

Лекция 14. Извлечение знаний из эксперта.

Проблемы извлечения знаний из эксперта.

1. Невозможность сформулировать интуитивные знания;
2. Невозможность задать большое число вопросов. Например, для 11 бинарных диагностических признаков сгруппированных кальцинозов есть ($2^{11} = 2\,048$) комбинаций признаков, каждый из которых представляет новый случай. Лобовой метод потребовал бы опроса радиолога для каждой из этих 2 048 комбинаций;
3. Возможная сложность правил: Обычно порядка 60–70 % времени при разработке системы, основанной на правилах, тратится на извлечение знаний; Эксперт обдумывает альтернативные сценарии и, говорит: «Я думаю, что при обстоятельствах, X, наиболее вероятное заключение – Y, но если есть дополнительный факт, скажем F, то более вероятное заключение могло бы быть P»;
4. Сложность точной формулировки правил ввиду массы *неявных предположений*;
5. Необходимость *отладки экспертных систем*;
6. *Экспертное мнение субъективно*, что бы проверить его на *объективность и дополнить* экспертное знание, надо применить систему Discovery и обнаружить правила на данных;
7. *Идентифицировать противоречия* между экспертными правилами и правилами, извлеченными из данных:
 - a. обнаруженное правило похоже на экспертное. Эксперт может проверить:
 - i. Подтверждает ли правило существующее экспертное знание?
 - ii. Если правило содержит меньше признаков, чем экспертное, то эксперт может найти, что правило совместимо с его/ее предыдущим опытом, но он/она хотел ли бы, чтобы оно было более надежным.
 - iii. Имеющихся данных недостаточно для достоверного обнаружения правила и возможно эксперт преувеличивает значимость признака либо у него есть дополнительная информация, не содержащаяся в данных, для такого правила;
 - iv. Если правило содержит признаки, не содержащиеся в экспертных правилах, то эксперт либо не учитывает некоторые важные данные, либо сами данные односторонни и их надо расширить;
 - b. обнаружение правил в данных, не обнаруженных в процессе опроса эксперта;
 - i. мнение эксперта односторонне и требует пересмотра. Система улучшает опыт эксперта;
 - ii. данные собраны односторонне и рассматриваемые случаи надо расширить или классифицировать;
 - c. обнаружены правила, которые противоречат его/ее знанию или пониманию:
 - i. правило было обнаружено путем использования вводящих в заблуждение случаев. Правило должно быть отклонено и обучающиеся данные должны быть расширены;
 - ii. Эксперт может признать, что его/ее знания не имеют под собой реального основания и основаны на некоторых теоретических соображениях, требующих дополнительную проверку или пересмотра. Система улучшает опыт эксперта;
8. *Устранить противоречия* и получить полную и совместную базу знаний и экспертную систему, в которой улучшены как экспертные правила, так и правила извлеченные из данных.
9. Если эксперт может ясно сформулировать процесс принятия решений, то подход, основанный на правилах, подходит для создания экспертной системы.

Метод извлечения диагностических правил из эксперта.

Иерархический подход. Опрос эксперта основанный на оригинальном методе восстановления монотонных Булевых функций.

Можно попросить эксперта, что бы он описал случаи множеством бинарных признаков.

Типичный вопрос будет иметь следующий формат:

Если признак 1 имеет значение V_1 , признак 2, имеет значение V_2 ..., признак n имеет значение V_n , то нужно ли рекомендовать:

- биопсию или нет?
- либо набор значений признаков соответствует случаю подозрительному к раку или нет?

Мы строим иерархию медицинских интерпретируемых признаков, начиная с обобщенного уровня до все менее обобщенного уровня.

Эта иерархия начинается с определения 11 медицинских бинарных признаков.

Медик-эксперт определил, что первичные 11 бинарных признаков $w_1, w_2, w_3, u_1, u_2, u_3, u_4, u_5, x_3, x_4, x_5$ могут быть организованы в иерархию с добавлением двух новых обобщенных признаков x_1 и x_2 : $x_1 - w_1, w_2, w_3$; $x_2 - u_1, u_2, u_3, u_4, u_5$.

Новый обобщенный признак:

x_1 – «Количество и объем кальцинозов» со стадиями (0 – «доброкачественный» и 1 – «рак») был введен на основании признаков:

- w_1 – количество кальцинозов / см^3 ,
- w_2 – объем кальциноза, см^3 и
- w_3 – общее количество кальцинозов.

Мы рассматриваем признак x_1 как функцию $v(w_1, w_2, w_3)$, которую надо определить.

Аналогично, новый признак:

x_2 – «Форма и плотность кальциноза» со значениями: (1) – «рак» и (0) – «доброкачественная» является обобщением признаков:

- u_1 – «Нерегулярность в форме индивидуальных кальцинозов»,
- u_2 – «Изменение в форме кальцинозов»,
- u_3 – «Изменение в размере кальцинозов»,
- u_4 – «Изменение в плотности кальцинозов»,
- u_5 – «Плотность кальцинозов».

Мы рассматриваем x_2 как функцию $x_2 = \psi(u_1, u_2, u_3, u_4, u_5)$, которая должна быть идентифицирована для диагностики рака.

Мы рассматриваем пять бинарных признаков x_1, x_2, x_3, x_4 и x_5 , на уровне 1.

В результате мы получили декомпозицию задачи $f(x_1, x_2, x_3, x_4, x_5)$

$f(x_1, x_2, x_3, x_4, x_5) = f(v(w_1, w_2, w_3), \psi(u_1, u_2, u_3, u_4, u_5), x_3, x_4, x_5)$.

Будем предполагать, что наши признаки сводятся к следующим пяти признакам, имеющим значения 0 – «доброкачественный», 1 – «рак»:

- x_1 – [количество и объем, занятый кальцинозами];
- x_2 – [форма и плотность кальцинозов];
- x_3 – [ориентация протоков];
- x_4 – [сравнение с предыдущей экспертизой];
- x_5 – [ассоциированные результаты исследования].

Свойство монотонности

Если радиолог правильно диагностировал набор (10100) как злокачественный, то, используя свойство монотонности, мы можем также заключить, что клинический случай (10110) также должен быть злокачественным.

Медику-эксперту представили идеи относительно монотонности функций, как было определено выше. Кроме того, диалог, который следовал, подтверждал законность этого предположения. Точно так же функция $x_2 = \psi(y_1, y_2, y_3, y_4, y_5)$ для x_2 была подтверждена как монотонная Булева функция.

Булева функция – компактное представление набора диагностических правил. Булева дискриминантная функция может быть представлена в форме множества ЕСЛИ–ТО правил.

Таким образом, **основными шагами извлечения правил из медика-эксперта** являются следующие:

- разработать иерархию понятий и представить их как ряд монотонных Булевых функций;
- восстановить каждую из этих функций с минимальной последовательностью вопросов эксперту;
- объединить обнаруженные функции в полную диагностическую функцию;
- представить полную функцию как традиционный набор простых диагностических правил вида: *Если A и B и ... F TO Z.*

Опишем восстановления каждой монотонной Булевой функции с минимальной динамической последовательностью вопросов.

Эта последовательность основана на фундаментальной лемме Hansel [108]. т. е. минимальное количество вопросов обязано восстанавливать самую сложную монотонную Булеву функцию с n аргументами.

Табл. иллюстрирует процедуру.

Столбцы 3 и 4 представляют собой значения определенных выше функций f и ψ . Мы опускаем восстановление функции $v(w_1, w_2, w_3)$, потому что нужно немного вопросов для восстановления этой функции.

В таблице первый вопрос: «Представляет ли последовательность (01100) доброкачественный случай?» Здесь, $x_1 = 0$ и $(01100) = (x_1, x_2, x_3, x_4, x_5)$. Если ответ «да» (1), то следующий вопрос будет о доброкачественности случая (01010). Если ответ «нет» (0), то следующий вопрос будет о доброкачественности для случая (11100).

Все 32 возможных случая с пятью бинарными признаками (x_1, x_2, x_3, x_4, x_5) представлены в столбце 1 табл. Они сгруппированы в группы называемыми цепями Hansel.

Таблица. Динамическая последовательность интервью с экспертом

Дело	f Рак	ψ Форма и плотность кальцино- зов	Монотонное удлинение		Цепь	Дело
			$1 \rightarrow 1$	$0 \rightarrow 0$		
1	3	4	5	6	7	8
(01100)	1*	1*	1.2;6.3;7.3	7.1;8.1	Цепь 1	1.1
(11100)	1	1	6.4;7.4	5.1;3.1		1.2
(01010)	0*	1*	2.2;6.3;8.3	6.1;8.1	Цепь 2	2.1
(11010)	1*	1	6.4;8.4	3.1;6.1		2.2
(11000)	1*	1*	3.2	8.1;9.1	Цепь 3	3.1
(11001)	1	1	7.4;8.4	8.2;9.2		3.2
(10010)	0*	1*	4.2;9.3	6.1;9.1	Цепь 4	4.1
(10110)	1*	1	6.4;9.4	6.2;5.1		4.2
(10100)	1*	1*	5.2	7.1;9.1	Цепь 5	5.1
(10101)	1	1	7.4;9.4	7.2;9.2		5.2
(00010)	0	0*	6.2;10.3	10.1	Цепь 6	6.1
(00110)	1*	0*	6.3;10.4	7.1		6.2

(01110)	1	1	6.4;10.5			6.3
(11110)	1	1	10.6			6.4
(00100)	1*	0*	7.2;10.4	10.1	Цепь 7	7.1
(00101)	1	0*	7.3;10.4	10.2		7.2
(01101)	1	1*	7.4;10.5	8.2;10.2		7.3
(11101)	1	1	5.6			7.4
(01000)	0	1*	8.2	10.1	Цепь 8	8.1
(01001)	1*	1	8.3	10.2		8.2
(01011)	1	1	8.4	10.3		8.3
(11011)	1	1	10.6	9.3		8.4
(10000)	0	1*	9.2	10.1	Цепь 9	9.1
(10001)	1*	1	9.3	10.2		9.2
(10011)	1	1	9.4	10.3		9.3
(10111)	1	1	10.6	10.4		9.4
(00000)	0	0	10.2		Цепь 10	10.1
(00001)	0*	0	10.3			10.2
(00011)	1*	0	10.4			10.3
(00111)	1	1*	10.5			10.4
(01111)	1	1	10.6			10.5
(11111)	1	1				10.6
Вопросов	13	12				

Чтобы строить цепи, представленные в табл. (с пятью измерениями, например x_1, x_2, x_3, x_4, x_5 или y_1, y_2, y_3, y_4, y_5), используется последовательный процесс.

Каждый шаг порождения цепи состоит в использовании текущей i -размерной цепи и построения $(i + 1)$ -размерной цепи. Поколение цепей для следующего измерения $(i + 1)$ появляется в результате следующего процесса.

- Мы клонируем i -пространственную цепь, например, имея 1-мерную цепь $(0) < (1)$ мы производим ее копию: $(0) < (1)$.
- После этого мы наращиваем цепь добавляя второе измерение.
- Цепь 1 : $(00) < (01)$.
- Цепь 2 : $(10) < (11)$.
- Затем мы отделяем главный случай (11) от цепи 2 и добавляем его в качестве головы к цепи 1, создавая две 2-мерные цепи:
Новая цепь 1 – $(00) < (01) < (11)$ и
Новая цепь 2 – (10) .
- Затем снова клонируем цепи и добавляем 0 и 1. Получим:
 $(000) < (001) < (011)$ и
 $(100) < (101) < (111)$ и
 $(010) < (110)$
- Отделяем главный случай:
 $(000) < (001) < (011) < (111)$
 $(100) < (101)$
 $(010) < (110)$
- Продолжаем процесс клонирования:
 $(0000) < (0001) < (0011) < (0111)$
 $(1000) < (1001) < (1011) < (1111)$

(0100) < (0101)

(1100) < (1101)

(0010) < (0110)

(1010) < (1110)

- Отделяем главный случай:

(0000) < (0001) < (0011) < (0111) < (1111)

(1000) < (1001) < (1011)

(0100) < (0101) < (1101)

(1100)

(0010) < (0110) < (1110)

(1010)

- Продолжаем процесс клонирования:

(00000) < (00001) < (00011) < (00111) < (01111)

(10000) < (10001) < (10011) < (10111) < (11111)

(01000) < (01001) < (01011)

(11000) < (11001) < (11011)

(00100) < (00101) < (01101)

(10100) < (10101) < (11101)

(1100) < (1100)

(0010) < (0110) < (1110)

(0010) < (0110) < (1110)

(1010) < (1010)

- Отделяем главный случай:

(00000) < (00001) < (00011) < (00111) < (01111) < (11111)

(10000) < (10001) < (10011) < (10111)

(01000) < (01001) < (01011) < (11011)

(11000) < (11001)

(00100) < (00101) < (01101) < (11101)

(10100) < (10101)

(01100) < (11100)

(00010) < (00110) < (01110) < (11110)

(10010) < (10110)

(01010) < (11010)

- Группируем в цепи:

Цепь 1: (01100) < (11100)

Цепь 2: (01010) < (11010)

Цепь 3: (11000) < (11001)

Цепь 4: (10010) < (10110)

Цепь 5: (10100) < (10101)

Цепь 6: (00010) < (00110) < (01110) < (11110)

Цепь 7: (00100) < (00101) < (01101) < (11101)

Цепь 8: (01000) < (01001) < (01011) < (11011)

Цепь 9: (10000) < (10001) < (10011) < (10111)

Цепь 10: (00000) < (00001) < (00011) < (00111) < (01111) < (11111)

Табл. представляет результат этого процесса.

Цепи пронумерованы от 1 до 10, каждый случай имеет свой номер в цепи.

Например, 1.2 означает второй случай в первой цепи.

Знак « * » в столбцах 2, 3 и 4 маркируют ответы, полученные от эксперта.

Например, 1* для случая (01100) в столбце 3 означает, что эксперт ответил «да».

Остающиеся ответы для той же самой цепи в столбце 3 автоматически получены, используя монотонность. Признак $f_1(01100) = 1$ для случая 1.1 расширен для случаев 1.2, 6.3. и 7.3.

Аналогично вычисляются значения третьей монотонной Булевой функции ψ , используя таблицу. Признаки в последовательности (10010) интерпретируются как y_1, y_2, y_3, y_4, y_5 вместо x_1, x_2, x_3, x_4, x_5 которые использовались для f_1 и f_2 . Цепи Hansel те же самые, так как количество признаков то же самое.

В столбцах 5 и 6 выписаны случаи, расширяющие значения функций, без опроса эксперта. Столбец 5 предназначен для расширения значений функции с 1 до 1, столбец 6 для расширения значений с 0 до 0.

Если бы эксперт дал противоположный ответ ($f(01100) = 0$) по сравнению с представленным в табл. для функции f и случая 1.1 (01100), то значения 0 могут быть расширены в столбце 2 для случаев 7.1 (00100) и 8.1 (01000). Эти случаи перечислены в столбце 6 для случая (01100). Тогда нет необходимости спрашивать эксперта о случаях 7.1 (00100) и 8.1 (01000).

Общее количество случаев со знаком «*» для столбцов 3,4 равны соответственно 13 и 12.

Эти количества показывают, что 13 вопросов необходимы для восстановления функции f от x_1, x_2, x_3, x_4, x_5 и 12 вопросов необходимы для восстановления функции ψ от y_1, y_2, y_3, y_4, y_5 . Это только 37.5 % из 32 возможных вопросов.

Полное восстановление функции f с 11 аргументами без оптимизации процесса интервью потребовало бы до $2^{11} = 2048$ вопросов к медику-эксперту.

Иерархия уменьшает максимальное количество вопросов для восстановления монотонных Булевых функций.

Обнаружение диагностических правил на данных

Правила были извлечены с использованием 156 случаев (73 злокачественный, 77 доброкачественный, 2 очень подозрительны и 4 со смешанным диагнозом).

Таблица 1. Примеры извлеченных диагностических правил

Диагностическое правило	F-критерий		Значение F-критерия			Точность диагноза на контроле
			0.01	0.05	0.1	
IF NUMber of calcifications per cm ² is between 10 and 20 AND VOLume > 5 cm ³ THEN Malignant	NUM	0.0029	+	+	+	93.3%
	VOL	0.0040	+	+	+	
IF TOTAl # of calcifications >30 AND VOLume > 5 cm ³ AND DENSITY of calcifications is moderate THEN Malignant	TOT	0.0229	-	+	+	100.0%
	VOL	0.0124	-	+	+	
	DEN	0.0325	-	+	+	
IF VARiation in shape of calcifications is marked AND NUMber of calcifications is between 10 and 20 AND IRregularity in shape of calcifications is moderate THEN Malignant	VAR	0.0044	+	+	+	100.0%
	NUM	0.0039	+	+	+	
	IRR	0.0254	-	+	+	
IF variation in SIZE of calcifications is moderate AND Variation in SHAPE of calcifications is mild AND IRregularity in shape of calcifications is mild THEN Benign	SIZE	0.0150	-	+	+	92.86%
	SHAPE	0.0114	-	+	+	
	IRR	0.0878	-	-	+	

Мы рассмотрели три уровня 0.7, 0.85 и 0.95.

Более высокий уровень условной вероятности уменьшает количество правил и диагностированных пациентов, но увеличивает точность диагноза.

Было обнаружено 44 статистически значительных диагностических правила при 0.05 уровне F-критерия с условной вероятностью, не меньшей, чем 0.75.

Было обнаружено 30 правил с условной вероятностью, не меньшей, чем 0.85.

Было обнаружено 18 правил с условной вероятностью, не меньшей, чем 0.95.

Точность диагноза по правилам с условной вероятностью, не меньшей, чем 0.75 при скользящем контроле – 82 %. Ошибка первого рода была 6.5 % (9 злокачественных случаев были диагностированы как доброкачественные); ошибка второго рода была 11.9 % (16 доброкачественных случаев были диагностированы как злокачественные).

Точность диагноза по правилам с условной вероятностью, не меньшей, чем 0.85 дали точность 90 %,

Точность диагноза по правилам с условной вероятностью, не меньшей, чем 0.95 дали точностью 96.6 %, только с тремя ошибками второго рода (3.4 %).

Извлечение правил из монотонных Булевых функций

Мы получили Булево выражение для формы и плотности кальциноза $x_2 = \psi(y_1, y_2, y_3, y_4, y_5)$ из информации в столбцах 1 и 4, следуя следующим шагам:

- i) найти все максимальные нижние единицы для цепей в виде элементарных конъюнкций;
- ii) исключить избыточные термины (конъюнкции) из окончательной формулы.

Таким образом, из столбца 4 мы получим

$$x_1 = v(w_1, w_2, w_3) = w_2 \vee w_1 w_3;$$

$$x_2 = \psi(y_1, y_2, y_3, y_4, y_5) = y_2 y_3 \vee y_2 y_4 \vee y_1 y_2 \vee y_1 y_4 \vee y_1 y_3 \vee y_2 y_3 y_5 \vee y_2 \vee y_1 \vee y_3 y_4 y_5 = \\ = y_2 \vee y_1 \vee y_3 y_4 y_5.$$

Из столбца 3 мы получим компоненты функции от переменных x_1, x_2, x_3, x_4, x_5 следующим образом:

$$f(x) = x_2 x_3 \vee x_1 x_2 x_4 \vee x_1 x_2 \vee x_1 x_3 x_4 \vee x_1 x_3 \vee x_3 x_4 \vee x_3 \vee x_2 x_5 \vee x_1 x_5 \vee x_4 x_5 =$$

$$x_1 x_2 \vee x_3 \vee (x_2 \vee x_1 \vee x_4) x_5 = (w_2 \vee w_1 w_3)(y_1 \vee y_2 \vee y_3 y_4 y_5) \vee x_3 \vee (y_1 \vee y_2 \vee y_3 y_4 y_5) \vee (w_2 \vee w_1 w_3 \vee x_4) x_5.$$

Сравнение экспертных и извлеченных из данных правил. Комментарии радиолога:

ЕСЛИ общее количество кальцинозов > 30

И объем $> 5 \text{ см}^3$

И плотность кальцинозов умеренна,

ТО злокачественная.

F-критерий значим при уровне 0.05. Точность диагноза на контроле – 100 %.

Комментарий радиолога – это правило обещающее, но я считаю это рискованным.

ЕСЛИ изменение в форме кальцинозов отмечено

И количество кальцинозов между 10 и 20

И неисправность в форме кальцинозов умеренна,

ТО – злокачественная.

F-критерий значим при уровне 0.05. Точность диагноза на контроле – 100 %.

Комментарий радиолога – я доверял бы этому правилу.

ЕСЛИ изменение в размере кальцинозов умеренно

И изменение в форме кальцинозов умеренно

И неисправность в форме кальцинозов умеренна,

ТО – доброкачественная.

F-критерий значим при уровне 0.05. Точность диагноза на контроле – 92.86%.

Комментарий радиолога – я доверял бы этому правилу.