

Новые методы сжатия временных рядов экологических показателей

Чуприн В.И., Родригес Залепинос Р.А.
Донецкий национальный технический университет
chuprin.vladislav@gmail.com, rodrigues@csm.donntu.edu.ua

Abstract

Chuprin V.I., Rodrigues Zalipynis R.A. "New time series compression methods of ecologic parameters". The paper analyzes storage peculiarities of satellite Earth remote sensing data time series. We propose methods for their compression based on the discovered peculiarities exploiting different schemes of Huffman coding. One of the proposed methods reaches 6% increase in the compression ratio (93%) in contrast to the deflate method used in Java SE6 (87%), for a time series of aerosol optical thickness derived from MODIS radiometer of TERRA satellite. Further improvement can be achieved by using the entropy coding of floating point numbers.

Keywords: time series, lossless compression, Huffman coding, arithmetic compression, floating-point data, Aura, TERRA, MODIS, OMI.

Введение

Анализ временных рядов является одним из важнейших инструментов оценки и прогнозирования состояния окружающей среды [1].

Новые средства мониторинга, генерируют большие объемы данных [2]. Например, космический исследовательский аппарат Aura предоставляет 26 гигабайт информации ежедневно [3].

Большие объемы данных являются причиной увеличения времени считывания и передачи информации по сети. Возрастают требования к размеру дисковых накопителей. Появляется необходимость использования полос высоких пропускных способностей, что часто приводит к повышению финансовых расходов. Одним из решений приведенных выше проблем является использование методов сжатия данных.

В статье проанализированы особенности временных рядов оптической толщины аэрозоля, полученной с радиометра MODIS спутника TERRA [4], а также концентрации озона, определенной радиометром OMI спутника Aura [5]. Предложены модификации существующих методов с использованием приведенных особенностей, за счет чего достигнуто увеличение степени сжатия.

Постановка задачи

Для решения поставленных проблем с использованием сжатия предлагаются методы, удовлетворяющие следующим требованиям:

– обеспечение сжатия данных без потерь точности;

– принятие во внимание особенностей рассмотренных экологических показателей, включая представление значений числами с плавающей точкой;

– наличие возможности использования блочных методов;

– допустимо использование методов с ресурсоемкой декомпрессией, так как нагрузка переносится на клиентские машины, которые обладают достаточными вычислительными мощностями.

Следует определить основные источники, порождающие избыточность, и подобрать нужный комплекс методов, удовлетворяющий поставленным требованиям.

Обзор существующих методов

Литература по вопросу сжатия временных рядов с плавающей точкой довольно разрежена. Сообщество разработало ряд схем сжатия неструктурированных данных [6], которые являются эффективными для 3D моделей основанных на точечном представлении. Приведенные алгоритмы не оправдывают себя при сжатии временных рядов экологических показателей, так как не учитывают существующие в них закономерности.

Одним из требований к сжатию временных рядов, является отсутствие потерь, что обусловлено накоплением большой погрешности при анализе. В большинстве случаев гораздо важнее получить точные показатели, чем сэкономить время при вычислениях. В статье не рассмотрены методы, которые хорошо зарекомендовали себя при сжатии многомерных данных с потерями, такие

как метод выделения структур, фрактальные методы, дискретное преобразование Фурье, дискретное косинусное преобразование.

Структура и особенности временных рядов

Временные ряды доступны для каждой ячейки регулярной широтно-долготной решетки земной поверхности за определенный период с заранее заданным шагом по временной шкале. Обладают низкой энтропией. Наблюдается дублирование данных как для близких по времени значений, так и при изменении положения внутри плоскости выбранного участка. Большая часть данных является неопределенной, по причине наличия облачности во время сканирования территории радиометром. Для хранения показателей обычно используются числа с плавающей точкой одинарной точности.

В большинстве случаев обрабатываются заранее подготовленные данные, при этом не требуется передача показателей в реальном времени. Потребность в накоплении данных делает более приоритетным использование блочных алгоритмов сжатия. Однако, дополнительно некоторый прирост коэффициента сжатия, можно получить используя поточные алгоритмы сжатия чисел с плавающей точкой [7]. Степень прироста напрямую зависит от использованного блочного алгоритма.

Благодаря низкой энтропии данных эффективными будут статистические методы. Прирост может дать предварительная модификация структуры при помощи преобразующих методов. Чем сложнее структура временного ряда, тем сильнее оптимальное преобразование данных улучшит сжатие. Следует рассмотреть метод разделения мантисс и экспонент для увеличения однородности показателей. Увеличение длины одинаковых фраз можно достигнуть благодаря использованию различных способов обхода координатной сетки. При этом, следует выделить обход методом зигзаг-сканирования.

Метод Хаффмана

Классический блочный алгоритм сжатия на основе статистических показателей, разработан Дэвидом Хаффманом в 1952 году. Идея алгоритма состоит в том, что символы с наибольшей частотой получают короткие коды, а с наименьшей частотой более длинные. Этот метод, требует дополнительного прохода по входному блоку для сбора статистики. Существуют адаптивные варианты, не требующие дополнительного прохода [8]. Алгоритм построен так, что присваивает

каждому символу код с целым количеством бит. Это порождает наилучшие коды переменной длины, когда вероятности символов алфавита являются степенями числа 2.

Благодаря однозначному отображению исходных элементов на биты сжатой последовательности, метод позволяет использовать оптимизацию, основанную на преобразовании цепочек. Цепочка кодируется идентификатором начала, длиной и символом, из которого эта цепочка состоит. Маркер начала не должен встречаться в исходном алфавите. Когда сжатая последовательность занимает больше места, чем исходная цепочка, в выходной поток записываются исходные символы.

Арифметический метод

Арифметический метод присваивает код не отдельным символам, а всей входной последовательности. Это позволяет решить проблему эффективности, когда относительные частоты не являются степенями числа 2 [9].

При данном подходе длинные цепочки одинаковых элементов, лишь косвенно влияют на коэффициент сжатия, путем увеличения вероятности встречаемого символа. Метод может быть обобщен для алфавита с произвольной длиной символа.

Преобразующие оптимизации

Оптимальное преобразование данных дает возможность увеличить эффективность сжатия, за счет увеличения однородности. Приведенные в данной статье оптимизации позволяют достичь прирост коэффициента сжатия при использовании метода Хаффмана с кодированием цепочек. Для арифметического метода результаты останутся неизменными в силу описанных выше особенностей.

В связи с малой энтропией данных для большинства показателей наблюдается одинаковое значение экспонент. Благодаря хранению экспонент отдельно от мантисс удалось получить длинные цепочки одинаковых значений.

Методы обхода координатной сетки

Данные временного ряда экологических показателей можно представить в виде трехмерной сетки значений для пространственно-временных координат. Указанная сетка ограничена прямоугольным параллелепипедом. Стороны соответствуют долготе, широте и значению времени для показателя.

В исходном виде, выполняется обход по временной шкале для определенного значения

координат. Использование иных обходов позволяет достигнуть большей однородности обрабатываемых данных.

Обход строками по координатной плоскости

В сравнении с временной шкалой, данные показателей, с большей вероятностью, имеют одинаковые значения для близких точек внутри координатной плоскости широты-долготы. Указанную закономерность подтверждает рисунок 1.

На рисунке 1b приведена координатная сетка для значения 0,0782 временного ряда концентрации озона радиометра OMI спутника Ауга. Наблюдаемая закономерность выражается в виде вытянутых вдоль плоскости цилиндрических скоплений (рис. 1a). С учетом этого, целесообразно применить обход строками по координатной плоскости.

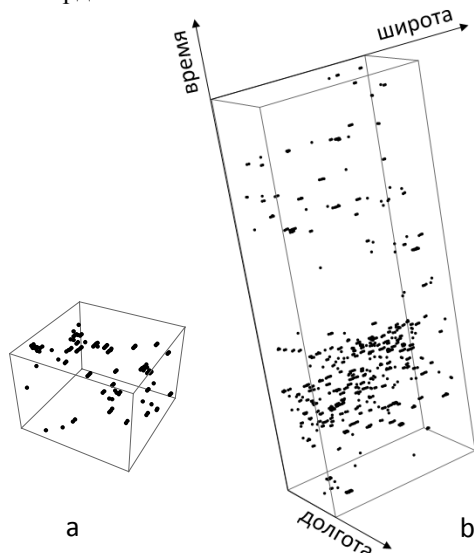


Рисунок 1. – Скопления одинаковых показателей

Метод зигзаг-сканирования

Метод эффективен в случаях, когда значения изменяются не вдоль строк, а по диагонали плоскости. Этот метод, используется в алгоритме сжатия изображений JPEG.

Обход массива показателей начинается с одного угла плоскости и заканчивается в противоположном (рис. 2).

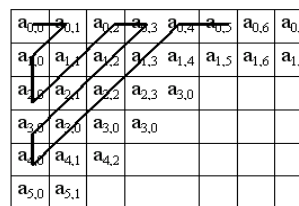


Рисунок 2. – Зигзаг-сканирование

Сравнение эффективности предложенных модификаций

Получены коэффициенты сжатия показателей для следующих временных рядов.

а) Оптическая толщина аэрозоля радиометра MODIS спутника TERRA для области: южная широта 30,0°, северная широта 60,0°, западная долгота 0,0°, восточная долгота 50,0° для уровня широтно-долготной решетки 1,0° x 1,0°, временным интервалом с 01.03.2000 по 01.09.2000 с шагом в 1 день.

б) Концентрация озона радиометра OMI спутника Ауга для области: южная широта 43,0°, северная широта 54,0° западная долгота 20,0°, восточная долгота 42,0° для уровня широтно-долготной решетки 0,25° x 0,25° и временным интервалом с 01.01.2005 по 31.07.2005 с шагом в 1 день.

Ниже приведены использованные для сравнения методы.

1. Встроенный в Java SDK алгоритм сжатия без потерь, который основан на методе deflate [10].
2. Арифметический метод для сжатия однобайтовых элементов.
3. Обобщенный арифметический метод для чисел с плавающей точкой одинарной точности.
4. Метод Хаффмана для однобайтовых элементов.
5. Обобщенный метод Хаффмана для чисел с плавающей точкой одинарной точности.
6. Обобщенный метод Хаффмана с кодированием одинаковых последовательностей элементов (обход вдоль временной шкалы).
7. Оптимизированный метод Хаффмана с обходом по координатной плоскости.
8. Метод Хаффмана с обходом по плоскостям с использованием зигзаг-сканирования.

Таблица 1. – Полученные коэффициенты сжатия временных рядов

	Характеристики данных				Результаты работы алгоритмов							
	уникальность (%)	min	max	NA* (%)	1	2	3	4	5	6	7	8
a	0,2	238,2	495,0	46,5	87,91	70,93	93,5	61,74	92,26	92,86	94,47	94,44
b	0,6	-0,05	5,0	84,98	74,28	49,66	79,45	45,66	76,39	79,56	80,79	80,04

* - количество неопределенных значений

Нижче приведені закономірності, виведені на основі отриманих результатів.

1. В загальному випадку, арифметичний метод виявляється більш ефективним, ніж метод Хаффмана. Це обумовлює велику його популярність на даний момент.

2. Методи адаптовані для чисел з плаваючою точкою працюють значно ефективніше, ніж ті ж методи для байтових елементів. По-перше, завдяки зменшенню кількості біт, необхідних для кодування кожного окремого символу. При цьому, приріст коефіцієнта стиснення досягається за рахунок зменшення кількості символів і збільшення ймовірності їх появи. По-друге, спостерігається зменшення таблиці перекодування, яка використовується в обох наведених методах.

3. Для великої кількості записів метод Хаффмана виявляється більш ефективним завдяки оптимізації, що дозволяє кодувати ці записи в більш короткі, виключаючи дублювання. Арифметичний метод лише косвенно враховує довгі послідовності повторюваних символів.

4. Приріст ефективності стиснення приносять методи обходу координатної площини на основі закономірностей представлення даних, що досягається оптимальним збереженням однорідних записів.

Висновки

Завдяки практичним дослідженням на основі даних оптичної товщини аерозолю, отриманих з радіометра MODIS супутника TERRA і концентрації озона, визначеної радіометром OMI супутника Аюга, продемонстровано ефективність підходу, заснованого на методі Хаффмана. Вказаний метод було узагальнено для елементів довільної довжини з метою застосування до 4х-байтних чисел з плаваючою точкою. Для мінімізації обсягу, який займають дубльовані дані, було запропоновано схему кодування повторюваних елементів. Було також апробовано метод обходу даних, що дозволяє збільшити довжини записів однакових елементів і тим самим підвищити однорідність стиснутих даних. Наведена схема дозволяє збільшити коефіцієнт стиснення без втрати точності вихідних даних.

В подальшому слід розглянути наступні оптимізації:

– різні методи обходу координатної сітки, які використовують статистичні показники, отримані на основі стиснутих даних за допомогою методів динамічного програмування;

– кодування ентропії чисел з плаваючою точкою існуючими поточними методами стиснення;

– застосування підходів, які б дозволяли попередньо оцінювати ефективність використання наведених методів. В разі високої неоднорідності даних, слід використовувати алгоритми більш ефективні в загальному випадку.

Наведені підходи можуть бути також ефективні для даних, що мають малу ентропію і складаються з елементів довільного розміру.

Література

- Rodrigues Zalipynis, R.A. Representing Earth remote sensing data as time series / R.A. Rodrigues Zalipynis // Donetsk National Technical University. Series: System analysis and information technology in environmental and social sciences, No.2, 2012. – 212 с. – С. 57–71.
- Родригес Заліпінос Р.А. Данні і методи інтелектуального аналізу даних для дослідження оточуючої природної середовища [Текст]: наук. журн. / Родригес Заліпінос Рамон Антоніо // Наук. праці Донецького національного технічного університету. – Сер.: Системний аналіз і інформаційні технології в науках про природу і суспільство, No.1, 2011. – 214 с. – С. 94–107.
- Up in the Air [Електронний ресурс] – Режим доступу: <http://www.iwu.edu/magazine/2004/winter/aura.html> (14.11.2012).
- MODIS Website [Електронний ресурс] – Режим доступу: <http://modis.gsfc.nasa.gov/> (10.10.2012).
- OMI Instrument Science [Електронний ресурс] – Режим доступу: <http://aura.gsfc.nasa.gov/instruments/omi.html> (10.09.2012).
- Devillers O., Gandoin P.-M. Geometric Compression For Interactive Transmission [Електронний ресурс] – Режим доступу: – <http://hal.inria.fr/docs/00/07/27/43/PDF/RR-3910.pdf> (18.10.2012).
- Lindstrom P., Isenburg M. Fast and Efficient Compression of Floating-Point Data [Електронний ресурс] – Режим доступу: <http://dl.acm.org/citation.cfm?id=1187859> (14.10.2012).
- Ватолин Д., Ратушняк А., Смирнов М. Методи стиснення даних. М.: Диалог-МИФИ, 2003. – 381 с. – С. 31.
- Сэлмон Д. Стиснення даних, зображень і звуку, М.: Техносфера, 2004. – 367 с. – С. 62.
- Class GZIPInputStream [Електронний ресурс] Режим доступу: <http://docs.oracle.com/javase/1.5.0/docs/api/java/util/zip/GZIPInputStream.html> (14.10.2012).