

ОГЛАВЛЕНИЕ

| | |
|--|----|
| Предисловие | 3 |
| Часть I. Релаксационные методы | 7 |
| Глава 1. Вводные сведения | 8 |
| 1.1. Основные обозначения и постановка задачи | 8 |
| 1.2. Метод Узавы — сопряженных градиентов | 10 |
| 1.3. Вспомогательные утверждения | 13 |
| 1.3.1. Две задачи на собственные значения | 14 |
| 1.3.2. Базис специального вида из собственных векторов | 15 |
| 1.3.3. Полезное начальное приближение | 17 |
| 1.4. Библиография и комментарии | 19 |
| Глава 2. Модифицированные методы релаксации. | |
| Общий анализ | 22 |
| 2.1. Сведения о методах релаксации | 22 |
| 2.1.1. Общие понятия | 22 |
| 2.1.2. Метод Якоби | 23 |
| 2.1.3. Метод SOR | 25 |
| 2.1.4. Метод SSOR | 26 |
| 2.2. Модифицированный метод Якоби (MJOR) | 28 |
| 2.2.1. Построение метода | 28 |
| 2.2.2. Спектр оператора перехода | 29 |
| 2.2.3. Условие сходимости | 31 |
| 2.2.4. Задача асимптотической оптимизации | 33 |
| 2.2.5. Оптимизация в подпространстве | 38 |
| 2.3. Модифицированный метод SOR (MSOR) | 40 |
| 2.3.1. Спектр оператора перехода | 40 |
| 2.3.2. Условие сходимости | 41 |
| 2.3.3. Задача асимптотической оптимизации | 43 |
| 2.4. Модифицированный метод SSOR (MSSOR) | 45 |
| 2.4.1. Спектр оператора перехода | 46 |

| | |
|--|-----------|
| 2.4.2. Условие сходимости | 49 |
| 2.4.3. Задача асимптотической оптимизации | 49 |
| 2.5. Трехпараметрический метод (3MSOR) | 50 |
| 2.5.1. Спектр оператора перехода | 50 |
| 2.5.2. Задача асимптотической оптимизации | 51 |
| 2.5.3. Частный случай: (β, τ) -метод | 52 |
| 2.5.4. Задача асимптотической оптимизации (β, τ) — метода | 54 |
| 2.6. Библиография и комментарии | 55 |
| Глава 3. Оценки погрешности методов MJOR и MSOR | 58 |
| 3.1. Оценки из общей теории | 59 |
| 3.1.1. Оптимальный одношаговый метод | 59 |
| 3.1.2. Циклический метод с чебышевскими параметрами | 59 |
| 3.1.3. Полуитерационный метод Чебышева | 60 |
| 3.1.4. Стационарный трехслойный метод | 61 |
| 3.1.5. Методы сопряженных направлений | 61 |
| 3.2. Погрешность метода MJOR в случае постоянных параметров | 63 |
| 3.2.1. Преобразование формул | 63 |
| 3.2.2. Начальное приближение | 64 |
| 3.2.3. Оценка погрешности с постоянными параметрами | 65 |
| 3.3. Погрешность метода MJOR в случае переменных параметров | 66 |
| 3.4. Погрешность метода MSOR в случае постоянных параметров | 71 |
| 3.4.1. Преобразование формул | 71 |
| 3.4.2. Начальное приближение | 73 |
| 3.4.3. Полином ошибки | 73 |
| 3.4.4. Оценка погрешности | 74 |
| 3.5. Погрешность метода MSOR в случае переменных параметров | 76 |
| 3.5.1. Преобразование формул | 76 |
| 3.5.2. Выбор параметров для p , как в циклическом методе | 78 |
| 3.5.3. Выбор параметров для p , как в трехслойных методах | 78 |
| 3.5.4. Выбор параметров для u , как в трехслойных методах | 80 |
| 3.6. Библиография и комментарии | 84 |

| | |
|--|-----|
| Глава 4. Релаксационные методы для системы с параметром | 86 |
| 4.1. Явный метод типа MSOR (MSORe) | 86 |
| 4.1.1. Построение метода | 86 |
| 4.1.2. Спектр оператора перехода | 87 |
| 4.1.3. Условие сходимости | 88 |
| 4.1.4. Задача асимптотической оптимизации | 90 |
| 4.2. Неявный метод типа MSOR (IMSORe) | 91 |
| 4.2.1. Построение метода | 91 |
| 4.2.2. Спектр оператора перехода | 91 |
| 4.2.3. Условие сходимости | 93 |
| 4.2.4. Задача асимптотической оптимизации | 95 |
| 4.3. Погрешность метода MSORe в случае постоянных параметров | 96 |
| 4.3.1. Преобразование формул | 96 |
| 4.3.2. Начальное приближение | 98 |
| 4.3.3. Полином ошибки | 98 |
| 4.3.4. Оценка погрешности | 99 |
| 4.4. Погрешность метода MSORe в случае переменных параметров | 100 |
| 4.4.1. Преобразование формул | 100 |
| 4.4.2. Выбор параметров для p , как в циклическом методе | 102 |
| 4.4.3. Выбор параметров для p , как в трехслойных методах | 103 |
| 4.4.4. Выбор параметров для u , как в трехслойных методах | 105 |
| 4.5. Библиография и комментарии | 108 |
| Глава 5. Методы для нормальных уравнений | 109 |
| 5.1. Оптимизация метода для базовой системы | 110 |
| 5.1.1. Спектр оператора равносильной задачи | 110 |
| 5.1.2. Минимизация числа обусловленности | 111 |
| 5.1.3. Наилучшая оценка погрешности | 112 |
| 5.2. Оптимизация метода для системы с параметром | 113 |
| 5.2.1. Спектр оператора равносильной задачи | 113 |
| 5.2.2. Минимизация числа обусловленности | 115 |
| 5.2.3. Наилучшая оценка погрешности | 118 |
| 5.3. Оптимизация метода для случая равносильной системы | 118 |
| 5.3.1. Спектр оператора равносильной задачи | 119 |
| 5.3.2. Минимизация числа обусловленности | 121 |

| | |
|---|-----|
| 5.3.3. Наилучшая оценка погрешности | 124 |
| 5.4. Библиография и комментарии | 124 |
| Часть II. Обобщенные методы | 125 |
| Глава 6. Предварительные результаты | 126 |
| 6.1. Классы оптимизации | 126 |
| 6.2. Что происходит с алгоритмом Узавы? | 128 |
| 6.3. Нерегулярные задачи с седловой точкой. | 131 |
| 6.4. Вспомогательные утверждения | 134 |
| 6.4.1. Сведения из анализа | 134 |
| 6.4.2. Спектральные свойства одного пучка операторов | 140 |
| 6.4.3. Свойства некоторых классов функций | 149 |
| 6.5. Ключевые этапы построения и анализа алгоритмов . | 156 |
| 6.6. Библиография и комментарии | 161 |
| Глава 7. Блочно треугольное предобусловливание (GMSOR) | 163 |
| 7.1. Формулировка метода и его свойства. | 163 |
| 7.2. Симметричная регулярная задача: оптимизация в классе \mathbb{K}_1 | 165 |
| 7.3. Симметричная регулярная задача: оптимизация в классе \mathbb{K}_2 | 171 |
| 7.4. Симметричная нерегулярная задача: оценка в классе \mathbb{K}_{2s} | 178 |
| 7.5. Несимметричная регулярная задача: оценка в классе \mathbb{K}_3 | 180 |
| 7.6. Библиография и комментарии | 184 |
| Глава 8. Блочно диагональное предобусловливание | 187 |
| 8.1. Обобщенный метод MJOR (GMJOR) | 187 |
| 8.1.1. Формулировка метода | 187 |
| 8.1.2. Связь спектров операторов перехода GMJOR и GMSOR | 188 |
| 8.1.3. Оптимизация метода в классах $\mathbb{K}_1, \mathbb{K}_2, \mathbb{K}_3, \mathbb{K}_{2s}$ | 189 |
| 8.1.4. Случай $\beta = 0$ | 190 |
| 8.2. Обобщенный метод Ланцоша (GMLan) | 194 |
| 8.2.1. Построение метода | 194 |
| 8.2.2. Оптимизация в классе \mathbb{K}_1 | 196 |
| 8.2.3. Оптимизация в классе \mathbb{K}_2 | 202 |
| 8.2.4. Оценка в классе \mathbb{K}_{2s} | 210 |
| 8.3. Библиография и комментарии | 211 |

| | |
|---|-----|
| Глава 9. Симметризация специального вида | 214 |
| 9.1. Обобщенный метод Bramble–Pasciak (GMBP) | 215 |
| 9.1.1. Построение метода | 215 |
| 9.1.2. Оптимизация метода в классе \mathbb{K}_1 | 218 |
| 9.1.3. Оптимизация метода в классе \mathbb{K}_2 | 221 |
| 9.1.4. Оценка в классе \mathbb{K}_{2s} | 224 |
| 9.2. Библиография и комментарии | 225 |
| Глава 10. Модельные седловые операторы | 228 |
| 10.1. Неконструктивный подход | 228 |
| 10.1.1. Построение методов | 228 |
| 10.1.2. Оценка в классе \mathbb{K}_1 | 230 |
| 10.1.3. Оценка в классе \mathbb{K}_2 | 232 |
| 10.1.4. Оценка в классе \mathbb{K}_{2s} | 235 |
| 10.2. Конструктивное предобусловливание | 236 |
| 10.2.1. Построение методов | 236 |
| 10.2.2. Оценка в классе \mathbb{K}_2 | 236 |
| 10.2.3. Оценка в классе \mathbb{K}_{2s} | 241 |
| 10.3. Библиография и комментарии | 242 |
| Глава 11. Методы попеременных симметричных и кососимметричных итераций | 244 |
| 11.1. Стационарный метод (GPHSSI) | 245 |
| 11.1.1. Формулировка метода | 245 |
| 11.1.2. Безусловная сходимость метода | 245 |
| 11.1.3. Оптимизация в классе \mathbb{K}_2 | 247 |
| 11.2. HSS-предобусловливание | 252 |
| 11.2.1. Построение методов | 252 |
| 11.2.2. Оценка спектра в классе \mathbb{K}_2 | 254 |
| 11.2.3. Анализ сходимости методов чебышевского типа и GMRES | 256 |
| 11.3. Библиография и комментарии | 263 |
| Глава 12. Нелинейные задачи и блочно треугольное предобусловливание | 265 |
| 12.1. Уравнения с кососимметричным возмущением | 265 |
| 12.1.1. Постановка задачи | 265 |
| 12.1.2. Оценка скорости сходимости | 268 |
| 12.2. Уравнения с сильно монотонным оператором | 277 |
| 12.2.1. Постановка задачи | 277 |
| 12.2.2. Вспомогательные факты и утверждения | 279 |
| 12.2.3. Оценка скорости сходимости | 281 |
| 12.3. Библиография и комментарии | 284 |

| | |
|---|-----|
| Часть III. Приложение к гидродинамике | 286 |
| Глава 13. Inf-sup неравенство и смежные вопросы | 289 |
| 13.1. О задаче Стокса и спектре Коссера | 290 |
| 13.2. Неравенства Фридрихса и Корна в двухмерном случае | 292 |
| 13.3. Точные значения и оценки снизу константы Ладыженской | 293 |
| 13.4. Анизотропия области | 297 |
| 13.5. Угловые точки на границе | 299 |
| 13.6. Разное | 303 |
| 13.6.1. Обобщенная задача Стокса | 303 |
| 13.6.2. Уравнения Ламе в теории упругости и слабо-сжимаемая жидкость | 305 |
| 13.6.3. Смешанный подход для эллиптических уравнений | 305 |
| 13.6.4. Другие применения | 307 |
| Глава 14. Численный анализ LBB-условия | 308 |
| 14.1. Задача с гладким решением | 308 |
| 14.2. Задача с негладким решением | 313 |
| 14.2.1. Схема 1 | 313 |
| 14.2.2. Схема 2 | 314 |
| 14.2.3. Вычислительные аспекты | 315 |
| 14.2.4. Расчеты для квадратной области | 316 |
| 14.2.5. Расчеты для прямоугольной области | 319 |
| Глава 15. Численный анализ роли оператора $\beta BC^{-1}B^T$ в сходимости GMSOR | 322 |
| 15.1. Алгоритм численной оптимизации | 323 |
| 15.2. Задача о квадратной каверне | 324 |
| 15.3. Обтекание тела в трубе | 326 |
| Список литературы | 329 |

ПРЕДИСЛОВИЕ

Не будет преувеличением сказать, что за последние 15–20 лет одной из самых значительных и популярных тем в численном анализе является исследование сеточных седловых задач, имеющих выраженную блочную структуру. Это объясняется двумя основными факторами: широтой приложений и новизной идей, что необходимо порождает разнообразие конкретных формулировок проблем и методов их решения.

Термины «седловая задача» или «оператор с седловой точкой» имеют происхождение из теории математического программирования ([1], с. 602). Применительно к системам линейных алгебраических уравнений это, как правило, означает симметрию матрицы и наличие у нее собственных значений разных знаков, хотя допустима и более широкая трактовка. Понятие «сеточная задача» относится к происхождению постановки, часто следующей из дискретизации дифференциальных уравнений в частных производных. Блочная структура по форме характеризует специфику задачи, а по сути делает постановку доступной для анализа.

Важным сигналом, что настало время провести определенную систематизацию результатов в этой области, стала публикация обзора Benzi M., Golub G., Liesen J. Numerical solution of saddle point problems (2005) [117], содержащего более пятисот ссылок на первоисточники. Конечно, эту работу следует воспринимать только в качестве путеводителя по тематике, так как изложение материала с необходимой полнотой потребовало бы не полторы сотни страниц, а во много раз больше.

Настоящая книга имеет целью осветить круг проблем вычислительной математики, традиционно считающихся одними из самых сложных не только при решении седловых задач. Здесь имеется в виду оптимизация итерационных методов для решения линейных систем, т. е. процесс наилучшего в некотором смысле выбора ограниченного множества скалярных параметров.

Теория итерационных методов для линейных алгебраических уравнений сама по себе достаточно обширна. Видимо, в настоящее время библиография работ в этой области уже не поддается учету. Поэтому отметим ее приоритетное направление — решение сеточных задач, отличительными особенностями матриц которых являются: большая размерность, плохая обусловленность и разреженность (ленточная структура). Достаточно полное представление о развитии и современном состоянии теории итерационных методов для решения таких задач можно получить из следующих монографий:

- Young D.M. *Iterative Solution of Large Linear Systems* (1971) [200];
- Самарский А. А., Николаев Е. С. *Методы решения сеточных уравнений* (1978) [76];
- Марчук Г. И., Лебедев В. И. *Численные методы в теории переноса нейтронов* (1981) [64];
- Hackbusch W. *Multi-Grid Methods and Applications* (1985) [159];
- Дьяконов Е. Г. *Минимизация вычислительной работы. Асимптотически оптимальные алгоритмы для эллиптических задач* (1989) [41];
- Saad Y. *Iterative methods for sparse linear systems* (1996) [183].

Решение седловых систем в них практически не затронуто. Поэтому предлагаемая книга является связующим звеном между ставшей уже классической теорией итерационных методов и необходимостью решать актуальные седловые задачи. Она написана на основе многочисленных журнальных публикаций авторов, а также с учетом сжатия и качественной переработки материала из монографии [96].

При решении линейных систем алгебраических уравнений ключевым понятием является предобусловливание, что связано с подбором некоторых матриц, удобных с вычислительной точки зрения. Поэтому к настоящему времени для изложения итерационной теории сложилась следующая структура. Сначала рассматриваются простейшие методы: простой итерации, Якоби, Зейделя, верхней релаксации (последние три традиционно называют релаксационными) и т. д., при этом объектом анализа являются условия сходимости и задача асимптотической оптимизации, т. е. выбор постоянных итерационных параметров, обеспечивающих наивысшую скорость сходимости. Кроме того, представляет интерес получение оценок погрешности метода для этого оптимального случая. Отличительной особенностью релаксационных методов является применение элементов матрицы исходной системы для построения предобусловливающих операторов (простейшее предобусловливание). Это приводит

к выявлению предельных, или асимптотически наилучших, характеристик методов. На следующем этапе изучаются возможности ускорения сходимости алгоритмов за счет использования переменных итерационных параметров, причем их выбор может осуществляться как по явным формулам (когда известными считаются постоянные из матричных неравенств), так и из вариационных принципов. И наконец, делается завершающий шаг, связанный с введением спектрально-эквивалентных операторов (нетривиальным предобуславливанием), т. е. производится обобщение наиболее важных полученных результатов на случай, когда матрицы в алгоритмах обращаются неточно.

Книга состоит из двух частей и приложения. Первая часть посвящена изложению теории релаксационных методов для решения седловых задач. Используя элементарные средства анализа, здесь удается показать базовые идеи и основные результаты, а также проследить преимущество с классической теорией. Этот материал вполне доступен для понимания молодыми исследователями, уверенно владеющими основами линейной алгебры, математического анализа и численных методов. Во второй части книги речь идет об обобщенных (спектрально-эквивалентных) методах, что требует более основательной подготовки от читателя. Здесь для цельности изложения исходная постановка переформулирована с необходимой общностью и с достаточной полнотой приведены вспомогательные утверждения. В этой части подвергнуты анализу итерационные методы, основанные на всех основополагающих идеях блочного седлового случая. В отличие от симметричной знакоопределенной сеточной задачи (типа дискретного аналога уравнения Пуассона) анализ обобщенных методов решения седловых задач не является формальной процедурой. Более того, введение спектрально-эквивалентных операторов приводит к такому расширению множества постановок оптимизационных задач, что влечет за собой необходимость разработки новой методологии исследования, включающей нестандартные технические элементы. Важно отметить, что главными качествами найденных решений задач по оптимизации методов являются неулучшаемость оценок и наличие явных формул для итерационных параметров.

Интерес к проблемам оптимизации алгоритмов для седловых задач возник у авторов при моделировании гидродинамических эффектов, в первую очередь, из необходимости численного решения уравнений Навье—Стокса. Именно поэтому приложение книги ориентировано на вычислительную гидродинамику. В рассматриваемом случае оценка некоторых величин в формулах для оптимальных

параметров тесно связана с определением констант в дискретном условии Ладыженской—Бабушки—Бреци и непрерывном $\inf\text{-sup}$ -неравенстве. Проблематика наилучших констант в функциональных неравенствах очень глубока и интересна сама по себе; нашей целью было зафиксировать содержательное родство между далекими на первый взгляд областями исследований.

В книге принята сквозная нумерация глав во всех частях, аналогично нумеруются формулы внутри главы. Имеется дополнительное разбиение на параграфы, в соответствии с ним организована нумерация утверждений. Для обозначения методов используются унифицированные англоязычные аббревиатуры, что способствует удобству в восприятии и обозначает преемственность в развитии теории.

Авторы глубоко признательны коллегам и учителям, общение с которыми позволило сформировать точку зрения на предмет; выделим среди них Е. Г. Дьяконова, В. И. Лебедева и Г. М. Кобелькова. Исключительная благодарность адресуется Н. С. Бахвалову, чья роль в обсуждении замысла и советах по написанию книги невозможно переоценить.

ЧАСТЬ I

**РЕЛАКСАЦИОННЫЕ
МЕТОДЫ**

ВВОДНЫЕ СВЕДЕНИЯ

В главе приводятся основные обозначения и постановка задачи, рассматривается базовый для решения блочных седловых задач метод Узавы и доказываются вспомогательные результаты, используемые в дальнейшем.

1.1. ОСНОВНЫЕ ОБОЗНАЧЕНИЯ И ПОСТАНОВКА ЗАДАЧИ

Обозначим через U и P евклидовы пространства векторов размерностей N_u и N_p соответственно, $Z = U \times P$. Запись вектора z в виде $z = \{u, p\} \in Z$ означает, что он состоит из двух компонент: $u \in U$, $p \in P$. Для системы линейных алгебраических уравнений $Lz = F$ это соответствует разбиению квадратной матрицы L на блоки L_{ij} ($1 \leq i, j \leq 2$):

$$L = \begin{pmatrix} L_{11} & L_{12} \\ L_{21} & L_{22} \end{pmatrix},$$

размерности которых определяются размерностями компонент вектора z : L_{11} является $N_u \times N_u$ матрицей, $L_{12} - N_u \times N_p$, $L_{21} - N_p \times N_u$, $L_{22} - N_p \times N_p$. Правая часть системы F имеет представление, аналогичное z : $F = \{f, \varphi\} \in Z$.

Если L_{11} невырождена, то можно определить матрицу

$$S = -L/L_{11} \equiv -(L_{22} - L_{21}L_{11}^{-1}L_{12}),$$

часто называемую дополнением Шура для матрицы L относительно L_{11} (для удобства взятым со знаком минус). Значимость S определяется следующим факторизованным представлением L :

$$L = \begin{pmatrix} I & 0 \\ L_{21}L_{11}^{-1} & I \end{pmatrix} \begin{pmatrix} L_{11} & 0 \\ 0 & -S \end{pmatrix} \begin{pmatrix} I & L_{11}^{-1}L_{12} \\ 0 & I \end{pmatrix}, \quad (1.1)$$

где I имеет смысл единичной матрицы соответствующего размера. Выражение (1.1) формально сводит решение системы $Lz = F$ к обращению двух подматриц L_{11} и S , а фактически — только S , так как в определении последней уже входит L_{11}^{-1} . Таким образом, этот прием

позволяет понизить размерность решаемой задачи путем сведения ее к равносильной системе с матрицей специальной структуры.

Далее мы будем иметь дело с вещественной системой линейных алгебраических уравнений $L_\varepsilon z = F$ с параметром $\varepsilon \geq 0$ наиболее распространенного вида:

$$L_\varepsilon z \equiv \begin{pmatrix} A & B \\ B^T & -\varepsilon D \end{pmatrix} \begin{pmatrix} u \\ p \end{pmatrix} = \begin{pmatrix} f \\ \varphi \end{pmatrix} \equiv F, \quad (1.2)$$

где $A = A^T > 0$, $D = D^T > 0$ — квадратные матрицы размеров $N_u \times N_u$ и $N_p \times N_p$, а B — прямоугольная, в общем случае, матрица размера $N_u \times N_p$. С другими постановками блочных седловых задач можно ознакомиться в [117].

Будем предполагать, что матрица L_ε невырождена при любом $\varepsilon \geq 0$. Это условие, в силу факторизации вида (1.1), означает невырожденность дополнения Шура

$$S_\varepsilon = B^T A^{-1} B + \varepsilon D.$$

По построению матрица S_0 симметрична и положительно полуопределена, т. е. $(S_0 p, p) \geq 0$ для произвольного $p \in P$. Поэтому при $\varepsilon > 0$ матрица S_ε (и следовательно, L_ε) невырождена всегда, а при $\varepsilon = 0$ условия невырожденности матриц L_0 и S_0 носят эквивалентный характер. Это означает исключительность ситуации с $\varepsilon = 0$, поэтому основное изложение будет посвящено именно этому случаю, а обобщение результатов для $\varepsilon > 0$ будет проводиться по мере необходимости.

Обратим внимание на соотношение между размерностями пространств U и P , следующее из условия невырожденности матрицы L_0 . Это — неравенство $N_u \geq N_p$. Действительно, рассмотрим в противном случае решение однородной системы $L_0 z = 0$ вида $z = \{0, p\}$, которое порождает для компоненты p условие $Bp = 0$, т. е. систему из N_u уравнений с $N_p > N_u$ неизвестными. Такая однородная система всегда имеет нетривиальное решение. Поэтому вместо неконструктивного предположения о невырожденности L_0 часто требуют, чтобы матрица B имела полный ранг при указанном выше условии на размерности.

Задачу (1.2) часто называют задачей с седловой точкой (или задачей с седловым оператором). В этом случае имеют в виду, что функция Лагранжа

$$\begin{aligned} \Phi(u, p) &\equiv (L_\varepsilon z, z) - 2(F, z) = \\ &= (Au, u) + 2(B^T u, p) - \varepsilon(Dp, p) - 2(f, u) - 2(\varphi, p) \end{aligned}$$

имеет седловую точку $z^* = (u^*, p^*)$, совпадающую с решением (1.2), т. е. справедливы равенства

$$\Phi(u^*, p^*) = \min_{u \in U} \Phi(u, p^*) = \max_{p \in P} \Phi(u^*, p).$$

Содержательными по смыслу следствиями этого факта являются следующие свойства симметричной матрицы L_ε : фиксированная блочная структура и принципиальное отсутствие знакоопределенности.

1.2. МЕТОД УЗАВЫ – СОПРЯЖЕННЫХ ГРАДИЕНТОВ

Рассмотрим для задачи $L_0 z = F$ в покомпонентной форме

$$\begin{cases} Au + Bp = f, \\ B^T u = \varphi \end{cases}$$

процедуру исключения Гаусса. Выделим компоненту решения u из первого уравнения

$$u = A^{-1}(f - Bp) \quad (1.3)$$

с последующей подстановкой во второе. В результате получим систему уравнений $S_0 p = b$ следующего вида:

$$S_0 p \equiv B^T A^{-1} B p = B^T A^{-1} f - \varphi \equiv b. \quad (1.4)$$

Отсюда имеем, что если матрица L_0 невырождена, то для нахождения $z = \{u, p\}$ достаточно сначала решить систему (1.4) с симметричной положительно определенной матрицей S_0 , а затем определить недостающую компоненту u по формуле (1.3). Под методом Узавы в самом широком смысле, как правило, понимают реализацию этого подхода.

Для решения (1.4) в качестве предобусловливателя матрицы S_0 обычно вводится матрица $C = C^T > 0$ размерности $N_p \times N_p$ (в простейшем случае $C = I$ — единичная матрица) и затем применяется обобщенный метод сопряженных градиентов. Приведем одну из его наиболее распространенных версий.

Зададим начальный вектор p^0 и последовательно вычислим векторы:

$$r^0 = S_0 p^0 - b, \quad w^0 = C^{-1} r^0, \quad s^0 = w^0,$$

затем с помощью величины

$$a_1 = \frac{(w^0, r^0)}{(S_0 s^0, s^0)}$$

определим следующее приближение

$$p^1 = p^0 - a_1 s^0.$$

Далее для $k = 1, 2, \dots$ формулы имеют вид:

$$\begin{aligned} r^k &= r^{k-1} - a_k S_0 s^{k-1}, & w^k &= C^{-1} r^k, \\ d_k &= \frac{(w^k, r^k)}{(w^{k-1}, r^{k-1})}, & s^k &= w^k + d_k s^{k-1}, \\ a_{k+1} &= \frac{(w^k, r^k)}{(S_0 s^k, s^k)}, & p^{k+1} &= p^k - a_{k+1} s^k. \end{aligned}$$

Хорошо известно (см., например, [15], с. 296), что при отсутствии ошибок округлений метод сопряженных градиентов приводит к точному решению за число итераций, не превышающее размерность системы. Однако на практике критерием останова часто служит абсолютная или относительная малость невязки

$$r^k = S_0 p^k - b$$

в какой-либо норме.

Для наших целей важной является оценка погрешности (ошибки) метода сопряженных градиентов. Обозначим через γ и Γ постоянные в матричном неравенстве (здесь и далее будем считать их точными)

$$\gamma C \leq S_0 \leq \Gamma C, \quad 0 < \gamma, \quad (1.5)$$

что эквивалентно принадлежности спектра матрицы

$$C^{-\frac{1}{2}} S_0 C^{-\frac{1}{2}}$$

отрезку $[\gamma, \Gamma]$. Тогда для любой итерации с номером k будет справедливо неравенство

$$\|p^k - p\|_{S_0} \leq \frac{2q_0^k}{1 + q_0^{2k}} \|p^0 - p\|_{S_0}, \quad q_0 = \frac{1 - \sqrt{\xi}}{1 + \sqrt{\xi}}, \quad \xi = \frac{\gamma}{\Gamma}. \quad (1.6)$$

Здесь и далее выражение $\|r\|_A$ будет использоваться для обозначения нормы вектора r , порожденной симметричной положительной матрицей A , в данном случае

$$\|p\|_{S_0} = (S_0 p, p)^{1/2}.$$

Приведенная оценка означает, что норма ошибки $p^k - p$ асимптотически убывает как геометрическая прогрессия со знаменателем q_0 . Отметим также, что в процессе реализации алгоритма Узавы в качестве промежуточных величин получаются приближения к u по формуле

$$u^{k+1} = A^{-1}(f - Bp^k),$$

для которых также справедливо неравенство вида (1.6)

$$\|u^{k+1} - u\|_A \leq \frac{2q_0^k}{1 + q_0^{2k}} \|u^1 - u\|_A.$$

При фиксированных γ и Γ величина

$$\frac{2q_0^k}{1 + q_0^{2k}}$$

является неувлучшаемой, поскольку определяется нормой так называемого оптимального полинома (в данном случае — чебышевского). Этот факт лежит в основе распространенной точки зрения:

Если матрица A исходной задачи (1.2) легко обратима, то метод Узаваы — сопряженных градиентов является наилучшим для ее решения.

Конечно, так дело обстоит не всегда (в комментариях к главе приведен пример задачи, «неудобной» для этого алгоритма), однако отрицать приоритетность метода Узаваы — сопряженных градиентов для решения очень широкого класса седловых задач было бы неправильно. Вследствие этого актуальной является разработка преобусловливателей C для дополнения Шура S_0 (или S_ε) в конкретных задачах, сводящихся к системам с седловой точкой (1.2). Использование специфики постановок, следующей, например, из исходных дифференциальных уравнений или геометрии областей, делает эту тематику практически неисчерпаемой.

Имеется еще одно важное следствие оценки (1.6). Поскольку в ней явно присутствует зависимость от γ и Γ (констант спектральной эквивалентности матриц в (1.5)), то для представления об эффективности любого другого итерационного метода решения системы (1.2) желательно иметь оценку его погрешности, выраженную в тех же величинах. Это приводит к принципиально новым постановкам задач для оптимизации итерационных методов в отличие от классической теории.

Закончим раздел указанием на границу применимости алгоритма Узаваы — сопряженных градиентов. Метод существенно теряет свою привлекательность, если требуются большие вычислительные затраты для обращения матрицы A (или, что эквивалентно, отсутствует хороший преобусловливатель для нее). В такой ситуации даже использование внутренних итераций для приближенного вычисления величины $A^{-1}x$ делает алгоритм малоэффективным (см. также раздел 6.2).