

Параллельная реализация алгоритма предсказания с помощью модели градиентного бустинга деревьев решений*

П.Н. Дружков, Н.Ю. Золотых, А.Н. Половинкин

Нижегородский государственный университет им. Н.И. Лобачевского

Приводится описание реализации одного из алгоритмов обучения с учителем – градиентного бустинга деревьев решений (Gradient Boosting Trees). Представлены различные варианты параллельной реализации алгоритма предсказания с использованием библиотеки Intel Threading Building Blocks. Приводятся результаты экспериментального сравнения и анализ производительности различных подходов к распараллеливанию.

1. Введение

Машинное обучение является подразделом весьма обширной области науки, изучающей искусственный интеллект. Алгоритмы, относящиеся к данному направлению, используются при решении задач, для которых зачастую сложно или невозможно придумать явный алгоритм решения: предсказание погоды, прогнозирование экономических и социальных процессов, медицинская диагностика, детектирование объектов на фото или видео, распознавание текста, речи, создание антивирусных программ и алгоритмов фильтрации рекламы и спама.

В настоящее время известно достаточно много алгоритмов обучения с учителем, предназначенных для решения задачи восстановления регрессии или классификации: машина опорных векторов [13], метод K ближайших соседей [9], нейронные сети [9], AdaBoost [9], деревья решений [2] и их различные ансамбли (случайные деревья [1], полностью случайные деревья [8], градиентный бустинг деревьев решений [6, 7]). В рамках данной работы рассматривается программная реализация одного из наиболее перспективных алгоритмов обучения с учителем – алгоритма градиентного бустинга деревьев решений (GBT – gradient boosting trees), которая является первой полнофункциональной C/C++ реализацией данного метода с открытым кодом. Результаты вычислительного эксперимента, проведенного с использованием широко распространенных наборов реальных данных, взятых из репозитория UCI [12], свидетельствуют о конкурентоспособности предлагаемой реализации по сравнению с реализациями других алгоритмов. Разработанный код интегрирован в одну из наиболее известных свободно распространяемых библиотек компьютерного зрения OpenCV [11].

Необходимо отметить, что многие из решаемых в настоящее время практических задач машинного обучения и компьютерного зрения требуют обработки значительного объема входных данных. В частности, каждый исследуемый объект может быть описан вектором признаков, содержащим сотни или даже тысячи переменных, а обучающая и тестовая выборки могут содержать десятки тысяч описаний объектов. Наглядным примером такой задачи может служить детектирование пешеходов [4], где в зависимости от выбираемых параметров необходимо классифицировать от 10000 до 185000 объектов на одно изображение. В связи с этим, наряду с качеством предсказания на одно из первых мест встает вопрос производительности используемого алгоритма. В работе рассматриваются аспекты параллельной реализации алгоритма обучения модели, а также предлагаются и анализируются различные подходы к распараллеливанию алгоритма предсказания на новых данных.

* Авторы благодарят И.Б. Меерова (ННГУ) и В. Л. Ерухимова (компания ITSeez) за ценные замечания и полезные обсуждения. Работа выполнена при поддержке федеральной целевой программы «Научные и научно-педагогические кадры инновационной России», госконтракт 02.740.11.5131.

2. Градиентный бустинг деревьев решений

2.1 Постановка задачи

Одной из задач, изучаемой в машинном обучении, является *задача обучения с учителем*. В рамках этой задачи дано некоторое множество *объектов* X . Каждому объекту $x \in X$ поставлена в соответствие величина y , называемая *выходом*, или *ответом*, и принадлежащая множеству допустимых ответов Y . Упорядоченная пара «объект-ответ» (x, y) , где $x \in X$, $y \in Y$ называется *прецедентом*. Требуется восстановить зависимость между входом и выходом, основываясь на данных о конечном наборе прецедентов, называемом *обучающей выборкой*: $\{(x_i, y_i) \mid x_i \in X, y_i \in Y, i = 1, \dots, n\}$. Другими словами, задача состоит в построении функции f из некоторого множества K , которая, получив на вход x , предсказала бы значение ответа y как можно точнее. В случае конечного Y , говорят о *задаче классификации*, если $Y = \mathbf{R}$ – *задаче восстановления регрессии* [9]. Процесс нахождения f называется *обучением* (*тренировкой*, *настройкой*) модели, процесс определения выхода по некоторому входу с помощью уже построенной модели – *предсказанием*.

2.2 Метод решения

Один из общих подходов решения задач обучения заключается в комбинировании моделей. Две основные конкурирующие идеи данного подхода – бэггинг (*bagging* от *Bootstrap Aggregating*) [3] и бустинг (*boosting*) [5]. Первая из них состоит в построении множества независимых между собой моделей с дальнейшим принятием решения путем голосования в случае задачи классификации и усреднения в случае регрессии. Данный подход реализован в алгоритме случайных деревьев (*random trees* или *random forest*). Бустинг, в противоположность бэггингу, обучает каждую следующую модель с использованием данных об ошибках предыдущих моделей.

Алгоритм градиентного бустинга деревьев решений является развитием бустинг-идеи. Он позволяет строить аддитивную функцию в виде суммы деревьев решений итерационно по аналогии с методом градиентного спуска. Данный подход позволяет расширить круг решаемых этим алгоритмом задач, а также зачастую получить выигрыш в точности предсказания.

Далее приводится краткое описание алгоритма обучения градиентного бустинга деревьев решений для задачи восстановления регрессии в случае использования квадратичной функции потерь. Пусть обучающая выборка содержит n прецедентов, y_i – значение ответа для i -го прецедента, $T_j(x_i)$ – значение, предсказанное j -м деревом в ансамбле для i -го объекта, ν – коэффициент масштабирования. Тогда псевдоостатком для i -го объекта на k -м шаге алгоритма обучения называется значение $\tilde{y}_i = y_i - T_0(x_i) - \nu \cdot \sum_{m=1}^{k-1} T_m(x_i)$, что соответствует разности между

истинным значением ответа и значением, предсказанным ансамблем деревьев решений, построенным на $(k - 1)$ -м шаге алгоритма обучения, M – общее число деревьев в ансамбле. Тогда общая схема тренировки модели может быть сформулирована следующим образом:

1. Обучить дерево T_0 на исходном наборе данных $(x_i, y_i), i = 1, \dots, n$
2. Для каждого $m = 1, 2, \dots, M$
 - для всех объектов в обучающей выборке вычислить псевдоостатки $\tilde{y}_i, i = 1, \dots, n$
 - добавить в ансамбль новое дерево, обученное на наборе данных $(x_i, \tilde{y}_i), i = 1, \dots, n$

Детальное описание алгоритма обучения, а также подробности, связанные с тренировкой отдельных деревьев решений, и особенности реализации алгоритма для задач восстановления регрессии с другими функциями потерь, а также классификации с двумя и более классами можно найти в [6].

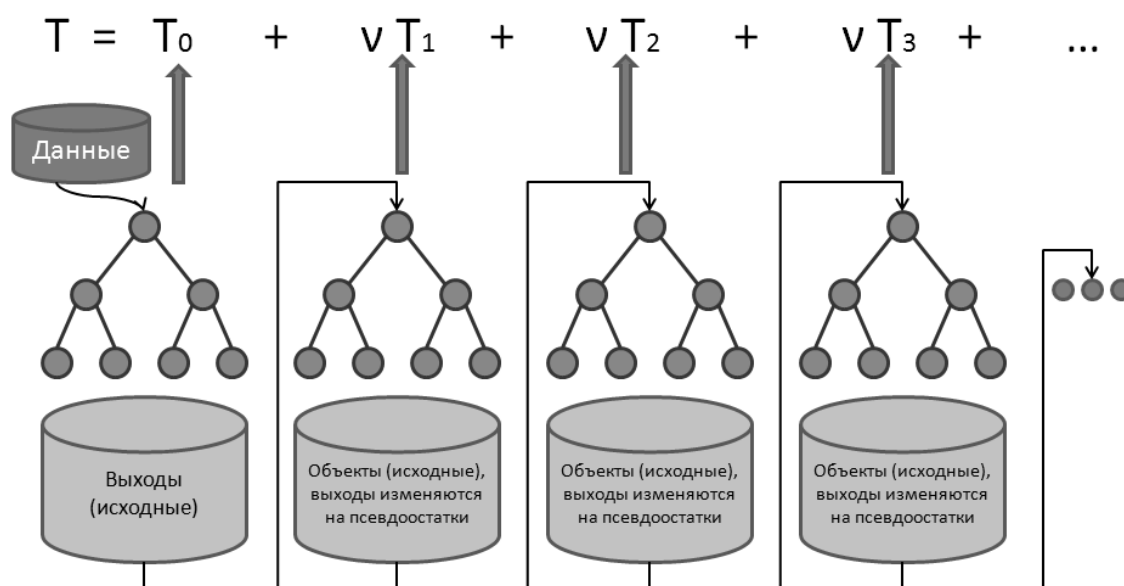


Рис. 1. Алгоритм обучения модели градиентного бустинга деревьев решений

Таким образом, на выходе алгоритма обучения мы получаем набор из M деревьев решений, и для осуществления предсказания, т. е. определения выхода y для нового объекта x , следует вычислить сумму $y = T_0(x) + \nu \cdot \sum_{m=1}^M T_m(x)$.

2.3 Реализация алгоритма градиентного бустинга деревьев решений

Авторами данной работы была выполнена программная реализация алгоритма градиентного бустинга деревьев решений, включающая как алгоритм тренировки модели, так и ее дальнейшее использование для предсказания. В данном разделе приведены некоторые экспериментальные результаты, показывающие достоинства и недостатки метода градиентного бустинга. Наряду с описываемым подходом были рассмотрены конкурирующие алгоритмы: одиночные деревья решений (алгоритм CART) [2], случайные деревья (случайные леса) [1, 8], машина опорных векторов [13]. Программой основой проведенных экспериментов является открытая библиотека компьютерного зрения OpenCV: все результаты, относящиеся к конкурирующим алгоритмам, были получены непосредственно с помощью ее компонентов: CvDTree, CvRTrees, CvERTrees и CvSVM. Эксперименты проводились на наборах реальных данных, взятых из репозитория UCI, при этом мерой качества модели считалась средняя абсолютная ошибка 10-кратного перекрестного контроля.

Таблица 1. Результаты экспериментального сравнения алгоритмов обучения с учителем

Название набора данных	Градиентный бустинг (GBT)	Дерево решений (CvDTree)	Случайные деревья (CvRTrees)	Случайные деревья (CvERTrees)	Машина опорных векторов (CvSVM)
auto-mpg	2	2.238	1.879	2.147	2.981
Computer hardware	12.62	15.62	11.62	9.631	37
Concrete slump	2.257	2.923	2.6	2.359	1.767
Forestfires	18.74	17.26	17.79	16.64	12.9
Boston housing	2.033	2.602	2.135	2.196	4.049
imports-85	1306	1649	1290	1487	1787
Servo	0.238	0.258	0.247	0.42	0.655
Abalone	1.47	1.604	1.492	1.498	2.091

3. Параллельная реализация алгоритма градиентного бустинга деревьев решений

3.1 Алгоритм обучения

Обучение модели представляет собой достаточно трудоемкий процесс. Более того, для осуществления подбора наилучших параметров модели (количество деревьев в ансамбле, масштабный параметр, ограничение на размер деревьев) требуется неоднократное повторение процесса, что обуславливает необходимость оптимизации по скорости, в том числе за счет применения техник параллельных вычислений. Согласно результатам профилировки последовательной версии алгоритма обучения (использовался набор данных `spambase` из библиотеки UCI: размерность пространства признаков равна 56, число прецедентов в обучающей выборке – 2300), приведенным на рис. 2, основное время работы тратится на обучение одиночных деревьев решений.

main	100.0%
CvMLData::read_csv	1.2%
CvGBTrees::train	98.5%
CvGBTrees::train	98.5%
CvDTree::predict	2.2%
CvGBTrees::find_gradient	1.2%
CvGBTrees::change_values	1.5%
CvGBTrees::do_subsample	0.3%
CvDTree::train	93.3%
CvDTree::do_train	93.3%
CvDTreeTrainData::subsample_data	12.5%
CvDTree::try_split_node	80.8%

Рис. 2. Результаты профилировки алгоритма обучения модели градиентного бустинга деревьев решений

К сожалению, бустинг-алгоритмы являются плохо распараллеливаемыми: подход, связанный с одновременным построением нескольких классификаторов в ансамбле, здесь невозможен в силу того, что тренировка каждого следующего дерева решений требует результатов предсказания всех предыдущих. Однако некоторые этапы построения одиночного дерева решений могут быть выполнены независимо друг от друга: например, нахождение разбиения во внутреннем узле дерева, которое осуществляется путем вычисления для каждой переменной уменьшения значения функции неоднородности [2] и выбора оптимального значения среди всех переменных. Именно на этом уровне выполнено распараллеливание обучения деревьев решений в библиотеке `OpenCV`.

Для изучения влияния параллельного подхода к отысканию оптимальной переменной каждого разбиения на скорость работы всего алгоритма обучения модели градиентного бустинга была проведена серия экспериментов.

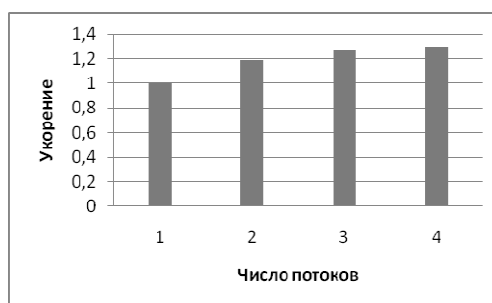


Рис. 3. Эффективность распараллеливания алгоритма обучения модели градиентного бустинга (на наборе данных `spambase`)

Как видно из рис. 3, полученное ускорение значительно ниже линейного. Это связано с тем, что суммарное время работы алгоритма поиска оптимального разбиения составляет (по результатам профилировки) лишь около 40% времени работы всего алгоритма обучения. Таким образом, максимальное ускорение, которое может быть получено, согласно закону Амдала равно 1.67.

3.2 Алгоритм предсказания

Несмотря на то, что предсказание в алгоритме градиентного бустинга деревьев решений не столь трудоемкий процесс по сравнению с построением модели, время работы этого алгоритма зачастую также является критичным. На это есть несколько причин. Обучение выполняется на «оффлайн» этапе, в то время как дальнейшее (и основное) применение модели заключается в осуществлении с ее помощью предсказаний на новых данных. Более того, на практике часто приходится выполнять не единичные предсказания, а целые серии: например, в некоторых алгоритмах, решающих задачу детектирования объектов на изображении, может потребоваться выполнение десятков тысяч предсказаний на одно изображение. Все это накладывает достаточно жесткие временные рамки на время вычисления предсказания – кроме того, в некоторых прикладных задачах требуется работа в режиме реального времени.

В рамках данной статьи рассматривается два различных подхода к распараллеливанию алгоритма предсказания с использованием обученной модели градиентного бустинга деревьев решений. Напомним, что для получения выхода y по некоторому входу x необходимо вычислить сумму предсказаний всех деревьев из имеющегося ансамбля. В отличие от процесса обучения, здесь значение каждого слагаемого может быть получено независимо от остальных – отсюда возникает первая параллельная схема предсказания, в которой для *одного* объекта x

значения предсказаний отдельных деревьев $T_m(x)$ в сумме $y = T_0(x) + \nu \cdot \sum_{m=1}^M T_m(x)$ вычисля-

ются параллельно несколькими потоками. Другой рассматриваемый подход: распараллеливание по данным – основан на необходимости одновременного предсказания для большого количества новых объектов. В данном случае выполняется параллельное вычисление значений

$y_i = T_0(x_i) + \nu \cdot \sum_{m=1}^M T_m(x_i)$ для нескольких объектов x_i . Программная реализация предложен-

ных подходов к распараллеливанию выполнена с использованием технологии Intel Threading Building Blocks [10].

Как видно из рис. 4–5, оба рассматриваемых подхода дают существенное ускорение процесса предсказания. Однако метод распараллеливания по данным обладает лучшей масштабируемостью, так как число объектов, для которых требуется выполнить предсказание, как правило, значительно превосходит количество деревьев в модели градиентного бустинга.

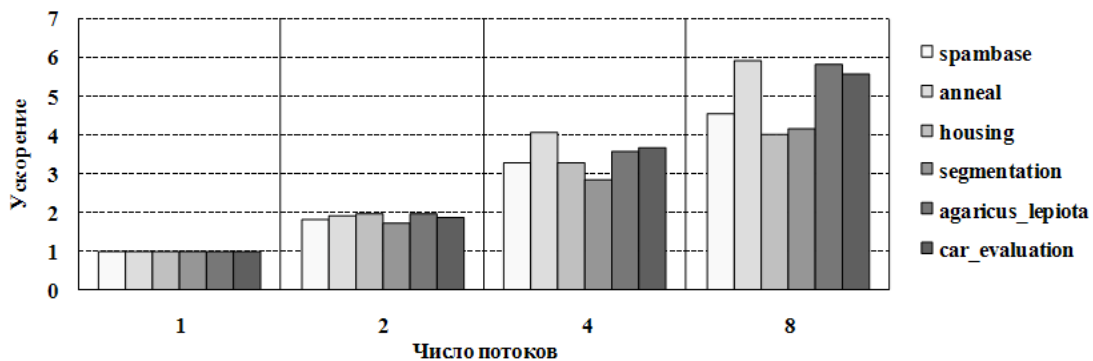


Рис. 4. Эффективность распараллеливания алгоритма предсказания по деревьям в ансамбле с использованием модели градиентного бустинга деревьев решений

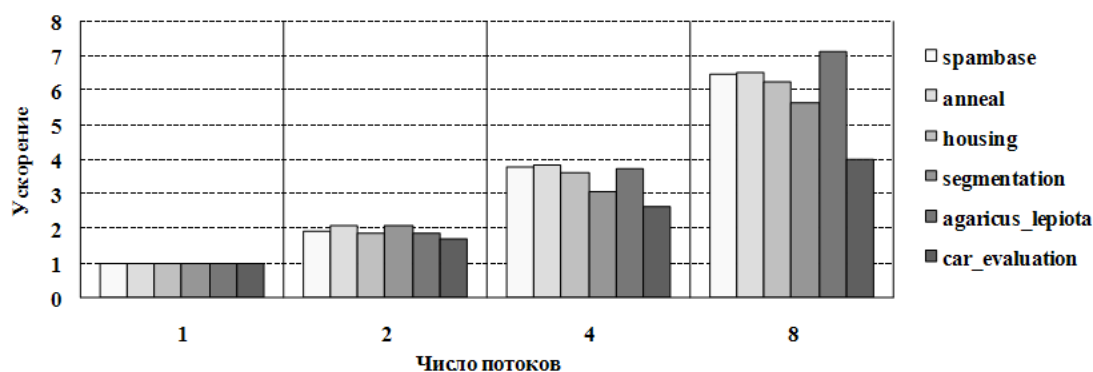


Рис. 5. Эффективность распараллеливания алгоритма предсказания по данным с использованием модели градиентного бустинга деревьев решений

4. Заключение

В статье рассмотрены аспекты параллельной реализации одного из наиболее перспективных алгоритмов обучения с учителем – градиентного бустинга деревьев решений. Проанализирована программная реализация обучения рассматриваемой модели, основанная на параллельной версии алгоритма тренировки отдельных деревьев решений, входящей в состав библиотеки OpenCV. Также в рамках данной работы были предложены и реализованы с использованием технологии Intel Threading Building Blocks две схемы распараллеливания алгоритма предсказания с помощью построенной модели.

Литература

1. Breiman L. Random Forests. // *Machine Learning*. 2001. V. 45, №. 1, P. 5-32.
2. Breiman L., Friedman J., Olshen R., Stone C. *Classification and Regression Trees*. Wadsworth, 1983.
3. Breiman L. Bagging predictors // *Machine Learning*. 1996. V. 26, № 2, P. 123-140.
4. Enzweiler M., Gavrilu D. M. Monocular Pedestrian Detection: Survey and Experiments // *IEEE Transactions on Pattern Analysis and Machine Intelligence*. 2009. V.31, № 12. P. 2179–2195.
5. Freund Y., Schapire R. Experiments with a New Boosting Algorithm // *Machine Learning: Proceedings of the Thirteenth International Conference*. 1996.
6. Friedman J. H. Greedy Function Approximation: a Gradient Boosting Machine. Technical report. Dept. of Statistics, Stanford University, 1999.
7. Friedman J. H. Stochastic Gradient Boosting. Technical report. Dept. of Statistics, Stanford University, 1999.
8. Geurts P., Ernst D., Wehenkel L. Extremely Randomized Trees // *Machine Learning*. 2006. V. 36, № 1. P. 3–42.
9. Hastie T., Tibshirani R., Friedman J. *The Elements of Statistical Learning*. Springer, 2008.
10. Intel Threading Building Blocks.
URL: <http://www.threadingbuildingblocks.org> (дата обращения: 07.12.2010).
11. OpenCV Wiki.
URL: <http://opencv.willowgarage.com/wiki> (дата обращения: 07.12.2010).
12. UCI Machine Learning Repository.
URL: <http://archive.ics.uci.edu/ml> (дата обращения: 07.12.2010).

13. Вапник В. Н. Восстановление зависимостей по эмпирическим данным. М.: Наука, 1979. 448 с.