., Barranquero, J., del Coz, J. J., & Hüllermeier, E. (2012a). *Rectifying classifier chains for multi-label classification*. Technical Report.
- Senge, R., del Coz, J. J., & Hüllermeier, E. (2012b). On the problem of error propagation in classifier chains for multi-label classification. In: *Conference of the German classification society on data analysis, machine learning and knowledge discovery*.
- Tsoumakas, G., & Vlahavas, I. (2007). Random k-labelsets: An ensemble method for multilabel classification. In: *ECML/PKDD'07, LNCS* (pp. 406–417). Springer.
- Zhang, M. L., & Zhou, Z. H. (2007). MI-knn: A lazy learning approach to multi-label learning. *Pattern Recognition*, 40(7), 2038–2048.