

ПОЛНАЯ ВЗАИМНАЯ ИНФОРМАЦИЯ КАК ОБОБЩЕННЫЙ ПОКАЗАТЕЛЬ КАЧЕСТВА АССОЦИАТИВНЫХ ЗАВИСИМО- СТЕЙ С БИНАРНЫМИ ПРИЗНАКАМИ

к.т.н. Д.Э. Ситников, Е.В. Титова
(представил д.т.н., проф. О.И. Сухаревский)

Проводится анализ характеристик ассоциативных правил при выделении знаний из данных (data mining). Доказывается необходимость и достаточность трех показателей: уровня поддержки, доверия и улучшения для определения взаимодействия двух частей ассоциативного правила. Предлагается интегральный показатель качества ассоциации – полная взаимная информация.

Введение. Обнаружение знаний в данных определяют как процесс выявления и идентификации в больших массивах информации новых, потенциально полезных и понятных зависимостей. В настоящее время это направление (англ. knowledge discovery in databases – KDD and data mining) интенсивно развивается, расширяется сфера его применений: в здравоохранении, биохимии, в социологии, финансовой, кредитной и маркетинговой деятельности, в правительственных службах, технической диагностике, в управлении производством и т.д.

Под термином data mining понимают центральный этап процесса выявления знаний в данных, т.е. собственно обнаружение зависимостей, регулярностей и отношений, который опирается на определенную форму повторяемости информации в базах данных. Операции этого этапа включают поиск и конструирование зависимостей, их оценку и отсеивание.

Одной из типичных задач data mining является задача открытия (генерации) правил ассоциации. Ассоциативное правило определяется как утверждение импликативного вида $A \rightarrow B$, где A и B – некоторые множества признаков объектов (записей в базе данных), причем данное утверждение должно иметь меру определенности не ниже заданного уровня.

Постановка проблемы. Одной из проблем, возникающих на этапе обработки полученных зависимостей – очистка результатов, их оценка, удаление частных и избыточных знаний. Иначе существует угроза астрономического изобилия обнаруживаемых зависимостей, подавляющая

часть которых оказывается ненужной [1].

Анализ литературы показывает, что предлагаемые меры качества ассоциативных правил: "интерес", корреляция [1], обобщенная мера ассоциации [2] и др. не всегда в достаточной мере характеризуют ассоциативную зависимость и имеют свои недостатки.

Цель статьи. Выработка интегрального показателя качества ассоциативных зависимостей при выделении знаний из данных, опирающегося на стандартные характеристики – уровень поддержки, доверие и улучшение.

Рассмотрим общепринятые характеристики ассоциативных зависимостей: поддержку (Support – сокр. Sup) и доверие (Confidence – сокр. Conf).

Поддержка определяется как отношение количества записей в базе данных (БД), удовлетворяющих этому правилу к общему количеству записей в БД (может определяться просто как количество записей в БД, удовлетворяющих данному правилу). С точки зрения теории вероятности поддержка определяет вероятность наличия у объекта (записи в БД) двух признаков А и В: $P(AB)$. Доверие (вероятность правила) определяют как отношение количества записей, удовлетворяющих правилу к количеству записей, удовлетворяющих посылке (левой части). С точки зрения теории вероятности уровень доверия определяет условную вероятность $P_A(B) = \frac{P(AB)}{P(A)}$. Если $A \rightarrow B$ – является

ассоциативным правилом, то должны выполняться следующие условия:

$$\text{Sup}(A \rightarrow B) \geq \min\text{Sup};$$

$$\frac{\text{Sup}(A \rightarrow B)}{\text{Sup}(A)} \geq \min\text{Conf},$$

где $\text{Sup}(A) = M(A) / M(\text{БД})$ – отношение мощности множества объектов в БД, обладающих признаком А к количеству объектов в БД; $\text{Sup}(A \rightarrow B) = \text{Sup}(AB) = M(A \cap B) / M(\text{БД})$ – отношение мощности множества объектов в БД, обладающих признаками А и В к количеству объектов в БД; $\min\text{Sup}$, $\min\text{Conf}$ – установленные минимальные уровни поддержки и доверия.

Однако, эти две характеристики, будучи очень важными, не позволяют в достаточной степени оценить качество ассоциативного правила. Третьей характеристикой, предлагаемой для оценки ассоциации, является Improvement (сокр. Imp) – дословно “улучшение”. Imp определяется как отношение вероятности правила к вероятности результата

$$\text{Imp}(A \rightarrow B) = \frac{P_A(B)}{P(B)} = \frac{\text{Conf}(A \rightarrow B)}{\text{Sup}(B)}.$$

Таким образом, Improvement характеризует степень повышения вероятности события В при условии наступления события А относительно простой

вероятности события В.

О важности Imp говорит хотя бы тот факт, что $\text{Imp}(A \rightarrow B) = 1$ показывает, что $P_A(B) = P(B)$, т.е. признаки А и В независимы. В этом случае даже при хороших показателях уровней поддержки и доверия ассоциативное правило является абсолютно бессмысленным.

Однако этой характеристике не всегда уделяется достаточно внимания в отличие от Sup и Conf (например, в одной из известных систем выделения знаний из данных WizWhy).

Покажем, что для полной оценки взаимодействия двух признаков (двух частей ассоциативного правила) необходимы и достаточны все три характеристики.

Система двух величин (в данном случае двух признаков, двух частей ассоциативного правила) с точки зрения теории вероятности характеризуется следующими значениями вероятностей: $P_{00}, P_{01}, P_{10}, P_{11}$.

Зная три характеристики ассоциативного правила $A \rightarrow B$: Support, Confidence и Improvement, мы можем однозначно определить P_{00}, P_{01}, P_{10} и P_{11} .

Для начала выразим Support, Confidence и Improvement правила $A \rightarrow B$ через вероятности:

$$\text{Sup}(A \rightarrow B) = P(AB) = P_{11}; \quad (1)$$

$$\text{Conf}(A \rightarrow B) = \frac{\text{Sup}(A \rightarrow B)}{\text{Sup}(A)} = \frac{P(AB)}{P(A)} = \frac{P_{11}}{P_{11} + P_{10}}; \quad (2)$$

$$\text{Imp}(A \rightarrow B) = \frac{\text{Sup}(A \rightarrow B)}{\text{Sup}(A)\text{Sup}(B)} = \frac{P(AB)}{P(A)P(B)} = \frac{P_{11}}{(P_{11} + P_{10})(P_{11} + P_{01})}. \quad (3)$$

Отсюда:

$$P_{11} = \text{Sup}(A \rightarrow B); \quad (4)$$

$$P_{10} = \frac{\text{Sup}(A \rightarrow B)}{\text{Conf}(A \rightarrow B)} - \text{Sup}(A \rightarrow B) = \text{Sup}(A \rightarrow \bar{B}); \quad (5)$$

$$P_{01} = \frac{\text{Conf}(A \rightarrow B)}{\text{Imp}(A \rightarrow B)} - \text{Sup}(A \rightarrow B) = \text{Sup}(\bar{A} \rightarrow B); \quad (6)$$

$$\begin{aligned} P_{00} &= 1 - P_{11} - P_{10} - P_{01} = \\ &= 1 + \text{Sup}(A \rightarrow B) - \frac{\text{Sup}(A \rightarrow B)}{\text{Conf}(A \rightarrow B)} - \frac{\text{Conf}(A \rightarrow B)}{\text{Imp}(A \rightarrow B)} = \text{Sup}(\bar{A} \rightarrow \bar{B}). \end{aligned} \quad (7)$$

Таким образом, для определения всех значений $P_{00}, P_{01}, P_{10}, P_{11}$ требуются все три характеристики ассоциативного правила: Sup, Conf, Imp и при отсутствии хотя бы одной из них все вероятности $P_{00}, P_{01}, P_{10}, P_{11}$ не могут быть получены.

Следует отметить, что на численные значения всех трех величин

Sup, Conf и Imp накладываются определенные ограничения:

- 1) $0 \leq \text{Sup}(A \rightarrow B) \leq 1$.
- 2) $\text{Sup}(A \rightarrow B) \leq \text{Conf}(A \rightarrow B) \leq 1$.

Действительно, $\text{Sup}(A \rightarrow B) = P_{11}$; $\text{Conf}(A \rightarrow B) = \frac{P_{11}}{P_{11} + P_{10}}$.

Но $P_{11} \leq \frac{P_{11}}{P_{11} + P_{10}}$, т.к. $P_{11} + P_{10} \leq 1$.

- 3) при заданных значениях $\text{Sup}(A \rightarrow B)$ и $\text{Conf}(A \rightarrow B)$

$$\frac{(\text{Conf}(A \rightarrow B))^2}{\text{Sup}(A \rightarrow B) \cdot \text{Conf}(A \rightarrow B) + \text{Conf}(A \rightarrow B) - \text{Sup}(A \rightarrow B)} \leq \leq \text{Imp}(A \rightarrow B) \leq \frac{\text{Conf}(A \rightarrow B)}{\text{Sup}(A \rightarrow B)}.$$

Доказательство третьего неравенства будем вести, исходя из определения Improvement

$$\text{Imp}(A \rightarrow B) = \frac{P_A(B)}{P(B)} = \frac{\text{Conf}(A \rightarrow B)}{\text{Sup}(B)} = \frac{\text{Sup}(AB)}{\text{Sup}(A)\text{Sup}(B)}.$$

Если $\text{Sup}(B)$ – минимальна, т.е. $\text{Sup}(B) = \text{Sup}(AB)$, то $\text{Imp}(A \rightarrow B)$ – максимален и равен

$$\text{Imp}(A \rightarrow B) = \frac{1}{\text{Sup}(A)} = \frac{\text{Conf}(A \rightarrow B)}{\text{Sup}(A \rightarrow B)}.$$

Если $\text{Sup}(B)$ – максимальна, т.е. $\text{Sup}(B) = 1 - \text{Sup}(A) + \text{Sup}(AB)$, то $\text{Imp}(A \rightarrow B)$ – минимален и равен

$$\text{Imp}(A \rightarrow B) = \frac{\text{Sup}(AB)}{\text{Sup}(A)(1 - \text{Sup}(A) + \text{Sup}(AB))}.$$

Но $\text{Sup}(A) = \frac{\text{Sup}(A \rightarrow B)}{\text{Conf}(A \rightarrow B)}$, откуда

$$\text{Imp}(A \rightarrow B) = \frac{(\text{Conf}(A \rightarrow B))^2}{\text{Sup}(A \rightarrow B) \cdot \text{Conf}(A \rightarrow B) + \text{Conf}(A \rightarrow B) - \text{Sup}(A \rightarrow B)}.$$

В процессе обработки полученных правил встает вопрос оценки качества ассоциаций. Каждая из приведенных выше величин, безусловно, характеризует полученную зависимость, однако необходима выработка некоторой интегральной характеристики, которая учитывала бы все три параметра. Как сравнить, например, два ассоциативных правила, у одного из которых больше Sup и Imp, а у другого – Conf? Какое из них “лучше” и насколько?

В [2] предложена обобщенная мера ассоциации

$$\Delta = \frac{P_{00} \cdot P_{11}}{P_{01} \cdot P_{10}}, \quad (8)$$

которая характеризует отношение вероятности сходства признаков к вероятности их отличия.

Безусловно, эта мера может быть выражена через три характеристики ассоциативного правила. Подставляя в формулу (8) значения для P_{00} , P_{01} , P_{10} , P_{11} из формул (4 – 7), получим

$$\Delta = \frac{\text{Sup}(A \rightarrow B) \cdot \left(1 + \text{Sup}(A \rightarrow B) - \frac{\text{Sup}(A \rightarrow B)}{\text{Conf}(A \rightarrow B)} - \frac{\text{Conf}(A \rightarrow B)}{\text{Imp}(A \rightarrow B)} \right)}{\left(\frac{\text{Sup}(A \rightarrow B)}{\text{Conf}(A \rightarrow B)} - \text{Sup}(A \rightarrow B) \right) \cdot \left(\frac{\text{Conf}(A \rightarrow B)}{\text{Imp}(A \rightarrow B)} - \text{Sup}(A \rightarrow B) \right)}.$$

Однако, предлагаемая мера обладает существенным недостатком: она не всегда может быть рассчитана, например, если хотя бы одна из вероятностей P_{00} , P_{01} , P_{10} , P_{11} равна нулю. Действительно, пусть:

$$P_{00} = 0,4;$$

$$P_{01} = 0,1;$$

$$P_{10} = 0,2;$$

$$P_{11} = 0,3.$$

Рассчитываем характеристики Sup , Conf , Imp , Δ , пользуясь формулами (1, 2, 3, 8): $\text{Sup} = 0,3$; $\text{Conf} = 0,6$; $\text{Imp} = 1,5$.

$$\Delta = \frac{P_{00} \cdot P_{11}}{P_{01} \cdot P_{10}} = \frac{0,3 \cdot 0,4}{0,1 \cdot 0,2} = 6.$$

Пусть теперь:

$$P_{00} = 0,4;$$

$$P_{01} = 0,1;$$

$$P_{10} = 0,1;$$

$$P_{11} = 0,4.$$

Отметим, что ассоциация усилилась: вероятность сходства признаков возросла, соответственно возросли все три характеристики.

$$\text{Sup} = 0,4; \text{Conf} = 0,8; \text{Imp} = 1,6.$$

$$\Delta = \frac{P_{00} \cdot P_{11}}{P_{01} \cdot P_{10}} = \frac{0,4 \cdot 0,4}{0,1 \cdot 0,1} = 16.$$

До сих пор предложенная мера ассоциации работает нормально: она увеличивается.

Теперь еще усиливаем ассоциацию:

$$P_{00} = 0,4;$$

$$P_{01} = 0;$$

$$P_{10} = 0,1;$$

$$P_{11} = 0,5.$$

$$\text{Sup} = 0,5; \text{Conf} = 0,8(3); \text{Imp} = 1,0(6).$$

Теперь обобщенная мера ассоциации Δ не может быть рассчитана из-за равенства нулю знаменателя. Простое отбрасывание нулевой вероятности приводит к абсурдным результатам: $\Delta = \frac{P_{00} \cdot P_{11}}{P_{10}} = \frac{0,4 \cdot 0,5}{0,1} = 2$,

т.е. при усилении ассоциации обобщенная мера уменьшилась.

Таким образом, необходим такой показатель качества ассоциации, который, во-первых, включал бы в себя все три характеристики ассоциативного правила, во-вторых, мог бы быть рассчитан при любых значениях вероятностей P_{00} , P_{01} , P_{10} , P_{11} , в-третьих, его численное значение должно возрастать при усилении ассоциации и, соответственно, уменьшаться при ослаблении.

Показатель, который отвечает всем поставленным требованиям – взаимная информация. Действительно, ассоциативное правило $A \rightarrow B$ с точки зрения теории информации можно рассматривать, как информацию о событии B , получаемую в результате сообщения о событии A .

Информация “от события к событию” определяется как $I_{A \rightarrow B} = \log \frac{P(AB)}{P(A)P(B)}$ [3], что есть не что иное, как $\log(\text{Imp}(A \rightarrow B))$.

Полная взаимная информация определяется [3]

$$I_{A \leftrightarrow B} = \sum_{i=1}^n \sum_{j=1}^m P_{ij} \log \frac{P_{ij}}{P_i P_j},$$

где $P_{ij} = P((A \sim a_i)(B \sim b_j))$ – вероятность того, что A находится в состоянии a_i , а B – в состоянии b_j ; $P_i = P(A \sim a_i)$ – вероятность того, что A находится в состоянии a_i ; $P_j = P(B \sim b_j)$ – вероятность того, что B находится в состоянии b_j .

Таким образом, для нашего случая формула полной взаимной информации примет вид

$$\begin{aligned} I = & \text{Sup}(A \rightarrow B) \cdot \log(\text{Imp}(A \rightarrow B)) + \\ & + \text{Sup}(A \rightarrow \bar{B}) \cdot \log(\text{Imp}(A \rightarrow \bar{B})) + \\ & + \text{Sup}(\bar{A} \rightarrow B) \cdot \log(\text{Imp}(\bar{A} \rightarrow B)) + \text{Sup}(\bar{A} \rightarrow \bar{B}) \cdot \log(\text{Imp}(\bar{A} \rightarrow \bar{B})). \end{aligned} \quad (9)$$

Как было показано ранее, три характеристики ассоциативного правила Sup , Conf и Imp полностью определяют систему двух величин. Следовательно, все параметры в формуле (9) могут быть выражены че-

рез них. Выражения для $\text{Sup}(\bar{A} \rightarrow \bar{B})$, $\text{Sup}(\bar{A} \rightarrow B)$, $\text{Sup}(A \rightarrow \bar{B})$ были получены в формулах (5 – 7). Выразим величины $\text{Imp}(\bar{A} \rightarrow \bar{B})$, $\text{Imp}(\bar{A} \rightarrow B)$, $\text{Imp}(A \rightarrow \bar{B})$ через известные характеристики ассоциативного правила:

$$\text{Imp}(\bar{A} \rightarrow B) = \frac{\text{Sup}(\bar{A} \rightarrow B)}{\text{Sup}(\bar{A})\text{Sup}(B)} = \frac{P_{01}}{(P_{01} + P_{00})(P_{01} + P_{11})}; \quad (10)$$

$$\text{Imp}(A \rightarrow \bar{B}) = \frac{\text{Sup}(A \rightarrow \bar{B})}{\text{Sup}(A)\text{Sup}(\bar{B})} = \frac{P_{10}}{(P_{10} + P_{11})(P_{10} + P_{00})}; \quad (11)$$

$$\text{Imp}(\bar{A} \rightarrow \bar{B}) = \frac{\text{Sup}(\bar{A} \rightarrow \bar{B})}{\text{Sup}(\bar{A})\text{Sup}(\bar{B})} = \frac{P_{00}}{(P_{00} + P_{01})(P_{00} + P_{10})}. \quad (12)$$

Подставим в формулы (10 – 12) выражения для P_{01} , P_{00} , P_{10} , P_{11} из формул (4 – 7) и проведем алгебраические преобразования. Получим:

$$\text{Imp}(\bar{A} \rightarrow B) = \frac{\text{Conf}(A \rightarrow B) - \text{Sup}(A \rightarrow B) \cdot \text{Imp}(A \rightarrow B)}{\text{Conf}(A \rightarrow B) - \text{Sup}(A \rightarrow B)}; \quad (13)$$

$$\text{Imp}(A \rightarrow \bar{B}) = \frac{\text{Imp}(A \rightarrow B) (1 - \text{Conf}(A \rightarrow B))}{\text{Imp}(A \rightarrow B) - \text{Conf}(A \rightarrow B)}; \quad (14)$$

$$\text{Imp}(\bar{A} \rightarrow \bar{B}) = \frac{1 + \text{Sup}(A \rightarrow B) - \frac{\text{Sup}(A \rightarrow B)}{\text{Conf}(A \rightarrow B)} - \frac{\text{Conf}(A \rightarrow B)}{\text{Imp}(A \rightarrow B)}}{\left(1 - \frac{\text{Sup}(A \rightarrow B)}{\text{Conf}(A \rightarrow B)}\right) \left(1 - \frac{\text{Conf}(A \rightarrow B)}{\text{Imp}(A \rightarrow B)}\right)}. \quad (15)$$

Теперь можно записать несколько громоздкую формулу для полной взаимной информации, которая выражается через три характеристики ассоциативного правила: Sup , Conf и Imp .

$$\begin{aligned} I = & \text{Sup}(A \rightarrow B) \cdot \log(\text{Imp}(A \rightarrow B)) + \left(\frac{\text{Sup}(A \rightarrow B)}{\text{Conf}(A \rightarrow B)} - \text{Sup}(A \rightarrow B) \right) \times \\ & \times \log \left(\frac{\text{Imp}(A \rightarrow B) (1 - \text{Conf}(A \rightarrow B))}{\text{Imp}(A \rightarrow B) - \text{Conf}(A \rightarrow B)} \right) + \left(\frac{\text{Conf}(A \rightarrow B)}{\text{Imp}(A \rightarrow B)} - \text{Sup}(A \rightarrow B) \right) \times \\ & \times \log \left(\frac{\text{Conf}(A \rightarrow B) - \text{Sup}(A \rightarrow B) \cdot \text{Imp}(A \rightarrow B)}{\text{Conf}(A \rightarrow B) - \text{Sup}(A \rightarrow B)} \right) + \end{aligned} \quad (16)$$

$$+ \left(1 + \text{Sup}(A \rightarrow B) - \frac{\text{Sup}(A \rightarrow B)}{\text{Conf}(A \rightarrow B)} - \frac{\text{Conf}(A \rightarrow B)}{\text{Imp}(A \rightarrow B)} \right) \times$$

$$\times \log \left(\frac{1 + \text{Sup}(A \rightarrow B) - \frac{\text{Sup}(A \rightarrow B)}{\text{Conf}(A \rightarrow B)} - \frac{\text{Conf}(A \rightarrow B)}{\text{Imp}(A \rightarrow B)}}{\left(1 - \frac{\text{Sup}(A \rightarrow B)}{\text{Conf}(A \rightarrow B)} \right) \left(1 - \frac{\text{Conf}(A \rightarrow B)}{\text{Imp}(A \rightarrow B)} \right)} \right).$$

Проверим, как ведет себя полная взаимная информация при разных значениях P_{01} , P_{00} , P_{10} , P_{11} . Обратимся к нашему примеру.

При $\text{Sup} = 0,3$; $\text{Conf} = 0,6$; $\text{Imp} = 1,5$ полная взаимная информация, рассчитанная по формуле (16), равна $I = 0,125$.

При $\text{Sup} = 0,4$; $\text{Conf} = 0,8$; $\text{Imp} = 1,6$ информация растёт: $I = 0,278$.

Когда $P_{01} = 0$ в формуле (16) слагаемое, соответствующее этой вероятности, перестаёт работать (отбрасывается). Формула принимает вид:

$$I = \text{Sup}(A \rightarrow B) \cdot \log(\text{Imp}(A \rightarrow B)) + \text{Sup}(A \rightarrow \bar{B}) \cdot \log(\text{Imp}(A \rightarrow \bar{B})) +$$

$$+ \text{Sup}(\bar{A} \rightarrow \bar{B}) \cdot \log(\text{Imp}(\bar{A} \rightarrow \bar{B}))$$

или

$$I = \text{Sup}(A \rightarrow B) \cdot \log(\text{Imp}(A \rightarrow B)) + \left(\frac{\text{Sup}(A \rightarrow B)}{\text{Conf}(A \rightarrow B)} - \text{Sup}(A \rightarrow B) \right) \times$$

$$\times \log \left(\frac{\text{Imp}(A \rightarrow B) (1 - \text{Conf}(A \rightarrow B))}{\text{Imp}(A \rightarrow B) - \text{Conf}(A \rightarrow B)} \right) +$$

$$+ \left(1 + \text{Sup}(A \rightarrow B) - \frac{\text{Sup}(A \rightarrow B)}{\text{Conf}(A \rightarrow B)} - \frac{\text{Conf}(A \rightarrow B)}{\text{Imp}(A \rightarrow B)} \right) \times$$

$$\times \log \left(\frac{1 + \text{Sup}(A \rightarrow B) - \frac{\text{Sup}(A \rightarrow B)}{\text{Conf}(A \rightarrow B)} - \frac{\text{Conf}(A \rightarrow B)}{\text{Imp}(A \rightarrow B)}}{\left(1 - \frac{\text{Sup}(A \rightarrow B)}{\text{Conf}(A \rightarrow B)} \right) \left(1 - \frac{\text{Conf}(A \rightarrow B)}{\text{Imp}(A \rightarrow B)} \right)} \right).$$

Для нашего примера при $P_{01} = 0$: $\text{Sup} = 0,5$; $\text{Conf} = 0,8(3)$; $\text{Imp} = 1,0(6)$.

$I = 0,612$ – взаимная информация, как показатель качества ассоциации, выросла.

Отметим, что при $\text{Imp}(A \rightarrow B) = 1$ (т.е. при независимости A и B) полная взаимная информация равна нулю, что является вполне логичным результатом.

Своего максимума $I = 1$ достигает при $P_{00} = 0,5$; $P_{01} = 0$; $P_{10} = 0$; $P_{11} = 0,5$: вероятность сходства максимальна $\text{Sup} = 0,5$; $\text{Conf} = 1$; $\text{Imp} = 2$. Безусловно,

в реальных базах данных получение таких показателей для ассоциативного правила вряд ли возможно.

Хорошими показателями поддержки и доверия, как показывает опыт, являются $\text{Sup} = 0,1 \div 0,3$; $\text{Conf} = 0,7 \div 0,8$. При $\text{Sup} = 0,3$; $\text{Conf} = 0,8$ и максимальном $\text{Imp} = 8/3$, полная взаимная информация достигает $I \approx 0,58$.

При среднем $\text{Imp}(A \rightarrow B) = \frac{\max \text{Imp}(A \rightarrow B) + 1}{2}$, где единица соответствует $I = 0$, значения полной взаимной информации находятся в пределах $I \approx 0,2 \div 0,4$.

Надо отметить, что вместо задания трех параметров при выявлении ассоциации в БД, можно выставлять требования на значения полной взаимной информации (т.е. производить отсев правил, информативность которых меньше заданного порога).

Выводы. В статье проанализированы общепринятые характеристики ассоциативных правил: уровень поддержки, доверие и улучшение. Показана важность третьей характеристики, которой не всегда уделяется достаточно внимания. Доказана необходимость и достаточность этих трех показателей для описания взаимодействия двух частей ассоциативной зависимости. Предложен интегральный показатель качества ассоциации – полная взаимная информация, который рассчитывается при любых показателях системы двух величин и позволяет количественно сравнивать полученные правила. Кроме этого, информативность правила (полная взаимная информация) может задаваться в качестве ограничения при генерации ассоциативных зависимостей во время выделения знаний из данных.

ЛИТЕРАТУРА

1. Балабанов А.С. Выделение знаний из баз данных – передовые компьютерные технологии интеллектуального анализа данных // Математичні машини і системи. – 2001. – № 1, 2. – С. 40 – 54.
2. Edwards A.W.F. The measure of association in 2x2 table // Journal of the Royal Statistical Society. ser. A29. – P. 109 – 114.
3. Вентцель Е.С. Теория вероятностей. – М.: Наука, 1964. – 576 с.

Поступила 17.12.2003

СИТНИКОВ Дмитрий Эдуардович, канд. техн. наук, доцент, зав. кафедрой информационно-документных систем ХГАК. В 1988 году окончил Харьковский институт радиоэлектроники. Область научных интересов – выделение знаний из данных, нечеткие множества.

ТИТОВА Елена Витольдиевна, младший научный сотрудник НЦ Войск ПВО, аспирантка кафедры информационно-документных систем ХГАК. В 1988 году окончила

Харьковский институт радиоэлектроники. Область научных интересов – выделение знаний из данных.
