
Разработка эффективной реализации методов, основанных на поиске оптимального баланса между дивергенцией и точностью аппроксимации

A Preprint

Аристархов Данила Дмитриевич
ВМК МГУ
aristarkhov.danila@yandex.ru

Сенько Олег Валентинович
ВМК МГУ
senkoov@mail.ru

Abstract

Рассмотрен новый метод построения ансамбля деревьев при решении задачи регрессии. При этом оптимизация производится исходя из одновременного достижения расходимости алгоритмов в пространстве прогнозов и хорошей аппроксимации данных отдельными алгоритмами ансамбля.

Ключевые слова регрессия · коллективные методы · бэггинг · градиентный бустинг

1 Введение

Методы решения задачи регрессии, основанные на вычислении более точного предсказания с помощью набора менее точных и более простых базовых алгоритмов, получили самое широкое распространение в современном машинном обучении. К числу таких методов может быть отнесен регрессионный случайный лес, а также методы, основанные на использовании адаптивного или градиентного бустинга. Важную роль при построении таких алгоритмов играет способ получения ансамбля так называемых слабых алгоритмов. Теоретический анализ показывает, что увеличение обобщающей способности может быть достигнуто за счет выбора ансамбля алгоритмов, обладающих не только высокой точностью, но и максимально расходящимися прогнозами [Докунин and Сенько, 2015]. В случайном лесе расходимость прогнозов достигается за счет обучения алгоритмов ансамбля на различных выборках, генерируемых из исходной выборки с использованием процедуры бутстрэпа [Breiman, 1996]. В методе градиентного бустинга [Hastie et al., 2009] ансамбль генерируется последовательно. При этом на каждой итерации в ансамбль добавляются деревья, аппроксимирующие первые производные функции потерь в точке, соответствующей текущему прогнозу алгоритма на данном шаге.

Построение ансамбля, имеющего как высокую точность предсказаний, так и сильное расхождение прогнозов отдельных алгоритмов, продемонстрировало улучшение качества модели [Журавлев et al., 2021]. Целью данной работы является продолжение исследования данного подхода.

2 Постановка задачи

Дана выборка $S = (X_1, y_1), \dots, (X_m, y_m)$, где X_i — вектор признаков описания объекта, y_i — метка объекта. Рассматривается задача регрессии: $X_i \in \mathbb{R}^n$, $y_i \in \mathbb{R}$. Требуется построить ансамбль базовых алгоритмов $A_1(X), \dots, A_k(X)$, предсказывающих значения метки по вектору признаков.

Список литературы

А. А. Докунин and О. В. Сенько. Регрессионная модель, основанная на выпуклых комбинациях, максимально коррелирующих с откликом. Ж. вычисл. матем и матем. физ., 53:530–544, 2015.

- L. Brieman. Bagging predictors. *Machine Learning*, 24:123–140, 1996.
- T. Hastie, R. Tibshirani, and J. H. Friedman. *The Elements of Statistical Learning* (2nd ed.). 2009.
- Ю. И. Журавлев, О. В. Сенько, А. А. Докунин, Н. Н. Киселева, and И. А. Саенко. Двухуровневый метод регрессионного анализа, использующий ансамбли деревьев с оптимальной дивергенцией. *Докл. РАН. Матем., информ., проц. упр.*, 499:63–66, 2021.