

**ПАРАЛЛЕЛЬНЫЕ ВЫЧИСЛЕНИЯ В ОБРАБОТКЕ ДАННЫХ
ДИСТАНЦИОННОГО ЗОНДИРОВАНИЯ ЗЕМЛИ**

Асмус В.В.¹, Бучнев А.А.², Пяткин В.П.²

¹ Научно-исследовательский центр космической гидрометеорологии “Планета”
Росгидромета РФ, г. Москва

² ИНСТИТУТ ВЫЧИСЛИТЕЛЬНОЙ МАТЕМАТИКИ И МАТЕМАТИЧЕСКОЙ ГЕОФИЗИКИ СО РАН

В последние годы характерными чертами дистанционного зондирования Земли являются, с одной стороны - повышение пространственного разрешения космических снимков, с другой - использование гиперспектральной съемки с большим числом спектральных диапазонов. Синтезированные многоспектральные данные, получаемые при совмещении изображений из нескольких спектральных каналов, могут иметь объемы в сотни мегабайт. Становится типичным объем данных из $\approx 10^8$ многоспектральных векторов [1]. Поскольку анализ и обработка аэрокосмических изображений подобного объема и сложной гиперспектральной структуры требуют использования всех ресурсов доступных аппаратно-программных средств, в ИВМиМГ СО РАН разрабатывается технология решения этих задач на высокопроизводительных компьютерах базовых вычислительных комплексов, основанная на распределенной и параллельной обработке аэрокосмической информации в локальных сетях суперкомпьютерных центров. На суперкомпьютерах RM600 и MBS-1000/M используется система параллельного программирования MPI (Message Passing Interface), которая обеспечивает переносимость пакетов программ параллельной обработки и анализа аэрокосмических изображений на большинство известных многопроцессорных параллельных ЭВМ. В обработке гиперспектральных данных дистанционного зондирования (ДДЗ) можно выделить следующие этапы, которые требуют параллельных версий соответствующих алгоритмов из-за большого объема вычислений.

Контролируемая классификация (классификация с обучением). Разработанная в НИЦ “Планета” и ИВМиМГ СО РАН система классификации состоит из семи классификаторов (одного поэлементного классификатора и шести объектных), основанных на использовании байесовской стратегии максимального правдоподобия для нормально распределенных векторов признаков, и двух объектных классификаторов, основанных на минимуме расстояния. Под элементом здесь понимается N -мерный вектор признаков $x = (x_1, \dots, x_N)^T$, где N - число спектральных диапазонов, а под объектом блок смежных векторов квадратной или крестообразной формы. Предполагается, что векторы x имеют в классе ω_i нормальное распределение $N(m_i, B_i)$ со средним m_i и ковариационной матрицей B_i . В этом случае стратегия максимального правдоподобия для поэлементного классификатора формулируется следующим образом.

Пусть $\Omega = (\omega_1, \dots, \omega_m)$ - конечное множество классов, $p(\omega_i)$ - априорная вероятность класса ω_i . Дискриминантная функция класса ω_i имеет вид

$$g_i(x) = \ln(p(\omega_i)) - 0.5(\ln(|B_i|) + (x - m_i)^T B_i^{-1}(x - m_i)). \quad (1)$$

Обозначим через T_i основанное на распределении χ^2 пороговое значение для отклоненных векторов класса ω_i : $T_i = \ln(p(\omega_i)) - 0.5A(N, Q) - 0.5\ln(|B_i|)$, где $A(N, Q)$ – критическое значение уровня Q распределения χ^2 . Тогда решающее правило для данного классификатора принимает следующий вид: элемент x заносится в класс ω_i , если $g_i(x) > g_j(x)$ для всех $j \neq i$ и $g_i(x) > T_i$. В противном случае элемент заносится в класс отклоненных векторов.

Поскольку физические размеры реально сканируемых пространственных объектов, как правило, больше разрешения съемочных систем, между векторами данных существуют взаимосвязи. Использование информации подобного рода дает возможность повысить точность классификации, если пытаться распознавать одновременно группу смежных векторов – объект в приведенном выше смысле. В частности, в систему классификации входят классификаторы OMARKC, OMARKS – объектные контролируемые классификаторы, основанные на стратегии максимального правдоподобия, для объектов в форме перекрестия и квадрата соответственно в предположении, что объект является случайным гауссовским марковским полем второго порядка со следующей разделимой корреляционной функцией

$corr(x_{ij}, x_{kl}) = \rho_h^{|i-k|} \cdot \rho_v^{|j-l|}$, где x_{ij} - вектор признаков на пересечении i -ой строки и j -ого столбца, ρ_h - коэффициент пространственной корреляции между двумя соседними векторами в

горизонтальном направлении, ρ_v - коэффициент пространственной корреляции между двумя соседними векторами в вертикальном направлении.

Пусть i - номер класса, I, J - размеры образа. Тогда $\rho_{hi} = \frac{1}{N} \sum_{k=1}^N \rho_{hi}^k$, где

$$\rho_{hi}^k = \frac{1}{I(J-1)} \sum_{t=1}^I \sum_{j=1}^{J-1} (x_{tj}^k - m_{li}^k)(x_{t,j+1}^k - m_{ri}^k) / \sqrt{\text{Var}_{li}^k \text{Var}_{ri}^k}, \quad 1 \leq k \leq N.$$

Соответственно

$$m_{li}^k = \frac{1}{I(J-1)} \sum_{t=1}^I \sum_{j=1}^{J-1} x_{tj}^k, \quad m_{ri}^k = \frac{1}{I(J-1)} \sum_{t=1}^I \sum_{j=2}^J x_{tj}^k, \quad \text{Var}_{li}^k = \frac{1}{I(J-1)-1} \sum_{t=1}^I \sum_{j=1}^{J-1} (x_{tj}^k - m_{li}^k)^2,$$

$$\text{Var}_{ri}^k = \frac{1}{I(J-1)-1} \sum_{t=1}^I \sum_{j=2}^J (x_{tj}^k - m_{ri}^k)^2.$$

Аналогичные формулы получаются для ρ_{vi} .

Дальнейшее изложение приводится для объектов размером 3×3 . Покажем явно размещение векторов

$$\text{в объектах: } X = \begin{pmatrix} x_5 & & \\ x_4 & x_1 & x_2 \\ & x_3 & \end{pmatrix} \text{ для перекрестия, } X = \begin{pmatrix} x_1 & x_2 & x_3 \\ x_4 & x_5 & x_6 \\ x_7 & x_8 & x_9 \end{pmatrix} \text{ для квадрата.}$$

Для классификатора OMARKC (форма объекта – перекрестие) имеем решающие функции вида

$$G_i(X) = g_i(x_1) - 0.5 \sum_{j=2}^5 g_i^1(x_j | x_1).$$

Здесь $g_i(x_1)$ - дискриминантная функция (1) поэлементного классификатора, а в качестве примера для функций $g_i^1(x_j | x_1)$ распишем функцию $g_i^1(x_2 | x_1)$ (остальные получаются аналогично)

$$g_i^1(x_2 | x_1) = N \ln(1 - \rho_{hi}^2) + \ln|B_i| + \frac{1}{1 - \rho_{hi}^2} (D_{22}^i - 2\rho_{hi} D_{21}^i + \rho_{hi}^2 D_{11}^i), \quad (2)$$

$$\text{Где } D_{kl}^i = (x_k - m_i)^T B_i^{-1} (x_l - m_i). \quad (3)$$

Центральный элемент x_1 объекта X относится в класс ω_i , если $G_i(X) > G_j(X)$ для всех $j \neq i$ и $G_i(X) > T_i^c$, где T_i^c - пороговое значение для отклоненных векторов класса ω_i :

$$T_i^c = \ln p(\omega_i) - 0.5(A(NL, Q) + (2[L/2] + 1) \ln|B_i| + N[L/2] \ln((1 - \rho_{hi}^2)(1 - \rho_{vi}^2))). \quad (4)$$

Здесь $A(NL, Q)$ - критическое значение уровня Q распределения χ^2 , а L - количество векторов в объекте.

Для классификатора OMARKS (форма объекта – квадрат) имеем следующие решающие функции

$$G_i(X) = g_i^2(x_1 | x_2, x_4, x_5) + g_i^2(x_2 | x_3, x_5, x_6) + g_i^2(x_4 | x_5, x_7, x_8) + g_i^2(x_5 | x_6, x_8, x_9) + g_i^1(x_3 | x_6) + g_i^1(x_6 | x_9) + g_i^1(x_7 | x_8) + g_i^1(x_8 | x_9) + g_i(x_9),$$

где $g_i(x_9)$ - дискриминантная функция (1) поэлементного классификатора, $g_i^1(x_j | x_l)$ - функции вида (2), а в качестве примера для функций g_i^2 распишем функцию $g_i^2(x_1 | x_2, x_4, x_5)$ (остальные получаются аналогично):

$$g_i^2(x_1 | x_2, x_4, x_5) = -0.5(\ln|B_i| + N \ln((1 - \rho_{hi}^2)(1 - \rho_{vi}^2))) + \frac{1}{(1 - \rho_{hi}^2)(1 - \rho_{vi}^2)} (D_{11}^i - 2\rho_{hi} D_{12}^i - 2\rho_{vi} D_{14}^i + 2\rho_{hi} \rho_{vi} D_{15}^i + \rho_{hi}^2 D_{22}^i + 2\rho_{hi} \rho_{vi} D_{24}^i - 2\rho_{hi}^2 \rho_{vi} D_{25}^i + \rho_{vi}^2 D_{44}^i - 2\rho_{hi} \rho_{vi}^2 D_{45}^i + \rho_{hi}^2 \rho_{vi}^2 D_{55}^i) \quad (5)$$

Центральный элемент x_3 объекта X относится в класс ω_i , если $G_i(X) > G_j(X)$ для всех $j \neq i$ и $G_i(X) > T_i^s$, где T_i^s - пороговое значение для отклоненных векторов класса ω_i :

$$T_i^s = T_i^c - 0.5(L^{0.5} - 1)(N \ln((1 - \rho_{hi}^2)(1 - \rho_{vi}^2)) + \ln|B_i|) \quad (T_i^c \text{ берется из (4)}).$$

Анализ формул (2) и (5) показывает, что определяющий вклад в нахождение значений решающих функций вносит вычисление квадратичной формы (3): для каждого вектора ее необходимо вычислять $13m$ раз для перекрестия и $53m$ раз для квадрата (m - количество классов). Для объекта размером 5×5 количество вычислений квадратичной формы для каждого вектора составит $25m$ и $186m$ раз соответственно. В качестве примера в табл. 1 приведены временные характеристики (в секундах) классификации на компьютере IBM PC с процессором AMD Athlon XP 3200+ 10^7 векторов размерности $N = 5$.

Таблица 1

Количество классов	Перекрестие		Квадрат	
	3×3	5×5	3×3	5×5
5	80.7	147.4	311.6	1049.1
10	150.8	291	610.1	2120.4

Альтернативным вариантом, позволяющим существенно уменьшить временные затраты, является создание системы классификации, распределенной между персональным компьютером и многопроцессорной ЭВМ МВС-1000/М. На персональном компьютере выполняется обучение классификатора, которое, по сути, заключается в формировании на основе обучающих полей векторов средних и ковариационных матриц классов. Результаты обучения вместе с классифицируемым набором данных передаются по протоколу FTP на ЭВМ МВС-1000/М, где запускается параллельная версия соответствующей программы. Распараллеливание состоит в распределении набора данных между заданным количеством логических процессов, каждый из которых результаты своей работы записывает в отдельный файл. Эти файлы, в свою очередь, передаются на персональный компьютер, где производится их "склеивание" и дальнейшая интерпретация результатов выполненной классификации.

Другими алгоритмами, требующими параллельных вычислений, являются алгоритмы обнаружения линейных объектов (линементов) и кольцевых структур. При обработке ДДЗ с целью обнаружения линейных и кольцевых структур в силу целого ряда причин предпочтение отдается статистическому подходу [2]. Основная причина состоит в том, что вследствие случайного характера природных процессов данные дистанционных измерений содержат много случайных вариаций, маскирующих различия значений данных в точках области объекта и в точках области фона. Присутствие объектов проявляется в том, что случайные величины, наблюдаемые в точках области объекта, стохастически больше (или меньше) величин, наблюдаемых в точках области фона. В такой ситуации для обнаружения объектов эффективны так называемые непараметрические критерии, так как распределения статистик, используемых в этих критериях, не зависят от (неизвестных наблюдателю) распределений наблюдаемых величин. Предлагается следующая схема обнаружения объектов. Анализируются (почти) все возможные положения объектов, интересующих исследователя. Для каждого возможного положения решение о наличии объекта принимается по результату проверки, с помощью подходящего критерия, гипотезы однородности величин, наблюдаемых соответствующим образом.

Удобный непараметрический критерий для проверки этой гипотезы можно построить по значениям двух пар статистик Манна-Уитни [2]. Пусть ξ_i , ζ_i , ψ_i значения пикселей на нормали i ($i = 1, \dots, k$), такие, что значения ζ_i принадлежат точкам проверяемого положения объекта, а значения ξ_i и ψ_i - точкам, расположенным по разные стороны от этого положения. Статистики Манна-Уитни определяются в этом случае следующим образом:

$$\mu_{i1}^+ = I \{ \zeta_i > \xi_i \}, \quad \mu_{i1}^- = I \{ \zeta_i < \xi_i \}, \quad \mu_{i2}^+ = I \{ \zeta_i > \psi_i \}, \quad \mu_{i2}^- = I \{ \zeta_i < \psi_i \},$$

где $I\{\cdot\}$ - индикатор события $\{\cdot\}$, который равен 1 или 0 в зависимости от значения события $\{\cdot\}$.

Анализ значений пикселей изображения вдоль нормалей к предполагаемым положениям объектов требует больших временных затрат при его выполнении на однопроцессорной ЭВМ. В связи с этим разработаны параллельные реализации алгоритмов обнаружения линейных и кольцевых структур. Алгоритм обнаружения линейных структур допускает распараллеливание несколькими способами. Представляющийся наиболее естественным способ распараллеливания по направлениям линейных объектов оказывается достаточно громоздким в реализации, т.к. в зависимости от длины объектов количество направлений может

быть более сотни, поэтому реализован другой способ распараллеливания. Каждый процесс, запрашивая у системы общее количество процессов и свой номер, определяет горизонтальную полосу изображения для обработки. При этом соседние полосы перекрываются и глубина этого перекрытия зависит от длины обнаруживаемых объектов. Главный процесс передаёт необходимую для работы информацию остальным процессам, которые, в свою очередь, закончив обработку, сообщают результаты главному процессу. Аналогичный механизм распараллеливания реализован и для обнаружения кольцевых структур.

Работа выполнена при частичном финансировании по проекту РФФИ 05-07-90057.

Литература

1. Асмус В.В. Программно-аппаратный комплекс обработки спутниковых данных и его применение для задач гидрометеорологии и мониторинга природной среды. Диссертация в виде научного доклада на соискание ученой степени доктора физико-математических наук. На правах рукописи. Москва – 2002.
2. Салов Г.И. О мощностях непараметрических критериев для обнаружения протяженных объектов на случайном фоне. - Автометрия. 1997, №3, с.60-75.

PARALLEL COMPUTING IN EARTH REMOTE SENSING DATA PROCESSING

Asmus V.¹, Buchnev A.², Pyatkin V.²

¹ Scientific Research Center of Space Hydrometeorology “Planeta”,
ROSHYDROMET, Russia, Moscow

² Institute of Computational Mathematics and Mathematical Geophysics of SB RAS

The technology of analysis and processing of hyperspectral remote sensing data (RSD), founded on distributed and parallel aerospace data processing in local networks of supercomputer centers, is being developed in ICMMG SB RAS. On supercomputers RM600 and MVS-1000/M the system of parallel programming MPI, which allows portability of parallel processing program packages to the most known multiprocessor parallel computers, is used. In RSD processing the following steps, which require the parallel versions of corresponding algorithms because of large amount of computations, can be picked out.

The system of supervised classification (the classification with training). This system consists of seven classifiers (one classifier is per-element and other six – object classifiers), based on the use of strategy of maximum likelihood for the normal distributed vectors data, and two object classifiers based on the minimum distance. For the object we mean the block of adjacent vectors of square or crosswise form. As the physical characteristics of really scanned spatial objects, as a rule, larger than the resolution of survey systems, connections between data vectors take place. The usage of such information gives the opportunity to increase the accuracy of classification, if we try to recognize the group of adjacent vectors simultaneously – the object mentioned in the sense as above. Variant, which allows decrease greatly time expenses, is the creation the system of classification distributed between personal computer and multiprocessor MVS-1000/M. The training of the classifier is carried out on personal computer, that is the forming of estimates of vectors of means and covariance matrices of classes on the base of training fields. Results of the training together with the data set being classified are passed by protocol FTP to MVS-1000/M, where the parallel version of corresponding program is started. Paralleling is the distribution of data set between specified number of logical processes; each of them stores the results of its work into separate file. These files, in their turn, are passed to a personal computer, where they are “stuck together” and the results of received classification are further interpreted.

Other algorithms, which require parallel computations, are the algorithms of linear and circular objects detection. Nonparametric criteria are efficient for object detection as statistic distributions used in these criteria don't depend on distributions of observed values (unknown to observer). For each possible position the decision of the object presence is made according to the result of the test by means of appropriate criterion or the hypothesis of homogeneity of the values observed in a suitable way. Nonparametric criterion to test this hypothesis can be built according the values of two pairs of Mann–Witney statistics. The parallel versions of the algorithms of linear and circular objects detection are developed. Among different methods of paralleling the simple one in implementing was chosen: the original image is divided into parts, the number of which coincides with the number of logical processes run. Obviously, the maximal system output will be achieved at equal number of logical processes and physical processors.

The work was carried out partly with financial sponsoring of RFFR project 05-07-90057.