

Лекция 4. Излечение теории из данных.

Пусть исследуемая реальность – предметная область представлена эмпирической системой $\mathfrak{I} = \langle A; \Omega_{\mathfrak{I}} \rangle$ сигнатуры Ω ,

A – множество объектов;

$$\Omega_{\mathfrak{I}} = \{P^{\mathfrak{I}}_1, \dots, P^{\mathfrak{I}}_n\}.$$

Будем предполагать, что теория $Th(\mathfrak{I})$ эмпирической системы \mathfrak{I} (совокупность всех истинных на \mathfrak{I} высказываний) представляет собой совокупность универсальных формул.

Будем предполагать, что существует N , такое что любая аксиома из $Th(\mathfrak{I})$ имеет не более чем N кванторов всеобщности.

Определим фрагмент языка первого порядка L сигнатуры Ω , включающий:

$X = \{x_1, x_2, \dots\}$ – множество свободных переменных;

$U(\Omega)$ – множество всех атомарных формул (атомов) вида $P(x_1, \dots, x_n)$, $x_1, \dots, x_n \in X$;

$\mathfrak{R}(\Omega)$ – множество утверждений языка L , полученное замыканием множества

$U(\Omega)$ относительно логических операций $\&, \vee, \neg$.

В рамках исчисления высказываний на элементах булевой алгебры $\mathfrak{R}(\Omega)$ определено тождество утверждений $A \equiv B$. Будем предполагать, что логические константы $I \equiv A \vee \neg A$ и $L \equiv A \& \neg A$ принадлежат $\mathfrak{R}(\Omega)$.

Известно, что совокупность универсальных формул логически эквивалентна совокупности правил вида

$$\forall x_1, \dots, x_k (A_1^{\varepsilon_1} \& \dots \& A_k^{\varepsilon_k} \Rightarrow A_0^{\varepsilon_0}), \quad k \geq 0 \quad (1)$$

где A_0, A_1, \dots, A_k – атомарные формулы, $A_j = P_j(x_1^j, \dots, x_{n_j}^j)^{\varepsilon_j}$, $j = 0, 1, \dots, k$;

$\varepsilon_0, \varepsilon_1, \dots, \varepsilon_k = 1(0)$, если атомарная формула берется без отрицания (1) или с отрицанием (0).

Задача: Обнаружить теорию $Th(\mathfrak{I})$ эмпирической системы \mathfrak{I} и, в частности, обнаружить систему аксиом эмпирической системы \mathfrak{I} .

Проанализируем эту задачу.

Что можно сказать об истинности высказываний из $Th(\mathfrak{I})$ на эмпирической системе \mathfrak{I} , опираясь только на логический анализ высказываний.

Можно сказать, во-первых, что правило $C = (A_1 \& \dots \& A_k \Rightarrow A_0)$ может быть истинным на эмпирической системе только потому, что посылка правила всегда ложна. На самом деле, как мы покажем, это означает, что на эмпирической системе истинно некоторое логически более сильное "подправило", связывающее между собой атомы посылки.

Во-вторых, правило C может быть истинно на эмпирической системе только потому, что некоторое его логически более сильное "подправило", содержащее только часть посылки и то же заключение, истинно на эмпирической системе.

Поэтому система аксиом может быть истинной на эмпирической системе потому, что истинна некоторая система подправил, из истинности которой в свою очередь следует истинность системы аксиом.

Выясним из истинности каких логически более сильных "подправил" на эмпирической системе \mathfrak{I} следует истинность самого правила. Тем самым мы получим определение "подправил" и определение закона для эмпирической системы \mathfrak{I} .

Теорема 1. Правило $C = (A_1 \& \dots \& A_k \Rightarrow A_0)$ логически следует (в исчислении высказываний) из любого правила вида:

1. $A_{i1} \& \dots \& A_{ih} \Rightarrow \neg A_{i0}$,

где $\{A_{i1}, \dots, A_{ih}, A_{i0}\} \subset \{A_1, \dots, A_k\}$, $0 \leq h < k$, т.е.

$(A_{i1} \& \dots \& A_{ih} \Rightarrow \neg A_{i0}) \vdash \neg(A_1 \& \dots \& A_k) \vdash (A_1 \& \dots \& A_k \Rightarrow A_0)$;

\vdash - доказуемость в исчислении высказываний.

2. $(A_{i1} \& \dots \& A_{ih} \Rightarrow A_0)$,

где $\{A_{i1}, \dots, A_{ih}\} \subset \{A_1, \dots, A_k\}$, $0 \leq h < k$, т.е.,

$(A_{i1} \& \dots \& A_{ih} \Rightarrow A_0) \vdash (A_1 \& \dots \& A_k \Rightarrow A_0)$.

Доказательство. 1. Докажем сначала первую цепочку выводов $(A_{i1} \& \dots \& A_{ih} \Rightarrow \neg A_{i0}) \equiv (\neg(A_{i1} \& \dots \& A_{ih}) \vee \neg A_{i0}) \equiv (\neg A_{i1} \vee \dots \vee \neg A_{ih} \vee \neg A_{i0}) \equiv \neg(A_{i1} \& \dots \& A_{ih} \& A_{i0})$. Так как $\{A_{i1}, \dots, A_{ih}, A_{i0}\} \subset \{A_1, \dots, A_k\}$, то конъюнкция $A_{i1} \& \dots \& A_{ih} \& A_{i0}$ является частью конъюнкции $A_1 \& \dots \& A_k$. Из аксиомы алгебры высказываний $A \& B \vdash A$ по правилу modus ponense следует, что $A_1 \& \dots \& A_k \vdash A_{i1} \& \dots \& A_{ih} \& A_{i0}$. Поэтому из формул $(A_{i1} \& \dots \& A_{ih} \Rightarrow \neg A_{i0})$, $A_1 \& \dots \& A_k$ выводится противоречие $\neg(A_{i1} \& \dots \& A_{ih} \& A_{i0}) \& (A_{i1} \& \dots \& A_{ih} \& A_{i0})$. Отсюда следует, что $(A_{i1} \& \dots \& A_{ih} \Rightarrow \neg A_{i0}) \vdash \neg(A_1 \& \dots \& A_k)$. Докажем, что $\neg(A_1 \& \dots \& A_k) \vdash (A_1 \& \dots \& A_k \Rightarrow A_0)$. Так как $\neg(A_1 \& \dots \& A_k) \equiv \neg A_1 \vee \dots \vee \neg A_k$, то по правилу алгебры высказываний $A \vdash A \vee B$ выводим, что $\neg(A_1 \& \dots \& A_k) \vdash (\neg A_1 \vee \dots \vee \neg A_k \vee A_0) \equiv (A_1 \& \dots \& A_k \Rightarrow A_0)$.

2. Докажем, что $(A_{i1} \& \dots \& A_{ih} \Rightarrow A_0) \vdash (A_1 \& \dots \& A_k \Rightarrow A_0)$, если $\{A_{i1}, \dots, A_{ih}\} \subset \{A_1, \dots, A_k\}$, $0 \leq h < k$. Так как $(A_{i1} \& \dots \& A_{ih} \Rightarrow A_0) \equiv (\neg A_{i1} \vee \dots \vee \neg A_{ih} \vee A_0)$, то из схемы аксиом $A \vdash A \vee B$ алгебры высказываний и правила modus ponens следует, что $(\neg A_{i1} \vee \dots \vee \neg A_{ih} \vee A_0) \vdash (\neg A_1 \vee \dots \vee \neg A_k \vee A_0) \equiv (A_1 \& \dots \& A_k \Rightarrow A_0)$ ◆

Из логики и методологии науки хорошо известно, что те высказывания следует считать **законами**, которые при одинаковой их подтверждённости на экспериментальных данных наиболее фальсифицируемы, просты и/или содержат наименьшее число "параметров".

В нашем случае все эти свойства, которые обычно трудно определить, следуют из определения логической силы высказывания.

"Подправило" одновременно является:

- логически более сильным высказыванием, чем само правило и более фальсифицируемым, так как содержит более слабую посылку и, следовательно, применимо к большему объему данных и тем самым в большей степени подвержено фальсификации;
- более простым, так как содержит меньшее число атомарных высказываний, чем правило;
- включает меньшее число "параметров", так как лишние атомарные высказывания тоже можно считать параметрами "подстройки" высказывания под данные.

Почему же закон должен быть наиболее фальсифицируемым, простым и содержать наименьшее число параметров? Разные авторы придерживаются различных мнений на этот счет, либо не объясняют этого вообще.

В нашем случае для гипотез вида (1) мы можем более точно ответить на этот вопрос. Так как все описанные свойства закона вытекают из логической силы высказывания, то поиск логически наиболее сильных "подправил", истинных на эмпирической системе, позволяет нам не только проверить гипотезу об истинности системы аксиом, но и решить другую принципиально более важную задачу:

выяснить, а какова на самом деле та наиболее сильная (логически) теория, вытекающая из этих правил, которая описывает наши данные и возможно лежит в основании неизвестного нам закона их порождения?

Решение этой задачи обнаружения закона в данных или, что то же самое, поиска сильнейшей теории в данных как раз и требует нахождения среди всех правил вида (1) логически наиболее сильных (среди истинных на эмпирической системе). Именно такие правила в соответствии с существующими представлениями следует считать законами эмпирической системы.

Следствие 1. Если некоторое подправило правила С истинно на эмпирической системе \mathfrak{I} , то и само правило С истинно на \mathfrak{I} .

Определение 1. Подправилом некоторого правила С будем называть любое из логически более сильных правил вида 1 или 2, определенных в теореме 1 для правила С.

Как легко видеть любое подправило также имеет вид (1).

Определение 2. Законом эмпирической системы \mathfrak{I} будем называть любое, истинное на \mathfrak{I} правило С вида (1), для которого каждое его подправило уже не истинно на \mathfrak{I} .

Обозначим через L множество всех законов эмпирической системы \mathfrak{I} .

Теорема 2. $L \vdash \text{Th}(\mathfrak{I})$.

Доказательство: Пусть есть некоторое высказывание В истинное на эмпирической системе \mathfrak{I} . По предположению оно универсально аксиоматизируемо. Это правило может быть преобразовано в совокупность правил (1) истинных на \mathfrak{I} . Найдем для каждого правила найдем минимальное подправило, для которого уже нет подправил истинных на \mathfrak{I} . Тогда эти подправила будут законами на \mathfrak{I} в силу определения закона ◆

Таким образом, обнаружение множества всех законов L решает задачу обнаружения теории эмпирической системы $\text{Th}(\mathfrak{I})$.

Задача извлечения знаний из данных. Примеры правил.

1. Монотонность: $\forall a, b (a <_x b \Rightarrow a <_y b)$;
2. Образы разделенные монотонными правилами:
 $\forall a, b (P(a) \& (a <_x b) \& (a >_y b) \Rightarrow P(b))$;
3. $\forall a, b (\text{weekday}(a) = \text{weekday}(b) = \langle d_1, \dots, d_5 \rangle) \& (a \#) (b))^{gl} \& \dots \& (a \#) (b))^{gk} \Rightarrow (\text{target}(a^5) \# \text{target}(b^5))^{g0}$;
4. $\forall a (\text{weekday_16_days_back}(a) = d_5 \& \text{weekday_13_days_back}(a) = d_3 \& \text{price}(\text{weekday_16_days_back}) < \text{price}(\text{weekday_13_days_back}) \& \text{weekday_5_days_back}(a) = d_1 \& \text{price}(\text{weekday_13_days_back}) > \text{price}(\text{weekday_5_days_back}) \Rightarrow \text{target_5_days_ahead}(a) \# \text{target_5_days_ahead}(b^5))$;
5. $(a < b) \& (b < c) \Rightarrow (a < c)$ транзитивность в психофизических экспериментах.