

О вычислительной сложности задач отбора объектов и признаков при построении монотонных классификаторов

Зухба А. В.



Светлогорск • 19–25 сентября 2015

Основные определения, обозначения, понятия

Дано множество **объектов** $\mathbb{X} = \{x_1, \dots, x_\ell\}$, называемых **обучающей выборкой**, и $\mathbb{Y} = \{0, 1\}$ — множество **классов**.

Объекты описываются **признаками** $\mathbb{F} = \{f_1, \dots, f_t\}$.

Каждый признак задает отображение $f_j: \mathbb{X} \rightarrow E_j$,
где E_j — линейно упорядоченное множество.

Любое непустое подмножество множества признаков $F \subseteq \mathbb{F}$
индуцирует отношение **частичного порядка** на \mathbb{X} :
 $x \leq x'$ тогда и только тогда, когда $f(x) \leq f(x')$ для всех $f \in F$.

Основные определения, обозначения, понятия

Множество объектов \mathbb{X} разбито на подмножества: $\mathbb{X} = \mathbb{A} \cup \mathbb{B}$
 \mathbb{A} — объекты класса **1** и \mathbb{B} — объекты класса **0**.

Пара объектов $(a, b) \in \mathbb{A} \times \mathbb{B}$ называется **монотонной**, если $a > b$.
Множество всех монотонных пар обозначается через M .

Пара объектов $(a, b) \in \mathbb{A} \times \mathbb{B}$ называется **дефектной**, если $a < b$.
Множество всех дефектных пар обозначается через D .

Множество пар, монотонных по признаку f , будем обозначать M_f ,
монотонных по совокупности признаков F — через M_F .

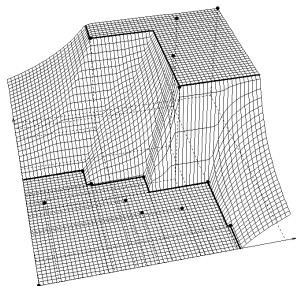
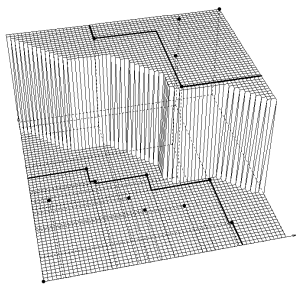
Множество пар, дефектных по признаку f , будем обозначать D_f ,
дефектных по совокупности признаков F — через D_F .

Задача построения монотонного классификатора

Функция $y: \mathbb{X} \rightarrow \mathbb{Y}$ монотонна, если для любых двух объектов $x, x' \in \mathbb{X}$ из $x < x'$ следует $y(x) \leq y(x')$.

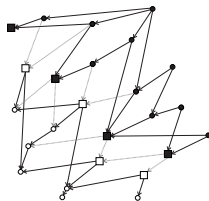
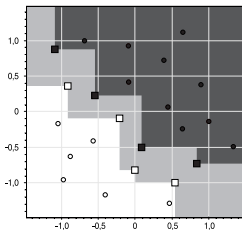
Задача

Приблизить неизвестную функцию $y^*: \mathbb{X} \rightarrow \mathbb{Y}$, заданную на объектах \mathbb{X} , монотонной функцией y .



Задача построения монотонного классификатора

Построение монотонного классификатора по монотонной выборке:



Представим выборку в виде ориентированного графа частичного порядка с раскрашенными в два цвета вершинами. Цвета соответствуют классам, ребра направлены от большего объекта к меньшему.

Утверждение

Построение множества эталонных объектов можно совершить за $O(|X| + |E|)$, где E — количество ребер.

Монотонные алгоритмы классификации

- линейные модели с неотрицательными коэффициентами
- монотонный метод ближайшего соседа¹
- монотонные решающие деревья²
- монотонные нейросети³

1. Воронцов К. В., Махина Г. А. Принцип максимизации зазора для монотонного классификатора ближайшего соседа. В: *15-я всероссийская конференция «Математические методы распознавания образов»*. М.:МАКС Пресс, 2011. С.

2. Kamp R., Feelders A., Barile N. Isotonic classification trees. In: *Proceedings of the 8th International Symposium on Intelligent Data Analysis: Advances in Intelligent Data Analysis VIII*. Berlin, Heidelberg: Springer-Verlag, 2009. P. 405–416.

3. Sill J. Monotonic networks. In: *Advances in Neural Information Processing Systems*. Ed. Jordan M. I., Kearns M. J., Solla S. A. Cambridge: MIT Press, 1998. P. 661–667.

Задача монотонизации выборки

Замечание

Среди объектов выборки \mathbb{X} могут присутствовать дефектные пары: $x < x'$, такие, что $y^(x) > y^*(x')$, то есть $x \in \mathbb{A}$, $x' \in \mathbb{B}$.*

Задача

Отобрать подмножества объектов $X \subseteq \mathbb{X}$ и признаков $F \subseteq \mathbb{F}$ так, чтобы монотонных пар M_F было как можно больше, а дефектных пар D_F — как можно меньше.

Функционалы качества

- Степень монотонности⁴ (degree of monotonicity)

$$DgrMon = \frac{|M|}{|M| + |D|}$$

- Эмпирический риск⁵

$$\sum_{x \in A} [y(x) = 0] + \sum_{x \in B} [y(x) = 1]$$

То есть для минимизации эмпирического риска необходимо минимизировать количество объектов, на которых происходит ошибка. Переформулируем в виде задачи **отбора объектов**:

$$|D| = 0, |X| \rightarrow \max$$

считая, что алгоритм отбрасывает объекты, на которых ошибается.

4. Marina Velikova, Hennie Daniels. On Testing Monotonicity of Datasets. In: *Learning monotone models*. 2009. P. 11–22.

5. Гуз И. С. Минимизация эмпирического риска при построении монотонных композиций классификаторов. *Труды МФТИ*. 2011. Т. 3, № 3(11). С.115–121.

Параметры оптимизации

Утверждение

Для произвольного подмножества признаков $F \subseteq \mathbb{F}$

$$M_F = \bigcap_{f \in F} M_f, \quad D_F = \bigcap_{f \in F} D_f.$$

Следствие

Для любых подмножеств признаков $F, G \subseteq \mathbb{F}$

$$F \subseteq G \Rightarrow |M_G| \leq |M_F|, \quad |D_G| \leq |D_F|.$$

$|D_F| \rightarrow \min$ можно заменить $|F| \rightarrow \max$ или $|X| \rightarrow \min$,
 $|M_F| \rightarrow \max$ можно заменить $|F| \rightarrow \min$ или $|X| \rightarrow \max$.

Систематизация задач монотонизации

Отбор признаков:

$|M| \rightarrow \max$, $|D| \rightarrow \min$,
 $|F| \rightarrow \max$ (максимум информации),
 $|F| \rightarrow \min$ (самая простая модель).

- $FS(|M| \geq m, |D| \leq d)$
- $FS(|M| \geq m, |F| \leq q)$
- $FS(|M| \geq m, |F| \geq q)$
- $FS(|F| \leq q, |D| \leq d)$
- $FS(|F| \geq q, |D| \leq d)$

Отбор объектов:

$|M| \rightarrow \max$, $|D| \rightarrow \min$,
 $|X| \rightarrow \max$ (фильтрация выбросов),
 $|X| \rightarrow \min$ (отбор эталонов).

- $PS(|M| \geq m, |D| \leq d)$
- $PS(|M| \geq m, |X| \leq n)$
- $PS(|M| \geq m, |X| \geq n)$
- $PS(|X| \leq n, |D| \leq d)$
- $PS(|X| \geq n, |D| \leq d)$

Всего 10 задач.

Систематизация задач монотонизации

Отбор признаков:

$|M| \rightarrow \max$, $|D| \rightarrow \min$,
 $|F| \rightarrow \max$ (максимум информации),
 $|F| \rightarrow \min$ (самая простая модель).

- $FS(|M| \geq m, |D| \leq d)$
- ~~$FS(|M| \geq m, |F| \leq q)$~~
- $FS(|M| \geq m, |F| \geq q)$
- $FS(|F| \leq q, |D| \leq d)$
- ~~$FS(|F| \geq q, |D| \leq d)$~~

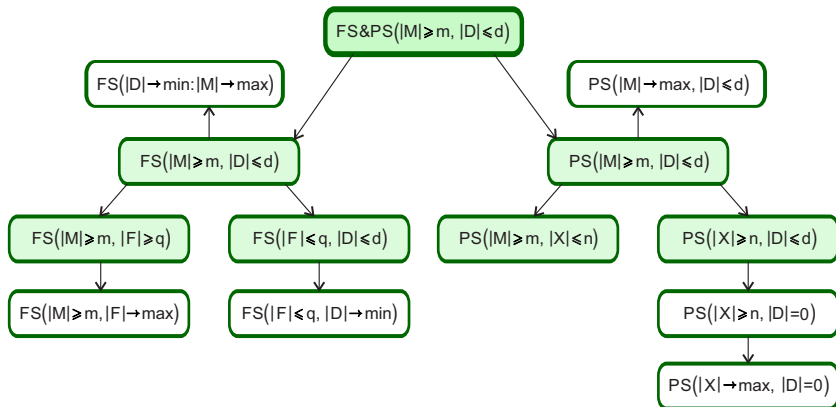
Отбор объектов:

$|M| \rightarrow \max$, $|D| \rightarrow \min$,
 $|X| \rightarrow \max$ (фильтрация выбросов),
 $|X| \rightarrow \min$ (отбор эталонов).

- $PS(|M| \geq m, |D| \leq d)$
- $PS(|M| \geq m, |X| \leq n)$
- ~~$PS(|M| \geq m, |X| \geq n)$~~
- ~~$PS(|X| \leq n, |D| \leq d)$~~
- $PS(|X| \geq n, |D| \leq d)$

Всего 10 задач. Из них 4 тривиальных и 6 содержательных.

Систематизация задач отбора объектов и признаков



Вычислительная сложность

Утверждение

Все предложенные постановки задач монотонизации кроме $PS(|X| \geq n, |D| \leq d)$, $PS(|X| \geq n, |D|=0)$, $PS(|X| \rightarrow \max, |D|=0)$ являются NP-трудными.

Для доказательства NP-трудности была построена полиномиальная сводимость к задачам монотонизации таких известных NP-полных задач, как:

- *задача о рюкзаке,*
- *задача о биклике,*
- *задача о минимальном покрытии множества подмножествами.*

Вычислительная сложность

Утверждение

Решение задачи в постановках $PS(|X| \rightarrow \max, |D|=0)$ и $PS(|X| \geq n, |D|=0)$ сводится к решению задачи поиска минимального вершинного покрытия для двудольного графа.

Построен алгоритм, решающий задачу в постановках $PS(|X| \rightarrow \max, |D|=0)$ и $PS(|X| \geq n, |D|=0)$ за время $O(|D_0||X_0|^{0,5})$, где $|X_0|$ — количество объектов, задействованных в дефектных парах, а $|D_0|$ — количество дефектных пар.

Для постановки $PS(|D| \leq d, |X| \geq n)$ построен алгоритм, решающий задачу за время $O(|D_0||X_0|^{2|D_0|+0,5})$.

Основные результаты и выводы

- Рассмотрены различные функционалы качества монотонных классификаторов.
- Рассмотрены различные критерии качества монотонизации выборки.
- Предложена систематизация постановок задач монотонизации выборки как задач дискретной оптимизации.
- Сделаны оценки вычислительной сложности предложенных постановок задачи монотонизации выборки.
- Предложен алгоритм построения монотонного классификатора по монотонной выборке методом отбора эталонных объектов.

Литература



Воронцов К. В., Махина Г. А. Принцип максимизации зазора для монотонного классификатора ближайшего соседа. В: *15-я всероссийская конференция «Математические методы распознавания образов»*. М.:МАКС Пресс, 2011. С.



Kamp R., Feelders A., Barile N. Isotonic classification trees. In: *Proceedings of the 8th International Symposium on Intelligent Data Analysis: Advances in Intelligent Data Analysis VIII*. Berlin, Heidelberg: Springer-Verlag, 2009. P. 405–416.



Sill J. Monotonic networks. In: *Advances in Neural Information Processing Systems*. Ed. Jordan M. I., Kearns M. J., Solla S. A. Cambridge: MIT Press, 1998. P. 661–667.



Marina Velikova, Hennie Daniels. On Testing Monotonicity of Datasets. In: *Learning monotone models*. 2009. P. 11–22.



Гуз И. С. Минимизация эмпирического риска при построении монотонных композиций классификаторов. *Труды МФТИ*. 2011. Т. 3, № 3(11). С.115–121.



Загоруйко Н. Г. *Прикладные методы анализа данных и знаний*. Новосибирск: ИМ СО РАН, 1999.



Кормен Т., Лейзерсон Ч., Ривест Р., Штайн К. *Алгоритмы: построение и анализ, 2-е издание.: Пер. с англ.* М.: Издательский дом «Вильямс», 2005.



Воронцов К. В. О проблемно-ориентированной оптимизации базисов задач распознавания. *ЖВМ и МФ*. 1998. Т. 38. № 5. С. 870–880.