

Inexact Field Learning Approach for Data Mining

Honghua Dai

Deakin University, Australia

INTRODUCTION

Inexact fielding learning (IFL) (Ciesieski & Dai, 1994; Dai & Ciesieski, 1994a, 1994b, 1995, 2004; Dai & Li, 2001) is a rough-set, theory-based (Pawlak, 1982) machine learning approach that derives inexact rules from fields of each attribute. In contrast to a point-learning algorithm (Quinlan, 1986, 1993), which derives rules by examining individual values of each attribute, a field learning approach (Dai, 1996) derives rules by examining the fields of each attribute. In contrast to exact rule, an inexact rule is a rule with uncertainty. The advantage of the IFL method is the capability to discover high-quality rules from low-quality data, its property of low-quality data tolerant (Dai & Ciesieski, 1994a, 2004), high efficiency in discovery, and high accuracy of the discovered rules.

BACKGROUND

Achieving high prediction accuracy rates is crucial for all learning algorithms, particularly in real applications. In the area of machine learning, a well-recognized problem is that the derived rules can fit the training data very well, but they fail to achieve a high accuracy rate on new unseen cases. This is particularly true when the learning is performed on low-quality databases. Such a problem is referred as the Low Prediction Accuracy (LPA) problem (Dai & Ciesieski, 1994b, 2004; Dai & Li, 2001), which could be caused by several factors. In particular, overfitting low-quality data and being misled by them seem to be the significant problems that can hamper a learning algorithm from achieving high accuracy. Traditional learning methods derive rules by examining individual values of instances (Quinlan, 1986, 1993). To generate classification rules, these methods always try to find cut-off points, such as in well-known decision tree algorithms (Quinlan, 1986, 1993).

What we present here is an approach to derive rough classification rules from large low-quality numerical databases that appear to be able to overcome these two problems. The algorithm works on the fields of continuous numeric variables; that is, the intervals of possible

values of each attribute in the training set, rather than on individual point values. The discovered rule is in a form called β -rule and is somewhat analogous to a decision tree found by an induction algorithm. The algorithm is linear in both the number of attributes and the number of instances (Dai & Ciesieski, 1994a, 2004).

The advantage of this inexact field-learning approach is its capability of inducing high-quality classification rules from low-quality data and its high efficiency that makes it an ideal algorithm to discover reliable knowledge from large and very large low-quality databases suitable for data mining, which needs higher discovering capability.

INEXACT FIELD-LEARNING ALGORITHM

Detailed description and the applications of the algorithm can be found from the listed articles (Ciesieski & Dai, 1994a; Dai & Ciesieski, 1994a, 1994b, 1995, 2004; Dai, 1996; Dai & Li, 2001; Dai, 1996). The following is a description of the inexact field-learning algorithm, the Fish_net algorithm:

Input: The input of the algorithm is a training data set with m instances and n attributes as follows:

Instances	X_1	X_2	...	X_n	Classes
Instance ₁	a_{11}	a_{12}	...	a_{1n}	γ_1
Instance ₂	a_{21}	a_{22}	...	a_{2n}	γ_2
...
Instance _m	a_{m1}	a_{m2}	...	a_{mn}	γ_m

(1)

Learning Process:

- **Step 1:** Work Out Fields of each attribute $\{x_i | 1 \leq i \leq n\}$ with respect to each class.

$$h_j^{(k)} = [h_{j_i}^{(k)}, h_{j_n}^{(k)}] \quad (k = 1, 2, \dots, s; j = 1, 2, \dots, n). \quad (2)$$

$$h_u^{(k)} = \max_{a_j \in a_j^{(k)}} \{a_{ij} \mid i = 1, 2, \dots, m\} \quad (3)$$

$$(k = 1, 2, \dots, s; j = 1, 2, \dots, n)$$

$$h_j^{(k)} = \min_{a_j \in a_j^{(k)}} \{a_{ij} \mid i = 1, 2, \dots, m\} \quad (4)$$

$$(k = 1, 2, \dots, s; j = 1, 2, \dots, n).$$

- **Step 2:** Construct Contribution Function based on the fields found in Step 1.

$$\mu_{c_k}(x_j) = \begin{cases} 0 & x_j \in \bigcup_{i \neq k}^s h_j^{(i)} - h_j^{(k)} \\ 1 & x_j \in h_j^{(k)} - \bigcup_{i \neq k}^s h_j^{(i)} \\ \frac{x_j - a}{b - a} & x_j \in h_j^{(k)} \cap \left(\bigcup_{i \neq k}^s h_j^{(i)}\right) \end{cases} \quad (5)$$

$$(k = 1, 2, \dots, s)$$

The formula (5) is given on the assumption that $[a, b] = h_j^{(k)} \cap \left(\bigcup_{i \neq k}^s h_j^{(i)}\right)$, and for any small number $\varepsilon > 0, b \pm \varepsilon \in h_j^{(k)}$ and $a + \varepsilon \notin \bigcup_{i \neq k}^s h_j^{(i)}$ or $a - \varepsilon \notin \bigcup_{i \neq k}^s h_j^{(i)}$. Otherwise, the formula (5) becomes,

$$\mu_{c_k}(x_j) = \begin{cases} 0 & x_j \in \bigcup_{i \neq k}^s h_j^{(i)} - h_j^{(k)} \\ 1 & x_j \in h_j^{(k)} - \bigcup_{i \neq k}^s h_j^{(i)} \\ \frac{x_j - b}{a - b} & x_j \in h_j^{(k)} \cap \left(\bigcup_{i \neq k}^s h_j^{(i)}\right) \end{cases} \quad (6)$$

$$(k = 1, 2, \dots, s)$$

- **Step 3:** Work Out Contribution Fields by applying the constructed contribution functions to the training data set.

- Calculate the contribution of each instance.

$$\alpha(I_i) = \left(\sum_{j=1}^n \mu(x_{ij})\right) / n \quad (7)$$

$$(i = 1, 2, \dots, m)$$

- Work out the contribution field for each class.
 $h^+ = \langle h_l^+, h_u^+ \rangle$

$$h_u^{(+)} = \max_{\alpha(I_i), I_i \in +} \{\alpha(I_i) \mid i = 1, 2, \dots, m\} \quad (8)$$

$$h_l^{(+)} = \min_{\alpha(I_i), I_i \in +} \{\alpha(I_i) \mid i = 1, 2, \dots, m\} \quad (9)$$

Similarly we can find $h^- = \langle h_l^-, h_u^- \rangle$

- **Step 4:** Construct Belief Function using the derived contribution fields.

$$B_i(C) = \begin{cases} -1 & \text{Contribution} \in \text{NegativeRegion} \\ 1 & \text{Contribution} \in \text{PositiveRegion} \\ \frac{c-a}{b-a} & \text{Contribution} \in \text{RoughRegion} \end{cases} \quad (10)$$

- **Step 5:** Decide Threshold. It could have 6 different cases to be considered. The simplest case is to take the threshold

$$\alpha = \text{midpoint of } h^+ \text{ and } h^- \quad (11)$$

- **Step 6:** Form the Inexact Rule.

$$\text{If } \bar{\alpha}(I) = \frac{1}{N} \sum_{i=1}^N \mu(x_i) > \alpha \quad (12)$$

$$\text{Then } \gamma = 1(B_i(c))$$

This algorithm was tested on three large real observational weather data sets containing both high-quality and low-quality data. The accuracy rates of the forecasts were 86.4%, 78%, and 76.8%. These are significantly better than the accuracy rates achieved by C4.5 (Quinlan, 1986, 1993), feed forward neural networks, discrimination analysis, K-nearest neighbor classifiers, and human weather forecasters. The fish-net algorithm exhibited significantly less overfitting than the other algorithms. The training times were shorter, in some cases by orders of magnitude (Dai & Ciesieski, 1994a, 2004; Dai 1996).

FUTURE TRENDS

The inexact field-learning approach has led to a successful algorithm in a domain where there is a high level of noise. We believe that other algorithms based on fields also can be developed. The β -rules, produced by the current FISH-NET algorithm involve linear combinations of attributes. Non-linear rules may be even more accurate.

2 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/inexact-field-learning-approach-data/10669

Related Content

An Introduction to Information Technology and Business Intelligence

Stephan Kudyba and Richard Hoptroff (2002). *Data Warehousing and Web Engineering* (pp. 1-21). www.irma-international.org/chapter/introduction-information-technology-business-intelligence/7860

Hyperbolic Space for Interactive Visualization

Jörg Andreas Walter (2005). *Encyclopedia of Data Warehousing and Mining* (pp. 575-581). www.irma-international.org/chapter/hyperbolic-space-interactive-visualization/10663

Analytical Customer Requirement Analysis Based on Data Mining

Jianxin ("Roger") Jiao, Yiyang Zhang and Martin Helander (2008). *Data Warehousing and Mining: Concepts, Methodologies, Tools, and Applications* (pp. 2798-2815). www.irma-international.org/chapter/analytical-customer-requirement-analysis-based/7801

Unsupervised Mining of Genes Classifying Leukemia

Diego Liberati, Sergio Bittanti and Simone Garatti (2005). *Encyclopedia of Data Warehousing and Mining* (pp. 1155-1159). www.irma-international.org/chapter/unsupervised-mining-genes-classifying-leukemia/10771

Business Data Warehouse: The Case of Wal-Mart

Indranil Bose, Lam Albert Kar Chun, Leung Vivien Wai Yue, Li Hoi Wan Ines and Wong Oi Ling Helen (2008). *Data Warehousing and Mining: Concepts, Methodologies, Tools, and Applications* (pp. 2762-2771). www.irma-international.org/chapter/business-data-warehouse/7798