

Московский государственный университет имени М. В. Ломоносова



Факультет Вычислительной Математики и Кибернетики

Кафедра Математических Методов Прогнозирования

КУРСОВАЯ РАБОТА СТУДЕНТА 317 ГРУППЫ

«Сравнительный анализ методов быстрого поиска ближайших соседей»

Выполнил:

студент 3 курса 317 группы

Федоров Илья Сергеевич

Научный руководитель:

д.ф.-м.н., профессор

Дьяконов Александр Геннадьевич

Заведующий кафедрой

Математических Методов

Прогнозирования, академик РАН

_____ Ю. И. Журавлёв

К защите допускаю

«_____» _____ 2020 г.

К защите рекомендую

«_____» _____ 2020 г.

Москва, 2020

Содержание

1	Введение	3
1.1	Определения и обозначения	3
1.2	Обзор литературы	4
2	Новые подходы и результаты	5
3	Вычислительные эксперименты	5
3.1	Исходные данные и условия эксперимента	5
3.2	Результаты эксперимента	5
3.3	Обсуждение и выводы	6
4	Заключение	6

Аннотация

Todo

1 Введение

Одним из наиболее простых и естественных методов машинного обучения является метод ближайшего соседа. Имея набор данных, представленных в виде точек в некотором многомерном пространстве, целевая величина (будь то класс или вещественное число) прогнозируется по значениям отклика на k ближайших к запросу точках из исходного набора данных. В то время как существует множество различных подходов к усреднению данных k значений, наиболее вычислительно затратной частью алгоритмов подобного типа является именно поиск ближайших соседей. Действительно, в современных задачах объемы данных достигают колоссальных размеров, что делает алгоритмы, основанные на полном переборе, неэффективными. Задача поиска ближайших к запросу точек в некотором наборе данных встречается не только в задачах прогнозирования. Примерами приложений могут служить задачи поиска дубликатов в больших объемах данных (или «почти» дубликатов), поиска похожих изображений и текстов. Целью данной работы является обзор современных подходов к решению задачи поиска ближайших соседей и её вариаций, а также сравнительный анализ эффективности тех или иных методов её решения в зависимости от особенностей пространства, в котором расположены данные. В исследовании представлены как классические подходы, основанные на формировании некоторых дополнительных структур данных, так и наиболее современные «приближенные» методы.

1.1 Определения и обозначения

Формализуем постановку задачи. Основным объектом нашего изучения будет пространство признаков вместе с функцией расстояния $\mathbb{X} = (\mathbb{R}^n, d)$. Важно заметить, что функция d в приложениях довольно часто может не удовлетворять формальному определению метрики, однако даже в этом случае в данной работе подобные функции, допуская некоторую вольность, будут называться метриками. К примеру, широко используемое в анализе текстов косинусное расстояние $d(x, y) = 1 - \frac{\langle x, y \rangle}{\|x\| \|y\|}$ не удовлетворяет неравенству треугольника. В данной работе в большинстве случаев будет использоваться евклидова метрика и описанное выше косинусное расстояние.

Будем обозначать $X \in \mathbb{R}^{l \times n}$ матрицу для выборки точек из \mathbb{X} , где строки соответствуют объектам, а столбцы признаками. Формально задача поиска k ближайших соседей ставится следующим образом: имея множество объектов X из пространства \mathbb{X} и запрос $q \in \mathbb{X}$, нужно найти в X k ближайших к q точек по метрике d . Более подробно, если посчитать расстояния между q и всеми объектами из X , а потом расположить их в отсортированном порядке

$$d(x_{i_1}, q) \leq d(x_{i_2}, q) \leq \dots \leq d(x_{i_l}, q),$$

то алгоритм должен выдать k объектов с минимальными расстояниями: $x_{i_1}, x_{i_2}, \dots, x_{i_k}$.

Как мы увидим далее, большинство современных методов быстрого поиска ближайших соседей на самом деле решают описанную выше задачу с некоторыми ослаблениями, что в сущности приводит к задаче «приближенного» поиска ближайших соседей. Эта задача не имеет общепризнанной формальной постановки, разные авторы могут по-разному понимать её. К примеру, довольно распространенной постановкой задачи является следующая формулировка: имея набор точек X из \mathbb{X} , запрос $q \in \mathbb{X}$, а также параметр $c \geq 1$, если существует точка $x \in X$, такая что $d(x, q) \leq r$, алгоритм должен вернуть точку $x^* \in X$, такую что $d(x^*, q) \leq cr$. В дальнейшем при описании конкретных методов, мы будем уточнять, какую именно задачу приближенного поиска ближайших соседей они решают.

1.2 Обзор литературы

Лучше, чтобы название этого подраздела было содержательным, например, общепринятым названием задачи, проблемы или метода, рассматриваемого в данной работе.

Перечисляются подходы, методы, факты, на которые существенно опирается данная работа, но которые могут быть не известны широкому кругу читателей. Здесь ссылки на литературу обязательны. Теоремы только формулируются, но не доказываются.

Данный раздел преследует две цели. Во-первых, сделать работу самодостаточной — дать необходимый минимум информации тем читателям, которые не очень хорошо ориентируются в теме, но желают поближе познакомиться именно с дан-

ной работой. Во-вторых, облегчить сопоставление полученных автором результатов с ранее известными.

2 Новые подходы и результаты

Название этого раздела обязательно надо заменить на содержательное. В этом разделе, как правило, много подразделов.

В дипломной работе не стоит делать более двух уровней, достаточно разделов и подразделов. Будете писать диссертацию или монографию — сделаете три уровня.

3 Вычислительные эксперименты

Цель данного раздела: продемонстрировать, что предложенная теория работает на практике; показать границы её применимости; рассказать о новых экспериментальных фактах.

Чисто теоретические работы могут вообще не содержать раздела экспериментов (не работает, ну и не надо — зато теория красивая). Кстати, теоретики имеют право не догадываться, где, кому и когда их теории пригодятся.

3.1 Исходные данные и условия эксперимента

Описывается прикладная задача, параметры анализируемых данных (например, сколько объектов, сколько признаков, каких они типов), параметры эксперимента (например, как производился скользящий контроль).

3.2 Результаты эксперимента

Результаты экспериментов представляются в виде таблиц и графиков. Объясняется точный смысл всех обозначений на графиках, строк и столбцов в таблицах.

3.3 Обсуждение и выводы

Приводятся выводы: в какой степени результаты экспериментов согласуются с теорией? Достигнут ли желаемый результат? Обнаружены ли какие-либо факты, не нашедшие объяснения, и которые нельзя списать на «грязный» эксперимент?

Обсуждаются основные отличия предложенных методов от известных ранее. В чем их преимущества? Каковы границы их применимости? Какие проблемы удалось решить, а какие остались открытыми? Какие возникли новые постановки задач?

4 Заключение

В квалификационных работах последний раздел нужен для того, чтобы конспективно перечислить основные результаты, полученные лично автором.

Результатами, в частности, являются:

- Предложен новый подход к...
- Разработан новый метод..., позволяющий...
- Доказан ряд теорем, подтверждающих (опровергающих), что...
- Проведены вычислительные эксперименты..., которые подтвердили / опровергли / привели к новым постановкам задач.

Цель данного раздела: доказать квалификацию автора. Даже беглого взгляда на заключение должно быть достаточно, чтобы стало ясно: автору удалось решить актуальную, трудную, ранее не решённую задачу, предложенные автором решения обоснованы и проверены.

Иногда в Заключении приводится список направлений дальнейших исследований.

Список литературы необходим в любой научной публикации. В дипломной работе он обязателен. Дурным тоном считается: ссылаться на работы только одного-двух авторов (например, себя или шефа); ссылаться на слишком малое число работ; ссылаться только на очень старые работы; ссылаться на работы, которых автор ни разу не видел; ссылаться на работы, которые не упоминаются в тексте или которые не имеют отношения к данному тексту.