

# 求解大型稀疏线性方程组的 Krylov 子空间方法的发展

李晓爱<sup>1</sup>, 陈玉花<sup>2</sup>, 张耘<sup>2</sup>, 王新苹<sup>2</sup>

1. 河南师范大学数学与信息科学学院, 河南新乡 453007
2. 北京联合大学应用科技学院, 北京 102200

**摘要** 求解大型稀疏线性方程组是许多科学和工程计算中最重要的问题之一, Krylov 子空间方法是求解这类线性方程组的一个研究热点。本文介绍了 Krylov 子空间方法及其分类, 例如正交投影方法(或 Ritz-Galerkin 方法), 正交化方法(或极小残差方法), 双正交化方法(或 Petrov-Galerkin 方法), 解法方程组的 CGNE 和 CGNR 方法等, 指出了这些方法在算法设计方面国内外研究现状和存在问题, 着重考虑稀疏矩阵向量乘积与内积计算方法的并行处理问题; 讨论了预条件与并行预条件技术, 残差磨光技术及其并行实现, 数据的合理分布问题, 内积瓶颈问题等方面研究的发展趋势, 希望有更多学者了解和研究这些方法。

**关键词** 大型稀疏线性方程组; 迭代法; Krylov 子空间方法; 预条件技术

**中图分类号** O242, TP301.6

**文献标志码** A

**doi** 10.3981/j.issn.1000-7857.2013.11.010

## Development of Krylov Subspace Methods for Solving Large Sparse Linear System of Equations

LI Xiaoi<sup>1</sup>, CHEN Yuhua<sup>2</sup>, ZHANG Yun<sup>2</sup>, WANG Xinping<sup>2</sup>

1. College of Mathematics and Information Science, Henan Normal University, Xinxiang 453007, Henan Province, China
2. Applied School of Science and Technology, Beijing United University, Beijing 102200, China

**Abstract** Solving a large sparse linear system of equations is one of the most important problems in scientific and engineering computations. The Krylov subspace methods are widely used in this respect. This paper first reviews the Krylov subspace methods and their various types, such as, the orthogonal projection method (Ritz-Galerkin method), the orthogonalization method (or the minimal residual method), the bi-orthogonalization method (Petrov-Galerkin method), and the CGNE and CGNR methods for normal systems. The advantages and shortcomings of these methods are analyzed. Especially, we focus on the parallel computation of the sparse matrix-vector multiplication and the inner product. Then, this paper discusses the development of the preconditioning and the parallel preconditioning technique, the residual smoothing technology with its parallel implementation, the reasonable distribution of data, the bottleneck problem of the inner product.

**Keywords** large sparse linear system of equations; iterative methods; Krylov subspace method; preconditioning techniques

### 0 引言

大型稀疏线性方程组的求解是许多科学和工程计算中最重要也是最基本的问题之一。在许多重要领域, 如高维微分方程数值解、核物理与流体力学计算、结构与非结构问题的有限元分析、电力系统的优化设计、石油地震数据处理及

数值天气预报等, 都遇到许多大型和超大型的科学与工程计算和数据处理问题, 这些问题的求解往往归结为大型稀疏线性方程组的求解问题。另一方面, 当前大型科学计算技术已进入大规模并行计算时代, 面向并行计算环境研究大型稀疏线性方程组的高效并行算法显得尤为重要。

收稿日期: 2012-04-05; 修回日期: 2013-01-08

基金项目: 国家自然科学基金项目(11171094, 11171368); 河南省基础与前沿技术研究计划项目(132300410285)

作者简介: 李晓爱, 副教授, 研究方向为最优化理论及应用, 电子信箱: lxa.hnsd@163.com

近年来,随着软、硬件技术的发展,高性能并行计算机不断出现。目前为止,对称多处理机系统(SMP)可构成每秒几十亿次运算的系统,流水线向量多处理机系统(VPP)和 workstation 机群(NOWs)可构成每秒几百亿次的运算系统,而分布式存储大规模并行处理机系统(MPP)可构成每秒万亿次运算或更高的系统。高性能并行计算机为大型稀疏线性方程组的求解提供了计算工具,但如何充分发挥并行计算机的潜在性能并对大型稀疏线性方程组进行高效求解,关键是依据并行计算机的特点进行并行算法的研究和并行程序的设计与实现。

求解线性方程组

$$Ax=b \quad (1)$$

其中,  $A$  为  $n$  阶实(或复)系数矩阵,  $x$  和  $b$  是长度为  $n$  的向量。求解方法通常被分为两类:直接方法和迭代方法。对于大型稀疏线性方程组,由于存储量和计算量的限制,常使用迭代法求解。

从迭代法的发展看,20世纪50—70年代,由于电子计算机的发展,人们开始考虑和研究在计算机上用迭代法求线性方程组的近似解,并发展了许多非常有效的方法,如 Jacobi 方法、Gauss-Seidel 方法、SOR 方法、SSOR 方法等,及这些方法的改进和加速形式。Young<sup>[1]</sup>和 Varga<sup>[2]</sup>对迭代法进行了细致的描述与探讨,并已成为迭代法方面的经典著作。但由于科技的发展,越来越多的非结构化的、特殊的、大型的、稀疏的问题摆在了计算数学工作者面前,仍用这些方法显得收敛太慢(甚至不收敛),且这些方法的优劣往往依赖于最优参数的选取,除了一些特殊方程组有复杂的计算公式可求最优参数外,一般方程组无法直接得到最优参数,从而限制了方法的高效使用。目前这些方法已很少用于直接求解大型稀疏线性方程组,而是作为预条件子和其他方法(如 Krylov 子空间方法)结合使用。

20世纪70年代以来,人们将研究的重点转移到了对大型稀疏线性方程组的高效求解上来,提出了许多可行且有效的方法,特别是随着向量计算机和并行计算机的出现,大型稀疏线性方程组的并行求解已受到越来越多的关注。目前这方面的研究热点是 Krylov 子空间方法类及其并行实现、与 Krylov 子空间方法类有关的预条件及并行预条件技术以及结合多层方法、多重网格方法、区域分解法等提出的一些混合方法等。这类方法已在文献中层出不穷,且各种方法所提出的背景与处理方式各有不同,因此对迭代方法进行分类与综述是有必要的。

## 1 Krylov 子空间方法及其分类

设  $K_m$  为一个  $m$  维子空间,一般投影方法是从  $m$  维仿射子空间  $x_0+K_m$  中寻找近似解  $x_m$ ,使相应的残差  $r_m=b-x_m$  满足 Petrov-Galerkin 条件:

$$r_m \perp L_m \quad (2)$$

其中,  $L_m$  为另一个  $m$  维子空间。如果  $K_m=K_m(A, r_0)$ ,即 Krylov 子空间,则上述投影方法就称为 Krylov 子空间方法,其中  $r_0$

为初始残差,  $K_m(A, r_0)$  定义为

$$K_m(A, r_0) = \text{span}\{r_0, Ar_0, A^2r_0, \dots, A^{m-1}r_0\} \quad (3)$$

理想的解线性方程组(1)的 Krylov 子空间方法应具有以下特征:

- (1) 最优性,极小残差性或极小误差性(保证收敛速度快);
- (2) 计算的高效性,每步迭代的计算量少、存储量小。

依据  $L_m$  和  $K_m$  的不同选取,可得不同类型的 Krylov 子空间方法,其大致可分4类。

### 1.1 正交投影方法(或 Ritz-Galerkin 方法)

取  $L_m=K_m=K_m(A, r_0)$ ,这类 Krylov 子空间方法称为正交投影方法。其中最重要的方法是共轭梯度(CG)方法,此时要求  $A$  是对称正定的。CG方法1952年由 Hestenes 和 Stiefel<sup>[3]</sup>提出,在没有舍入误差的情况下,至多用  $n$  步即可求出线性方程组的精确解。但 CG 方法起初并未引起太多关注,因为它被认为是直接方法,直到1971年经过 Reid<sup>[4]</sup>的工作才使得 CG 方法被看作是迭代法,1976年 Concus 等<sup>[5]</sup>提出了预条件 CG 方法,并指出 CG 方法可有效地求解线性方程组,至此,CG 方法才得到越来越多的关注。CG 方法是一种理想的 Krylov 子空间方法,它使得误差的能量范数( $A$  范数)  $\|x^*-x_m\|_A$  在子空间  $x_0+K_m$  中极小,其中  $x^*$  为方程组(1)的精确解,且具有短迭代计算公式,保证了计算实现的高效性。目前与各种预条件子相结合的 CG 方法是求解大型稀疏对称正定线性方程组的主要方法。

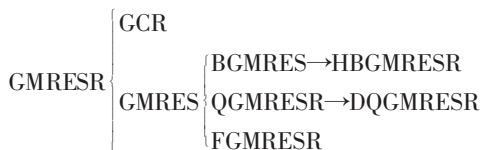
全正交化方法(FOM)及正交残差方法(ORTHORES)(不要求  $A$  是对称正定的)也是正交投影方法,也具有极小残差性,但不具有短迭代计算公式。

### 1.2 正交化方法(或极小残差方法)

取  $L_m=AK_m, K_m=K_m(A, r_0)$ ,这类方法具有残差极小性,即  $\|r_m\|_2 = \|b-Ax_m\|_2$  在子空间  $x_0+K_m$  中极小。GMRES 方法<sup>[6]</sup>是这类方法中的典型代表,由于它的高效性和良好的数值稳定性,近年来一直是研究的热点。但 GMRES 方法依据长迭代公式,在第  $m$  个迭代步需要  $m$  个内积及  $m$  个向量校正,因此每步迭代的计算量和存储量线性增长。近来,人们提出了许多它的有效变形,如采用 Householder 正交化的 GMRES 方法、拟 GMRES 方法、DQGMRES 方法、块 GMRES 方法、GMRES 类方法等。

共轭残差 CR(Conjugate Residual)方法也是一种正交化方法,GMRES 方法即是它的一种推广,其另一种推广是 GCR(Generalized CR)方法。1994年, Vorst 和 Vuik<sup>[7]</sup>提出了 GMRESR 方法,它涉及内积和外积两类方法,内积法即 GMRES 方法,外积法即 GCR 方法。1993年, Saad<sup>[8]</sup>提出了与此相类似的 FGMRES 方法,这种方法使得每步所用的预条件子可以不同,更具灵活性。1996年, Sturler<sup>[9]</sup>基于 GCR 方法给出了一大类方法,这类方法的特点是将 GCR 方法的正交性与其他算法,如 GMRES、BiCGSTAB 等方法结合产生新的收敛性好的方法。

属于正交化方法的主要 Krylov 子空间方法的相互关系可表示为



如 GMRES 方法,具有残差极小性,但随着迭代步数的增加,计算量和存储量线性增长,不具有计算实现的高效性且不易于并行化,因此常采用再启动 (restarted) 和截断 (truncated) 技术。在并行环境中,内积所需的通信可严重地降低这类方法的性能。

### 1.3 双正交化方法 (或 Petrov-Galerkin 方法)

取  $L_m=K_m(A^T, r_0), K_m=K_m(A, r_0)$ 。显然,若  $A$  为对称,则这类方法即是正交投影方法,因此这类方法常用于  $A$  是非对称的情形。对此情形,Faber 证明了一般不可能确定一个最优的  $x_m \in K_m(A, r_0)$  具有短递归计算公式。因此为了具有如 CG 中的短递归,可计算  $x_m \in K_m(A, r_0)$ ,使得  $b-Ax_m \perp K_m(A^T, r_0)$ ,导致 1952 年 Lanczos<sup>[10]</sup>提出的 BiCG (BiConjugate Gradient) 方法基于 Lanczos 双正交化算法。1975 年,Fletcher<sup>[11]</sup>对 BiCG 方法进行了改进。现在的许多双正交化方法均可由 BiCG 方法推出。但 BiCG 方法要使用  $A^T$  的运算,并在许多情形下可能出现中断现象,同时收敛性态不好。

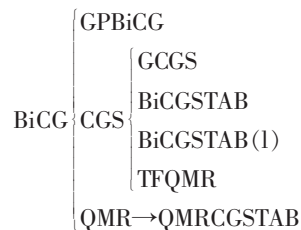
一个有趣的现象是 BiCG 方法中与  $A^T$  有关的运算可由与  $A$  本身的运算所代替,即  $\langle x, A^T y \rangle = \langle Ax, y \rangle$ 。由于 BiCG 中用  $A^T$  作用的函数仅在于维持残差被正交化的对偶空间,这个算子用  $A$  替代,扩充了 Krylov 子空间并寻找到更好的近似而每个迭代如 BiCG 具有相同的开销。1989 年,Sommeveld<sup>[12]</sup>提出的 CGS (Conjugate Gradient Squared) 方法避免了使用  $A^T$  的运算,无论方法收敛或发散,它比 BiCG 方法均快两倍。当用 BiCG 方法出现中断时,CGS 方法也中断。

CGS 方法是基于残差多项式平方的,在许多情况下会出现不规则的收敛性,为克服 CGS 方法和缺点,1992 年 Vorst 等<sup>[13]</sup>提出了 BiCGSTAB (BiConjugate Gradient STABlized) 方法,采用不同的方法来构造残差向量,其收敛性比 CGS 方法有所改善,同时不需要转置,但它仍未克服方法的中断性。

为改善 BiCG 方法的收敛性,1991 年 Freund 和 Nachtigal<sup>[14]</sup>提出了拟极小残差 QMR (Quasi-Minimal Residual) 方法,它具有局部残差极小性,收敛性态较好。QMR 方法采用 Lanczos 算法建立解空间的基,因此具有鲁棒性 (robustness),即不会产生中断现象,但它需要使用  $A^T$  的运算。为此,1993 年,Freund<sup>[15]</sup>又提出了 TFQMR (Trans-Free QMR) 方法。1994 年,Tong<sup>[16]</sup>提出了一类 QMR 方法,如 BQMR (2),BQMR (3) 方法,同时给出了一类不需要转置的 QMR 方法,如 QMRCSATB (2),QMRCSATB (4) 方法,这些方法具有更为光滑的收敛性,一般说来,具有更快的收敛速度。1996 年,Fokkema<sup>[17]</sup>等提出了 GCGS (Generalized CGS) 方法,将 CGS 和 BiCGSTAB 方法作为其特例并推出两种新的方法 CGS2 和位移 CGS (shifted CGS)

方法。1997 年,Zhang<sup>[18]</sup>提出了 GPBiCG (Generalized Product-type BiCG) 方法。

双正交化 Krylov 子空间方法还有很多,其主要方法之间的关系可表示为



双正交化方法的主要特点是采用短迭代计算公式,每步计算量和存储量不变,具有计算实现的高效性,但不具有极小残差性。QMR 类方法中要求具有某种近似极小性,从而保证了较好的收敛性。目前,如何得到具有更光滑收敛性的方法,即残差磨光技术,是一个研究热点,详见 1.4 节。BiCGSTAB (l) 方法为改进并行性能提供了某些可能。

### 1.4 与法方程组有关的方法

取  $L_m=K_m=K_m(A^T A, A^T r_0)$ ,这类方法主要是将 CG 方法应用到法方程组  $A^T A x = A^T b$  或  $AA^T u = b, x = A^T u$  上,由于  $A^T A$  的条件数比  $A$  的条件数更大,因此对大多数问题,它们无法与前面几类方法相竞争。1955 年,Craig 提出 CGNE (CG normal equations error minimizing) 方法,将 CG 方法应用到法方程组  $AA^T u = b, x = A^T u$  上。1987 年,Schönaauer<sup>[19]</sup>提出 CGNR (CG normal equations residual minimizing) 方法,将 CG 方法应用到法方程组  $A^T A x = A^T b$  上。这两种基于法方程组的方法利用短迭代公式,具有极小残差性,且从理论上讲总是收敛的,但由于  $A^T A$  和  $AA^T$  的特征值比  $A$  的特征值更分散,从而收敛速度慢,因而低效。近年来因为这类方法的低效性,对其研究相对较少,1997 年,Manneback<sup>[20]</sup>将 CGNR 方法用于求解一类不规则稀疏线性方程组,收到较好的效果。也有作者从理论上证明了 CGNE 和 CGNR 并非总是差的,对一些情况还是一种理想的迭代法。

根据问题类型的不同,可参考图 1 选择方法。

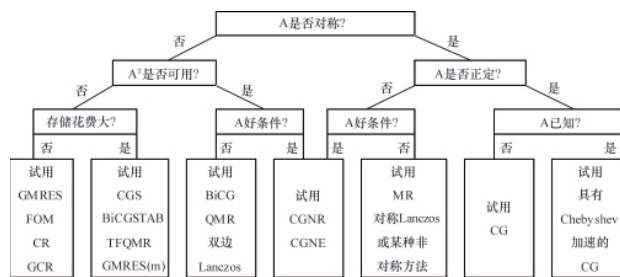


图 1 Krylov 子空间方法的选择  
Fig. 1 Choice of Krylov subspace methods

## 2 Krylov 子空间方法的并行方面

Krylov 子空间方法的主要计算核心是:

- (1) 稀疏矩阵向量乘积计算,即  $Ap$  和/或  $A^T p$  的计算,  $p$

由算法生成;

(2) 内积计算;

(3) 预条件步(若有的话),即若  $M$  为  $A$  的一个近似,则或者  $Mp$  的计算,或者  $Mz=p$  的计算;

(4) 向量校正。

Krylov 子空间方法的并行形式主要考虑的是以上 (1) 和 (2) 的并行处理。

## 2.1 稀疏矩阵向量乘积计算

关于稀疏矩阵向量乘积的并行处理已有许多文献论及,可以说目前仍是一个挑战,因为它通常没有标准的公开软件的支持,也没有统一的稀疏和并行 BLAS。起初,处理机的互连结构对程序员更透明的且重要的是使自己的代码适合特殊的互连结构,特别要避免在处理机间互连网内在长路径上发送数据。对向量多处理机情形,需要使用硬件的聚积-分散对矩阵向量乘积有效地向量化。对此,通常建议用 JAD (Jagged Diagonal) 格式存储稀疏系数矩阵,这种存储格式导致在共享存储向量计算机上有比较好的性能。如果可用多台处理机,则应采用 JAD 的块形式。在分布式环境中,为极小化通讯和提高负载平衡,矩阵向量乘积可用块行划分(可能结合一些数据的预处理)来并行化,还可考虑通信与计算重叠进行,矩阵在处理机上的分布尽量使得只有邻居-邻居的通信。

即使对串行执行,访存主存也是非常严重的问题。如果更多的数据可被保存在处理机的局部存储器(cache)中,则可能看到超线性加速,但对大多大型稀疏问题,所需向量不能整个地被保存在快速存储器(cache 或 registers)中,此时,对每个运算,向量必须从较慢的存储器内传到 cache 且由于没有足够的计算任务去做,从存储器最慢部分传送数据的速度变为主要成分,此时,建议重排迭代法中的运算,以使 3 层 BLAS 的操作成为可能。

还有许多处理稀疏矩阵向量乘积的方法,如对具有多个右端项的方程组,Krylov 子空间方法对不同的右端项生成不同的子空间,但可努力构造一个稍微大点的子空间来近似所有不同的,这种方法被称为块方法。虽然从计算复杂性的观点这不是十分有效,因为这样的子空间通常对任何一个各自方程组都不是最优的,但对并行处理和有限快速局部存储系统,许多计算可被局部组合,这导致更好地使用局部数据和更少的通信开销。

最近有个被称为 BLAS 技术论坛的活动,它试图对稀疏 BLAS 标准化并提供模型实现。它既包括稀疏三角矩阵的三角解法,也包括稀疏矩阵-向量乘积和稀疏矩阵-(稠密)矩阵乘积,现在的感觉是,只要适当编码,矩阵向量乘积在现代并行计算机上不会造成严重性能降低的通信问题,即使对相对小规模的问题也是这样。但内积计算并非如此。

## 2.2 内积计算

Krylov 子空间迭代法一般是对 Krylov 子空间生成一个正

交基,然后考虑给定式(1)关于这个基的限制。这些方法计算中集中的部分可如下表示。

设  $v_1, \dots, v_{m-1}$  为第  $m-1$  步处的正交基,则用修正的 Gram-Schmidt 构造  $v_m$  为

$$\begin{aligned} t &= Av_{m-1} \\ \text{for } j &= 1, 2, \dots, m-1 \text{ do} \\ t &= t - \langle t, v_j \rangle v_j \\ v_m &= t / \|t\|^2 \end{aligned} \quad (4)$$

在并行环境中,内积总是担当同步点并需要整体通信,在分布式存储计算机上,除了预条件外,它们成为一个主要的瓶颈。对于内积,既需要对约化操作也需要对聚集的内积的广播进行整体通信,因为所有处理机都需要知道结果。对一个有  $P(p \times p)$  台处理机的网络,这些通信花销与  $P$  成正比,这意味着当  $P$  充分大时通信花销将占主导地位,对大规模并行处理 MPP 系统这已成为获得高加速比的严重限制因素。例如 CG 方法每步只需两个内积,对每行只有 5 个元素的 90000P 阶矩阵,当处理机台数  $P > 400$  时,内积通信将占主导地位,因此解决内积通信问题在文献中受到极大的关注。

一些作者试图降低同步点数(并提高计算与存储开销的比值),Meurant<sup>[21]</sup>的方案是这些方案的代表:将两个分离的内积用 3 个连续的内积所代替,它们可被并行计算且通信可被组合,但代价是降低了数值稳定性。Bücher 和 Saurel<sup>[22]</sup>对 BiCG 和 QMR 方法提出了一个类似的方案。Demmel 等<sup>[23]</sup>提出了 CG 的另一个变形,其中有更多的可能重叠所有通信时间与有用的计算,其关键技巧是延迟解的校正一个迭代步,因为校正不必要等内积的完成,从而产生了重叠。

GMRES 方法是数值稳定的,但它每步内积的个数线性增加,常用再启动限制最大同步点数及通信开销。另一种方案是先计算一个合适但无须正交的基:

$$\{v_1, p_1(A)v_1, p_2(A)v_1, \dots, p_{m-1}(A)v_1\} \quad (5)$$

其中  $p_i$  为  $i$  次非退化多项式,且同时正交化  $\{p_i(A)v_1\}$ 。这提供了更多的并行性,一种产生这种基的合适的机制是用前一个 GEMRES( $m$ )循环的  $h_{ij}$ ,从而可用两层 BLAS 在这些基向量上执行 Gram-Schmidt 正交化,且所需通信可被大大重叠。显著的结果是,虽然 GMRES 中有大量内积,但可得到比 CG 方法(每步只需两个内积)更并行化的代码。

截止目前,内积计算的并行化、降低其通信开销仍是一个研究的热点与挑战。

## 3 发展趋势

### 3.1 预条件与并行预条件技术

在许多情形及应用中,直接使用迭代法不收敛或收敛非常慢,通常的补救方法是应用一个预条件或并行预条件技术改进,亦即代替解式(1)求解一个预条件系统。

寻求一个有效的预条件子(矩阵)的一般问题是找到一个  $M$  具有如下性质:

(1)  $M$  在某种意义上为  $A$  的一个好的近似。

(2) 构造  $M$  的花费是可以接受的(即花费不太大)。

(3) 系统  $My=z$  比原系统更容易求解。

这里有效性是指迭代法对预条件系统按 CPU 时间收敛得非常快。应用预条件的主要目的就是得到一个靠近单位阵的  $M^{-1}A$ , 对许多 Krylov 子空间方法, 期望  $M^{-1}A$  的条件数(远远)小于  $A$  的条件数, 或者  $M^{-1}A$  的特征值很强烈地聚集在某点(通常是 1)附近。应该认识到构造预条件子增加了迭代法的复杂性, 只有大量减少了迭代步, 使用预条件子才是值得的。

预条件子  $M$  的选取是多种多样的, 从可应用于一般矩阵的纯黑匣子代数技术到利用特殊问题类之特性的问题相关预条件子。现已发展出许多预条件与并行预条件技术: 如不完全分解技术, Meijerink 和 Vorst<sup>[24]</sup> 提出了不完全 LU(即 ILU) 分解技术。许多作者, 如 Axelsson<sup>[25]</sup>, Manteuffel<sup>[26]</sup>, Kershaw<sup>[27]</sup>, Saad<sup>[28]</sup> 等, 对 ILU 和 IC 分解进行了大量的研究和改进, 提出了许多更有效和并行化程度更高的 ILU 预条件子; 如改变计算次序技术; 重排未知量技术(如红黑排序); 多项式预条件技术; 稀疏近似逆技术(如 SPAI 方法、AINV 方法等)及块或区域预条件技术。

对一个给定的问题类, 通常难以构造一个有效的预条件子。若它们还必须是并行的则更难。起初为了并行性, 试图集中于对已存在的有效预条件子(如 ILU)进行最小的改变, 或利用矩阵  $A$  本身(如多项式预条件子), 这些努力有时已导出可量化的预条件子(例如超平面), 但对并行性而言, 它们的并行粒度太小。多项式预条件降低了内积和向量运算的比例, 但代价是迭代步数的增加。因此它通常对富含内积并在内积类运算比矩阵向量积运算相比更昂贵的平台上是有用的。最近注意力已经转向提供更粗粒度并行性的方法上: 区域分解法和稀疏近似逆方法。除了某些方面的成功, 稀疏近似逆技术仍处在刚起步时期。区域分解技术已被证明对大多与 PDEs 相关问题是成功的。

### 3.2 残差磨光技术及其并行实现

对双正交化方法, BiCG 用到  $A^T$  的运算, CGS 避免了  $A^T$  的运算, 但其收敛性和发散性同样突出, 造成残差范数的振荡很大, 解很不稳定。这种残差范数的不规则行为激励人们研究具有同样优点(不需要  $A^T$  的运算和短递归), 但可使残差范数具有更好的行为的方法, 如 BiCGSTAB 和 QMR 方法。

另一种产生好行为残差范数的方法最先由 Schönauer<sup>[29]</sup> 提出并由 Weiss<sup>[30]</sup> 进行了广泛的研究。Zhou 和 Walker<sup>[31]</sup> 将其称为迭代法的极小残差磨光 (Minimal Residual Smoothing, MRS) 技术, 并提出了可应用于任意迭代法的拟极小残差磨光 (QMRS) 技术, 指出对某些情形, MRS 的残差范数稍小于 QMRS 的残差范数且更稳定。Brezinski 等<sup>[32]</sup> 将其推广, 提出了解线性系统的混合过程, Cao 等<sup>[33]</sup> 对 QMRS 方法做了大量的理论研究, Zhang<sup>[34]</sup> 将 MRS 技术用于解偏微分方程的多重网格方法中, 以对多重网格方法进行加速。残差磨光技术是一个很有发展但研究刚起步的方法。

对此方面可进行如下研究:

(1) 对混合过程, 可以用多个因子组合多种迭代方法, 这样做本质是并行的, 它与并行多分裂方法有异曲同工之处, 即均是进行加权平均。借助于多分裂的思想可对此进行进一步的研究, 甚至提出更广泛与灵活的并行混合过程。当然, 目的是对残差进行磨光, 因此参数的选取方式也有待进一步研究。

对 MRS 及 QMRS 可类似考虑组合前面多个残差向量, 而得出一种多磨光因子的(拟)残差磨光方法。但这样做将导致长递归公式, 然而, 至少两个因子的三项递归乃至三因子的情形是可以考虑的。

(2) 更广泛的混合型过程。即用混合过程组合多于两种的不同方法, 首先组合前两种方法后, 可将所得结果与第 3 种方法相组合, 然后再将所得结果与第 4 种方法相组合, 如此等等。这样可得一种如楼梯状的混合过程。这种思想本质是串行的, 但其理论与实现均不困难。

(3) 二级混合过程以避免不规则性。Brezinski 曾指出当两种不同的方法相组合时, 混合过程可能是很不规则的。为避免这种情形, 可引入一个二级混合过程。即将混合过程所得的残差当作一个中间残差, 然后再用这个混合过程来磨光这个中间残差。也可结合(1)的研究提出并行二级混合过程。对 MRS 和 QMRS 亦可类似考虑。

(4) 方法的具体并行实现。在提出并研究并行混合过程之后, 必须研究其在高性能计算机上的具体实现。

### 3.3 数据的合理分布问题

目前, Cache 和多级存储与改进的编译器相结合的使用以最大程度地降低通信开销已成为关注的热点, 对此目前已有研究报告, 但尚不完全清楚。如何针对具体的并行机结构, 充分挖掘其潜力, 在算法设计中必须引起重视。尤其目前的趋势是用大的计算单位进行粗粒度的并行, 对稀疏矩阵计算如何合理分布数据以获得一个适当的负载平衡是分布式并行机必须解决的问题。这方面的解决需对并行机结构和改进的编译器有一定的了解, 需大量的并行编程经验, 可考虑如下操作以提高 Cache 命中率: 改变循环次序, 以使循环按列为主的形式进行, 以与矩阵的存储形式一致; 循环分块, 以提高数据的重复利用率; 循环展开、循环合并、改变循环次数(如矩阵加边等)和代码重写等。另外, 还应尽量保证连续的存储变量位于同一个线程, 并努力探索其他提高 Cache 命中率的方法。

### 3.4 内积瓶颈问题

在并行环境中, 内积总是担当同步点且其约化操作及广播都需要整体通信, 对 MMP 系统, 它是严重阻碍并行性的一个瓶颈问题。内积是 Krylov 子空间方法的主要计算成分之一, 它的任何并行改进均可得到一种新的并行 Krylov 子空间方法。以往的研究虽然很多, 但大多限于针对特殊的计算机平台或者将算法中几个不可并行计算的内积改为多个可并行计算的内积以减少同步点数。

对此,目前提出的最新思想是:将 Krylov 子空间视为一个多项式空间,用多项式空间中的一个离散内积代替欧氏内积,仅在离散内积中所用节点处计算这些多项式,在计算过程中,可期望使用一个正交多项式基,当节点集合为实时,这些正交多项式将满足一个三项递归,否则(即节点集合为复的),使用类似于 Arnoldi 过程的一个 Hessenberg 长递归。接着将这个新基正交化(可并行执行)且执行所得最小二乘问题的一个拟极小化求得一个拟极小残差解(如 QMR 那样),然后校正离散内积中所用节点集合(即  $A$  之特征值的近似值组成的集合),且重复以上过程。

这里对如何选择多项式内积并没有限制且可对方法做预条件处理。方法将标准的欧氏内积用多项式离散内积来代替,且新基的正交化可并行执行,这开辟了消除内积瓶颈问题的一个新的途径。

### 参考文献 (References)

- [1] Young D M. Iterative solution of large linear systems [M]. New York: Academic Press, 1971.
- [2] Varga R S. Matrix iterative analysis [M]. Englewood Cliffs, New Jersey: Prentice-Hall Inc, 1962.
- [3] Hesrenses M R, Stiefel E. Method of conjugate gradients for solving linear system [J]. Journal of Research of the National Bureau of Standards, 1952, 49(6): 409-436.
- [4] Reid J K. On the method of conjugate gradients for the solution of large sparse systems of linear equation [M]//Large Sparse Sets of Linear Equations. New York: Academic Press, 1971: 231-254.
- [5] Concus P, Golub G H, O'Leary D P. A generalized conjugate gradient method for the numerical solution of elliptic partial differential equations [M]//Bunch J R, Rose D J. Sparse Matrix Computations. New York: Academic Press, 1976: 309-332.
- [6] Saad Y. GMRES: A generalized minimal residual algorithm for solving nonsymmetric linear systems [J]. SIAM Journal on Scientific Statistical Computing, 1986, 7(3): 856-869.
- [7] van der Vorst H A, Vuik C. GMRESR: A family of nested GMRES methods [J]. Numerical Linear Algebra with Applications, 1994, 1(4): 369-386.
- [8] Saad Y. A flexible inner-outer preconditioned GMRES algorithm [J]. SIAM Journal on Scientific Computing, 1993, 14(2): 461-469.
- [9] de Sturler E. Nested Krylov methods based on GCR [J]. Journal of Computational and Applied Mathematics, 1996, 67(1): 15-41.
- [10] Lanczos C. Solution of systems of linear equations by minimized iteration [J]. Journal of Research of the National Bureau of Standards, 1952, 49(1): 33-53.
- [11] Fletcher R. Conjugate gradient methods for indefinite systems [C]//Proceedings of the Dundee Biennial Conference on Numerical Analysis 1974. New York: Springer Verlag, 1975: 73-88.
- [12] Sonneveld P. CGS, a fast Lanczos-type solver for nonsymmetric linear systems [J]. SIAM Journal on Scientific Statistical Computing, 1989, 10(1): 36-52.
- [13] de Sturler E. A performance model for Krylov subspace methods on mesh-based parallel computers [J]. Parallel Computing, 1996, 22(1): 57-74.
- [14] Freund R W, Nachtigal N M. QMR: A quasi-minimal residual method for non-Hermitian linear systems [J]. Numerische Mathematik, 1991, 60(1): 315-339.
- [15] Freund R W. A transpose-free quasi-minimal residual algorithm for non-Hermitian linear systems [J]. SIAM Journal on Scientific Computing, 1993, 14(2): 470-482.
- [16] Ton C H. A family of quasi-minimal residual methods for nonsymmetric linear systems [J]. SIAM Journal on Scientific Computing, 1994, 15(1): 89-105.
- [17] Fokkema D R, Sleijpen P, van der Vorst H A. Generalized conjugate gradient squared [J]. Journal of Computational and Applied Mathematics, 1996, 71(1): 125-146.
- [18] Zhang S L. GPBi-CG: Generalized product-type methods based on Bi-CG for solving nonsymmetric linear systems [J]. SIAM Journal on Scientific Computing, 1997, 18(2): 537-551.
- [19] van Duin V C N. Parallel sparse matrix computations [D]. Leiden: Leiden University, 1998.
- [20] Manneback P. Solving irregular sparse linear systems on a multicomputer using the CGNR method [J]. International Journal of High Performance Computing Applications, 1997, 11(3): 205-211.
- [21] Meurant G. The block preconditioned conjugate gradient method on vector computers [J]. BIT Numerical Mathematics, 1984, 24(4): 623-633.
- [22] Bücker H M, Sauren M. A parallel version of the unsymmetric Lanczos algorithm and its application to QMR [R]. Technical Report KFA-ZAM-IB-9605, Forschungszentrum Jülich GmbH, Jülich, Germany, 1996.
- [23] Demmel J W, Heath M T, van der Vorst H A. Parallel numerical linear algebra [M]. Cambridge University Press, Cambridge, 1993.
- [24] Meijerink J A, van der Vorst H A. An iterative solution method for linear systems of which the coefficient matrix is a symmetric M-matrix [J]. Mathematics of Computation, 1977, 31(137): 148-162.
- [25] Axelsson O. A general incomplete block-matrix factorization method [J]. Linear Algebra and its Applications, 1989, 74: 179-190.
- [26] Manteuffel T A. An incomplete factorization technique for positive definite linear systems [J]. Mathematics of Computation, 1980, 34(150): 473-497.
- [27] Kershaw D S. The incomplete Cholesky-conjugate gradient method for the iterative solution of systems of linear equations [J]. Journal of Computational Physics, 1978, 26: 43-65.
- [28] Saad Y. ILUM: A multi-elimination ILU preconditioner for general sparse matrices [J]. SIAM Journal on Scientific Computing, 1996, 17(4): 830-847.
- [29] Schönauer W. Scientific computing on vector computers [M]. New York: Elsevier Science Inc, 1987.
- [30] Weiss R. Convergence behavior of generalized conjugate gradient methods [D]. Karlsruhe: University of Karlsruhe, 1990.
- [31] Zhou L, Walker H F. Residual smoothing techniques for iterative methods [J]. SIAM Journal on Scientific Computing, 1994, 15(2): 297-312.
- [32] Brezinski C, Redivo-Zaglia M. Hybrid procedures for solving linear systems [J]. Numerische Mathematik, 1994, 67(1): 1-19.
- [33] Cao Z H. A note on convergence of quasi-minimal residual smoothing [J]. Applied Mathematics and Computation, 1999, 93: 289-297.
- [34] Zhang J. Multigrid acceleration techniques and applications to the numerical solution of partial differential equations [D]. Washington: The George Washington University, 1997.

(责任编辑 刘志远)