

Предсказание качества для процедуры выбора признаков

Аминов Тимур

Московский физико-технический институт

Курс: Численные методы обучения по прецедентам
(практика, В. В. Стрижов) Группа 674, весна 2019

Цель работы

Предложить алгоритм выбора оптимального подмножества признаков

Проблема

В случае избыточного признакового пространства предсказательная модель машинного обучения является неустойчивой

Метод решения

Использование методов выпуклой оптимизации для получения оптимального подмножества признаков

Пусть $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{m \times n}$ — заданная матрица, где $\mathbf{x}_j \in \mathbb{R}^m$ — j -ый признак, $\mathbf{y} \in \mathbb{R}^m$ — значения функции целевой функции

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} s(\mathbf{w}, \mathbf{X}_{\text{tr}}, \mathbf{y}_{\text{tr}}, \mathcal{A})$$

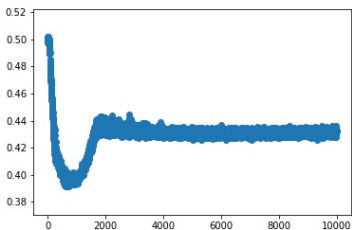
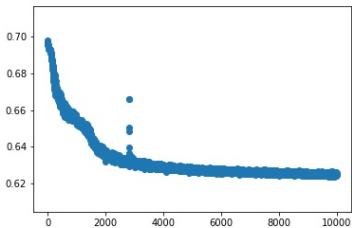
s — функция ошибки Множество $\mathcal{A} \subseteq \{1, \dots, n\}$ — индикатор $\{0, 1\}$.
Существует соответствие между множеством \mathcal{A} и двоичным вектором $\mathbf{a} \in \mathbb{B}^n$:

$$\mathcal{A} = \{j : a_j = 1\}.$$

Типы признаков

- информативные — существенно влияют на точность приближения целевого вектора
- шумовые — не влияют на точность приближения целевого вектора
- мультиколлинеарные — существует линейная зависимость между признаками, снижают устойчивость модели

$$\begin{array}{ccccc}
 a \in \mathbb{B}^n & \xrightarrow[\text{encoder}]{e} & u \in \mathbb{R}^h & \xrightarrow[\text{decoder}]{\delta} & a \in \mathbb{B}^n \\
 & & \downarrow \begin{array}{c} \text{performance} \\ \text{prediction} \end{array} P & & \\
 & & s \in \mathbb{R}^1 & &
 \end{array}$$



Проблема выбора признаков сведена к проблеме непрерывной оптимизации

Пока что сравнивать результаты невозможно