

**О ёмкости семейств характеристических функций,
обеспечивающих корректное решение
задач диагностического типа**

М.И.Забейайло

ФИЦ ИУ РАН

Интерполяционно-экстраполяционная схема

- ИЭ-схема для формализации задач машинного обучения (неклассическая *интерполяция* + *экстраполяция* по Ю.И.Журавлеву) :

Даны:

- множество Ω^+ (примеров) объектов, обладающих целевыми свойствами P_i^+ из некоторого множества $P = \{ P_{i_1}, P_{i_2}, \dots, P_{i_r} \}$, и
- множество Ω^- (контрпримеров) объектов, не обладающих целевыми свойствами из $P = \{ P_{i_1}, P_{i_2}, \dots, P_{i_r} \}$,
- *новый объект* O^t (или же некоторое явным образом представленное множество таких объектов Ω^t).

Требуется:

оценить *наличие* (или отсутствие) *целевых свойств* у нового объекта O^t (новых объектов из заданного множества Ω^t), т.е.

- дать соответствующий **прогноз** (о наличии целевых свойств) и
- предъявить **основания** (неоспоримые **аргументы**), позволяющие **принять** этот прогноз

- Пример: задачи **диагностического** типа

Эффект переобучения как критичная характеристика *ИЭ*-схем

(*К.В.Воронцов, Д.В.Виноградов*)

- **Воронцов К.В.** Комбинаторная теория надёжности обучения по прецедентам. – Диссертация на соискание учёной степени д.ф-м.н. по специальности 05.13.17 Теоретические основы информатики. – М: ВЦ РАН. - 273 С. – URL: <https://www.dissercat.com/content/kombinatornaya-teoriya-nadezhnosti-obucheniya-po-pretseidentam>
- **Виноградов Д.В.** Вероятностно-комбинаторный формальный метод обучения, основанный на теории решеток. – Диссертация на соискание учёной степени д.ф-м.н. по специальности 05.13.17 Теоретические основы информатики. – М.:ФИЦ ИУ РАН, 2018. – 131 с. – URL: http://www.frccsc.ru/diss-council/00207305/diss/list/vinogradov_dv

Можно ли бороться с переобучением?

Как бороться с переобучением? В чем его причина?

Возможные версии ответа:

- **Каузальность** vs **интерполируемость**
- Таблица прецедентов и «объясняющие» ее частично определенные функции для задач диагностического типа
- **Проблема устойчивости** интерполирующих ЭЗ при расширении БФ описаниями новых прецедентов
 - «универсальные» ЭЗ
 - контекстно-зависимые ЭЗ

Опыт поиска «универсальных» ЭЗ

Опыт поиска «универсальных» эмпирических зависимостей ЭЗ в ИИ-исследованиях 1990-2000-х годов (D.Lenat, Y.Zytkow, R.Michalski, ...)

- **Jan M.Zytkow** University of North Carolina, Dept. of Computer Science. Charlotte, NC 28223, USA
Principles of data mining and knowledge discovery : Proc. 3-d European conference “PKDD ’99”, Prague, Czech Republic, September 15 - 18, 1999. *Jan M. Zytkow ; Jan Rauch* (ed.). - *Lect. notes in comp. science.* - Vol. 1704 : *Lecture notes in artificial intelligence.* ISBN 3-540-66490-4. -Berlin et al.: Springer, 1999. -
<https://link.springer.com/content/pdf/bfm%3A978-3-540-48247-5%2F1.pdf>
- **D.Lenat** - https://en.wikipedia.org/wiki/Douglas_Lenat
Lenat D.B., Fishwick P.A., Modjeski R.B., Oresky C.M., Clarkson A., Kaisler S. (1991). "**STRADS: A Strategic Automatic Discovery System**". *Knowledge-based Simulation: Methodology and Application.* – Pp. 133-161. -
<https://www.springer.com/gp/book/9780387973746> (См. Главу 11: <https://www.springer.com/gp/book/9780387973746>)
- **R.Michalski** - https://en.wikipedia.org/wiki/Ryszard_S._Michalski
Michalski, R.S., Carbonell, J.G., Mitchell, T.M.: '**Machine learning: an artificial intelligence approach**' (Springer Science & Business Media, Berlin, Germany, 2013). - <https://www.springer.com/gp/book/9783662124079>

Проблема устойчивости ЭЗ при расширении анализируемой Базы Фактов

- ***Big Data***: эффекты ***Big*** и ***Open***
- Последовательности *расширяющихся* Баз Фактов
- Проблема *устойчивости* ЭЗ, «интерполирующих» данную БФ, в условиях характерного для *Big Data* эффекта *Open* – *открытости\пополняемости* текущей БФ новыми эмпирическими фактами (**В.К.Финн**)

Особые свойства задач диагностического типа

- Необходимость оперировать **выборками ограниченного размера** (по количеству прецедентов), что ставит под вопрос **адекватность** применения целого ряда методов **статистического анализа данных**
- Возможное отсутствие оснований рассматривать наблюдаемые (используемые для описания прецедентов) параметры как **независимые**
- Возникновение эффектов **Big Data** на уровне размеров множеств значений фиксируемых параметров
- Проблемы с **однородностью** анализируемых прецедентов (и их описаний)
- **Интерпретируемость** в содержательных терминах и понятиях исследуемой предметной области (а не сформированной для ее описания алгоритмической модели) **результатов** (заключений, диагнозов, ...), формируемых использованием тех или иных математических инструментов анализа данных
- Неформальная **объясняемость** результатов (заключений, диагнозов, ...)
- **Персонализация** заключений (диагнозов): возможности учесть *индивидуальные особенности* пациента
- Проблема **доверия - достаточности оснований** для **принятия** порождаемых заключений (диагнозов, врачебных *рекомендаций* и т.п.)
- Наличие в рассматриваемой предметной области **доказуемо трудно-разрешимых** переборных задач, накладывающее существенные ограничения на **практическую применимость** ряда математических методов анализа данных и поддержки принятия диагностических решений

Характеристические Функции на последовательностях расширяющихся БФ в задачах диагностического типа

- Характеристические Функции как (*каузально-ориентированный*) подкласс семейства частичных функций, интерполирующих выборки прецедентов – примеров и контрпримеров диагностируемого явления
- Характеристическая Функция (ХФ), принимает
 - значение «*истина*» на всех фактах ϕ (*примерах*) текущей базы фактов **БФ**, характеризующих наличие анализируемого целевого свойства и
 - значение «*ложь*» на всех фактах ϕ (*контрпримерах*) текущей базы фактов **БФ**, характеризующих наличие анализируемого целевого свойства

$$\forall \phi [(\phi \in \text{БФ}) \supset (\text{ХФ}(\text{БФ}) \mid - \phi)]$$

Процедура порождения Характеристических Функций

Прецеденты (формализованное описание)

=>

Сходство прецедентов (формализация)

=>

Классы сходства прецедентов

=>

Классы эквивалентности прецедентов

=>

Запрет на контрпримеры (условие **ЗКП**)

=>

Покрытия текущей БФ **ЗКП**-классами эквивалентности примеров

- Пусть $\Omega = \{ O_1, O_2, \dots, O_m \}$. Построим по Ω и \otimes множество $Dom(\Omega)$:

(i) $\Omega = \{ O_1, O_2, \dots, O_m \} \subset Dom(\Omega)$.

(ii) $\{ [A \in Dom(\Omega)] \ \& \ [B \in Dom(\Omega)] \ \& \ [(A \otimes B) \neq \emptyset] \} \rightarrow [(A \otimes B) \in Dom(\Omega)]$,

(iii) Других элементов в $Dom(\Omega)$ нет.

- Отношение **сходства** \mathbf{R}^\otimes :

$$O_{i1} \mathbf{R}^\otimes O_{i2} \text{ (т.е. } \langle O_{i1}, O_{i2} \rangle \in \mathbf{R}^\otimes) \quad \text{iff} \quad O_{i1} \otimes O_{i2} \neq \emptyset$$

- **Классы сходства** (толерантности) на Ω :

$$\mathbf{T}(O_{i1}) = \{ O_{i2} \mid O_{i1} \mathbf{R}^\otimes O_{i2} \}$$

- Фиксируя каждый конкретный результат $V=V_0$ вычисления операции \otimes сходства на элементах множества $Dom(\Omega)$:

$$A \otimes B = V_0$$

выделим соответствующие подклассы (**эквивалентности**) сформированных классов сходств:

$$\mathbf{E}_{V_0} = \{ \langle O_{i1}, O_{i2} \rangle \mid O_{i1} \otimes O_{i2} \otimes V_0 = V_0 \}$$

- **Запрет на контрпримеры** (условие **ЗКП**):

$$(\forall V_0) [(\mathbf{E}_{V_0}^+ = \{ \langle O_{i1}^+, O_{i2}^+ \rangle \mid (O_{i1}^+ \otimes O_{i2}^+ \otimes V_0 = V_0) \}) \ \& \ (O_{i1}^+ \in \Omega^+) \ \& \ (O_{i2}^+ \in \Omega^+)] \supset$$

$$\supset \neg (\exists O_{i0}^-) \{ (V_0 \otimes O_{i0}^- = V_0) \ \& \ (O_{i0}^- \in \Omega^-) \}]$$

Процедура формирования Характеристических Функций

Определение

Будем называть **характеристической функцией** $X\Phi_i(\mathbf{БФ})$ базы фактов $\mathbf{БФ}$ логическое условие, описывающее конкретное покрытие $Cov_i(\mathbf{БФ})$ по следующей схеме:

- для каждого класса эквивалентности прецедентов $GC\text{-}3КП_{ij}(\mathbf{БФ})$, входящего в покрытие $Cov_i(\mathbf{БФ})$ анализируемого множества прецедентов $\mathbf{БФ} = \Omega = \Omega^+ \cup \Omega^-$, (при этом, естественно, $\Omega^+ \cap \Omega^- = \emptyset$), формирующее его множество признаков $\{a^{ij}_1, a^{ij}_2, \dots, a^{ij}_{n(ij)}\}$ - образующих из исходного множества $U = \{a_1, a_2, \dots, a_n\}$ - порождает в сопоставляемой покрытию $Cov_i(\mathbf{БФ})$ характеристической функции $X\Phi_i(\mathbf{БФ})$ конъюнкцию переменных –

$$Q_{ij} = p^{ij}_1 \& p^{ij}_2 \& \dots \& p^{ij}_{n(ij)}, \quad a$$

- сама характеристическая функция $X\Phi_i(\mathbf{БФ})$ представляет собою дизъюнкцию по всем (сопоставляемым входящим в покрытие $Cov_i(\mathbf{БФ})$ классам эквивалентности $GC\text{-}3КП_{ij}(\mathbf{БФ})$) конъюнкциям Q_{ij} :

$$X\Phi_i(\mathbf{БФ}) = \bigvee_j (Q_{ij}).$$

по всем j

Характеристические Функции и каузальная репрезентативность анализируемой БФ

Утверждение 1 (о свойствах ХФ как характеристической функции, разделяющей примеры и контрпримеры с учетом причин возникновения анализируемого целевого эффекта)

Каждая построенная в соответствии с **Определением 1** функция $X\Phi_i(\mathbf{БФ})$ принимает на факте ϕ из текущей **БФ** значение

- «*истина*» тогда и только тогда, когда данный факт характеризуется наличием анализируемого целевого свойства;
- «*ложь*» тогда и только тогда, когда данный факт характеризуется отсутствием анализируемого целевого свойства,

и наоборот:

- на каждом факте ϕ из текущей **БФ**, характеризуемом наличием анализируемого целевого свойства, характеристическая функция $X\Phi_i(\mathbf{БФ})$ принимает значение «*истина*», а значение «*ложь*» тогда и только тогда, когда данный факт характеризуется отсутствием анализируемого целевого свойства



- **Непустота множества $X\Phi(\mathbf{БФ})$** как характеристика *каузальной репрезентативности* анализируемой БФ

Задача о емкости семейства всех ХФ, порождаемых на данной БФ

Задача о емкости семейства $\mathbf{ХФ(БФ)}$ всех ХФ, порождаемых на данной $\mathbf{БФ}$:

Дана: база фактов (множество прецедентов) $\mathbf{БФ} = \Omega = \Omega^+ \cup \Omega^-$,
объединяющая примеры Ω^+ и контрпримеры Ω^-

Требуется: определить, сколько элементов содержит множество $\mathbf{ХФ(БФ)}$
всех ХФ, порождаемых на данной $\mathbf{БФ}$

Утверждение 2

Задача о числе характеристических функций, формируемых на произвольной базе фактов $\mathbf{БФ}$, принадлежит классу $\mathbf{\#PC}$ перечислительно полных переборных задач



Задача о емкости семейства минимальных (по вложимости покрытий) ХФ

Задача о числе ХФ, *минимальных* по вложимости соответствующих им покрытий БФ ЗКП-классами эквивалентности примеров из БФ:

Дана: база фактов (множество прецедентов) $\mathbf{БФ} = \Omega = \Omega^+ \cup \Omega^-$,
объединяющая примеры Ω^+ и контрпримеры Ω^-

Требуется: определить, сколько элементов содержит множество $\mathbf{ХФ}_{min}(\mathbf{БФ})$
всех ХФ, порождаемых минимальными по вложению покрытиями
данной $\mathbf{БФ}$ ЗКП-классами эквивалентности примеров из Ω^+

Утверждение 3

Задача о числе характеристических функций, формируемых на произвольной базе фактов *минимальными* по вложению покрытиями данной $\mathbf{БФ}$ ЗКП-классами эквивалентности примеров принадлежит классу $\#P$ перечислительно полных переборных задач



«Быстрая» разрешимость задач о ХФ?

Проблема быстроты диагностирования:

- **каузальной репрезентативности** анализируемой **БФ**
- **непустоты подсемейства ХФ, наследуемых** (*сохраняющих способность корректно представлять и БФ и ее конкретное расширение БФ + ΔБФ*) при добавлении к текущей **БФ** новых прецедентов **ΔБФ**

Заключение

- **Каузальность** \neq **интерполируемость**
(особый статус задач диагностического типа)
- Комбинаторная «плата» за **точность интерполяции** в задачах **диагностического типа**

Спасибо за внимание



m.zabezhailo@yandex.ru