

**МОСКОВСКИЙ ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ  
имени М.В. ЛОМОНОСОВА  
Факультет вычислительной математики и кибернетики**

**Н.М. Андрушевский**

**Анализ устойчивости решений  
систем линейных алгебраических  
уравнений**

**Методическое пособие  
специального вычислительного практикума**

**Москва**

---

**2008**

**УДК 378(075.8):512.64: 514.12**

**ББК 22.143: 22.151.573**

**A66**

*Печатается по решению Редакционно-издательского совета*

*факультета вычислительной математики и кибернетики*

*Московского государственного университета имени М.В. Ломоносова*

**Рецензенты:**

Х.Д. Икрамов, д.ф.-м.н., профессор ВМиК МГУ;

М.В.Уфимцев, к.ф.-м.н., с.н.с. ВМиК МГУ

**Андрушевский Н.М.**

**A66 Анализ устойчивости решений систем линейных алгебраических уравнений:** Учебное пособие. – М.: Издательский отдел факультета ВМиК МГУ имени М.В. Ломоносова (лицензия ИД N 05899 от 24.09.2001 г.); МАКС Пресс, 2008. – 71 с.

ISBN 978-5-89407-329-3

ISBN 978-5-317-02592-2

Учебное пособие посвящено исследованию устойчивости решений систем линейных алгебраических уравнений. На основе использования метода сингулярного разложения матриц предложены методика классификации матриц относительно меры их чувствительности к погрешностям в исходных данных и алгоритмы устойчивого решения плохо-обусловленных систем уравнений произвольного порядка и ранга. Предложенный математический аппарат может применяться в задачах определения по методу наименьших квадратов оптимальных параметров и их радиусов доверительной области для линейных математических моделей, а также численного решения интегральных уравнений типа свертки. Изложенные в пособии теоретические сведения и задачи практикума могут быть рекомендованы студентам старших курсов, специализирующимся в области математического моделирования.

**ISBN 978-5-89407-329-3**

© Факультет вычислительной математики

**ISBN 978-5-317-02592-2**

и кибернетики МГУ имени М.В. Ломоносова, 2008

© Андрушевский Н.М., 2008

Учебное издание

АНДРУШЕВСКИЙ Николай Матвеевич

Анализ устойчивости решений  
систем линейных алгебраических уравнений

Методическое пособие  
специального вычислительного практикума

Издательский отдел  
Факультета вычислительной математики и кибернетики  
МГУ имени М.В. Ломоносова  
Лицензия ИД N 05899 о 24.09.01 г.

119992, ГСП-2, Москва, Ленинские горы, МГУ им. М.В. Ломоносова,  
2-й учебные корпус

Напечатано с готового оригинал-макета  
в издательстве ООО <<МАКС Пресс>>  
Лицензия ИД N 00510 от 01.12.99 г.  
Подписано к печати 12.11.2008  
Формат 60x90 1/16. Усл.печ.л. 3,75. Тираж 100 экз. Заказ 659.

119992, ГСП-2, Москва, Ленинские горы, МГУ им. М.В. Ломоносова,  
2-й учебный корпус, 627 к.  
Тел. 939- 3890. Тел./Факс 939-3891

## Введение

*Всякая лиса свой хвост хвалит.  
Пословица*

Задачи численного решения систем линейных алгебраических уравнений являются наиболее распространенными. Это связано с тем, что задачи определения оптимальных параметров математических моделей, статистической обработки данных, анализа изображений и сигналов, факторного анализа, экономики и другие при адекватном математическом описании сводятся к решению соответствующих систем линейных алгебраических уравнений. Кроме того, численное решение многих задач математической физики и математического моделирования сложных физических или химических процессов является неиссякаемым источником задач линейной алгебры.

В современных системах программирования научно-технических расчетов, таких как **MATLAB**, **MATCAD**, **MATHEMATICA** и других, имеются обширные библиотеки программ для решения разнообразных задач линейной алгебры. В библиотеках программ, подготовленных программистами-профессионалами, реализованы самые лучшие к настоящему времени алгоритмы, в основе которых лежат строгие математические обоснования. В целом этот программный продукт удовлетворяет самым высоким требованиям качества и освобождает пользователей от необходимости его усовершенствования. Однако это не освобождает пользователей от необходимости учитывать степень влияния погрешностей в исходных данных на результат получаемых решений даже в том случае, когда задачи линейной алгебры с неточными исходными данными решаются абсолютно точно. Суть в том, что влияние погрешностей в исходных данных имеет специфический *анизотропный* характер и, зачастую, даже незначительные погрешности случайного характера могут привести к тому, что полученные решения и истинные будут очень сильно отличаться. Благодушное забвение или игнорирование этих опасностей служит основой ложной интерпретации изучаемого явления или даже открытия научных псевдоэффектов.

В качестве иллюстрации рассмотрим простейший пример решения системы двух линейных уравнений:

$$\begin{cases} 3x_1 + 5x_2 = y_1 \\ 2x_1 + 4x_2 = y_2 \end{cases}.$$

Пусть в правой части находятся компоненты вектора  $\mathbf{y} = \begin{pmatrix} 8 \\ 6 \end{pmatrix}$ . Тогда

вектор  $\mathbf{x}_0 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$  – точное решение рассматриваемой системы. Для

вектора  $\mathbf{y}_1 = \begin{pmatrix} 8.16 \\ 6.12 \end{pmatrix}$  – точное решение  $\mathbf{x}_1 = \begin{pmatrix} 1.02 \\ 1.02 \end{pmatrix}$ , а для вектора

$\mathbf{y}_2 = \begin{pmatrix} 7.94 \\ 6.08 \end{pmatrix}$  – точное решение  $\mathbf{x}_2 = \begin{pmatrix} 0.68 \\ 1.18 \end{pmatrix}$ . Сопоставляя эти

результаты, получаем  $\|\mathbf{y}_0 - \mathbf{y}_2\| < \|\mathbf{y}_0 - \mathbf{y}_1\|$ ,  $\|\mathbf{x}_0 - \mathbf{x}_2\| \gg \|\mathbf{x}_0 - \mathbf{x}_1\|$ . Таким образом в рассматриваемом примере возмущение решения, вызванное меньшими погрешностями в правой части, значительно больше возмущения решения, вызванного большими погрешностями. Этот, на первый взгляд, парадоксальный факт обусловлен тем, что рассматриваемая *невырожденная* система относится к классу *плохо-обусловленных* систем. Именно для таких систем характерно появление больших возмущений в решениях при наличии малых возмущений в исходных данных. Поэтому при решении систем линейных алгебраических уравнений с погрешностями в исходных данных обязательно необходимо идентифицировать матрицы на предмет их хорошей или плохой обусловленности, определять количественную меру обусловленности и вычислять радиус доверительной области возмущенных решений в зависимости от меры чувствительности матриц к погрешностям и уровня погрешностей в исходных данных. Следует заметить, что используемое в данном пособии понятие плохой обусловленности системы отличается от того, которое используют специалисты по вычислительной линейной алгебре.

Суть этого различия в следующем. В практических приложениях часто приходится решать переопределенные системы линейных алгебраических уравнений  $\mathbf{A}(\mathbf{x} + \delta\mathbf{x}) = \mathbf{y} + \delta\mathbf{y}$  с погрешностями  $\delta\mathbf{y}$ . Для исследователей-практиков важное значение имеет коэффициент наибольшего возможного усиления *абсолютной* погрешности

$q = \frac{\|\delta \mathbf{x}\|}{\|\delta \mathbf{y}\|}$ , когда заданы матрица  $\mathbf{A}$  и уровень погрешности в правой части  $\|\delta \mathbf{y}\|$ . Естественно считать, что если  $q > 1$ , то система уравнений сильно чувствительна к погрешностям (*плохо обусловлена*), а если  $q \leq 1$ , то система уравнений слабо чувствительна (*хорошо обусловлена*). Поскольку  $\|\delta \mathbf{x}\| \leq \|\mathbf{A}^+\| \|\delta \mathbf{y}\| = (1/\sigma_{\min}) \|\delta \mathbf{y}_r\|$ , где  $\mathbf{A}^+$  – обобщенная (псевдообратная) матрица,  $0 < \sigma_{\min}$  – наименьшее сингулярное число в сингулярном разложении матрицы  $\mathbf{A}$ ,  $\|\delta \mathbf{y}_r\|$  – длина проекции вектора  $\delta \mathbf{y}$  на ранговое пространство матрицы  $\mathbf{A}$ , то индикатором плохой обусловленности служит величина  $0 < \sigma_{\min} < 1$ . Именно в этом смысле используется термин “*плохая обусловленность*”. В то же время специалисты по вычислительной линейной алгебре используют понятие плохой обусловленности при вычислении коэффициента максимального усиления *относительной* погрешности  $c = \text{cond}(\mathbf{A}) = \frac{\|\delta \mathbf{x}\|}{\|\mathbf{x}\|} / \frac{\|\delta \mathbf{y}\|}{\|\mathbf{y}\|}$ . Это число играет важную роль, в частности, для задачи обращения матриц. Однако, поскольку  $\text{cond}(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^+\| = \sigma_{\max} / \sigma_{\min} \geq 1$ , а  $\frac{\|\delta \mathbf{x}\|}{\|\mathbf{x}\|} \leq \text{cond}(\mathbf{A}) \frac{\|\delta \mathbf{y}\|}{\|\mathbf{y}\|}$ , то из последнего неравенства невозможно точно определить максимально возможное расстояние между приближенным и истинным решениями, т.к. точные значения  $\|\mathbf{x}\|$  и  $\|\mathbf{y}\|$  априори неизвестны. В то же время, величина  $\|\delta \mathbf{y}\|$  зачастую может быть надежно оценена из условий экспериментальных измерений, а так как  $\|\delta \mathbf{y}_r\| = \sqrt{\|\delta \mathbf{y}\|^2 - \|\delta \mathbf{y}_0\|^2}$ , где  $\|\delta \mathbf{y}_0\|$  – вычисляемая по методу наименьших квадратов длина проекции вектора  $\delta \mathbf{y}$  на ортогональное дополнение к ранговому пространству матрицы, то вычисление  $\|\delta \mathbf{x}\|$  является вполне реалистичной задачей. Именно поэтому мы акцентируем внимание на исследовании абсолютных погрешностей приближенных решений. В принципе не должно быть особой путаницы в использовании термина *плохо-обусловленная* система, если подчеркивать о каких погрешностях идет речь – абсолютных или относительных.

Данное пособие посвящено изложению этих вопросов применительно к ситуации, когда погрешности в исходных данных содержатся только в правой части системы, а матрица коэффициентов задана без погрешностей.

Пособие состоит из 3 частей. В первой части пособия подробно рассматривается метод *сингулярного разложения матриц* как наиболее эффективный метод анализа чувствительности к погрешностям матриц неполного ранга или близких к вырождению, а также его применение для приближенного решения плохо обусловленных систем уравнений большого порядка. Следует подчеркнуть, что традиционные курсы линейной алгебры, читаемые в высших учебных заведениях, мало затрагивают тему разнообразного применения сингулярного разложения матриц. На самом деле метод сингулярного разложения матриц является своеобразным *“томографом высокого разрешения”*. Этот метод позволяет рассмотреть структуру и спектральные характеристики матриц произвольного размера и ранга в мельчайших подробностях, определить характер и меру анизотропных деформаций при линейных отображениях конечномерных пространств, определить тип и меру обусловленности матриц, вычислить обобщенную обратную матрицу и оптимальное решение несовместных систем уравнений с матрицами неполного ранга. На основе спектральной классификации матриц можно сформулировать критерии слабой или сильной чувствительности систем алгебраических уравнений к погрешностям, вычислить меру обусловленности и радиус доверительной области (максимально возможное расстояние между приближенным и истинным решениями).

Для полноценного восприятия читателем очерченного круга вопросов в пособии излагаются формулировки ряда теорем вместе с их доказательствами. Сюда относятся: основная теорема линейной алгебры, теорема об основных свойствах матриц проектирования, теорема о сингулярном разложении матриц, теорема об анизотропном характере отображений конечномерных пространств матрицами произвольного размера и ранга. Доказательство других упоминаемых теорем можно найти в общедоступных учебниках линейной алгебры.

Теоретический материал первой части используется во второй части для построения алгоритмов определения оптимальных параметров и радиусов их доверительных областей для ряда

популярных линейных математических моделей с учетом обусловленности систем уравнений соответствующей задачи МНК, а также алгоритмы численного решения интегральных уравнений типа свертки.

В заключительной части приводится список заданий специального вычислительного практикума по указанной тематике. Индивидуальное выполнение заданий практикума предназначено для лучшего усвоения теоретического материала и подготовит заинтересованных читателей к осознанной работе с плохо-обусловленными системами линейных уравнений.

Изложенные в пособии теоретические сведения и задачи практикума могут быть рекомендованы студентам старших курсов и аспирантам, специализирующимся в области математического моделирования.

Автор выражает искреннюю благодарность коллегам – профессору Щедрину Б.М., профессору Икрамову Х.Д. и старшему научному сотруднику Уфимцеву М.В. – за доброжелательную критику и ценные замечания, способствующие улучшению изложения материала книги. Разумеется, это не снимает с автора ответственности за возможные упущения и недоработки. Критические замечания по поводу содержания и формы изложения теоретического материала пособия просим присылать по адресу [nandrush@cs.msu.su](mailto:nandrush@cs.msu.su)



## Часть 1

### §1. Основная теорема линейной алгебры

*Я видел Азии бесплодные пределы,  
Кавказа дальний край, долины обгорелы,  
Жилище дикое черкесских табунов,  
Подкумка знойный брег, пустынные вершины,  
Обвитые венцом летучих облаков,  
И закубанские равнины! ...*

*Пушкин А.С. Стихотворение*

Любая прямоугольная матрица вещественных чисел, размера  $n \times k$ ,

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1k} \\ a_{21} & a_{22} & \dots & a_{2k} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nk} \end{pmatrix}$$

может рассматриваться как математическая запись оператора, который осуществляет отображение абстрактного конечномерного вещественного пространства  $\mathbf{R}^k$  в другое конечномерное вещественное пространство  $\mathbf{R}^n$  следующим образом:

$$\begin{aligned} \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \cdot \\ \cdot \\ y_n \end{pmatrix} &= \begin{pmatrix} a_{11}x_1 + a_{12}x_2 + \dots + a_{1k}x_k \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2k}x_k \\ \cdot \\ \cdot \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nk}x_k \end{pmatrix} = \\ &= x_1 \begin{pmatrix} a_{11} \\ a_{21} \\ \cdot \\ \cdot \\ a_{n1} \end{pmatrix} + x_2 \begin{pmatrix} a_{12} \\ a_{22} \\ \cdot \\ \cdot \\ a_{n2} \end{pmatrix} + \dots + x_k \begin{pmatrix} a_{1k} \\ a_{2k} \\ \cdot \\ \cdot \\ a_{nk} \end{pmatrix}. \end{aligned} \quad (1)$$

Более кратко выражение (1) можно записать в матрично-векторном виде

$$\mathbf{R}^k \Rightarrow \mathbf{R}^n : \mathbf{y} = \mathbf{A}\mathbf{x} = x_1\mathbf{a}_1 + x_2\mathbf{a}_1 + \dots + x_k\mathbf{a}_k, \quad (2)$$

где  $\mathbf{x} \in \mathbf{R}^k$ ,  $\mathbf{y} \in \mathbf{R}^n$ ,  $\mathbf{a}_i \in \mathbf{R}^n, i = 1, 2, \dots, k$  – вектор-столбцы матрицы  $\mathbf{A}$ . (Далее везде мы закрепляем символ  $\mathbf{x}$  для векторов, принадлежащих “левому” пространству  $\mathbf{R}^k = \mathbf{X}$  и символ  $\mathbf{y}$  для векторов, принадлежащих “правому” пространству  $\mathbf{R}^n = \mathbf{Y}$ ). Представляя линейное преобразование (1) с помощью матрицы  $\mathbf{A}$  мы предполагаем, что как в  $\mathbf{X}$  так и в  $\mathbf{Y}$  заданы ортогональные системы координат.

Из формулы (2) следует, что каждый вектор  $\mathbf{x}$  однозначно отображается в некоторый вектор  $\mathbf{y}$ , который представляет собой линейную комбинацию векторов  $\mathbf{a}_i \in \mathbf{R}^n, i = 1, 2, \dots, k$ , где весовыми множителями этой комбинации служат координаты вектора  $\mathbf{x}$ . Аналогично, транспонированная матрица  $\mathbf{A}^T$ , размера  $k \times n$ , определяет оператор отображения в противоположном направлении конечномерного пространства  $\mathbf{R}^n \Rightarrow \mathbf{R}^k$ , т.е.

$$\mathbf{x} = \mathbf{A}^T \mathbf{y} = y_1\mathbf{b}_1 + y_2\mathbf{b}_2 + \dots + y_n\mathbf{b}_n, \quad (3)$$

где  $\mathbf{b}_j, j = 1, 2, \dots, n$  – вектор-столбцы матрицы  $\mathbf{A}^T$  (т.е. вектор-строки матрицы  $\mathbf{A}$ ).

Длиной вектора  $\mathbf{z} \in \mathbf{R}^l$  называется величина  $\|\mathbf{z}\| = \sqrt{z_1^2 + z_2^2 + \dots + z_l^2}$ . Так определяемую длину  $\|\mathbf{z}\|$  называют еще *евклидовой нормой* вектора.

Величина  $(\mathbf{y}, \mathbf{z}) = \mathbf{y}^T \mathbf{z} = \sum_{i=1}^n y_i z_i$  – *скалярное произведение* векторов.

Косинус угла  $\varphi$  между векторами  $\mathbf{y}$  и  $\mathbf{z}$  вычисляется по формуле  $\cos \varphi = (\mathbf{y}, \mathbf{z}) / (\|\mathbf{y}\| \|\mathbf{z}\|)$ . Векторы  $\mathbf{y}$  и  $\mathbf{z}$  *ортогональны*, когда  $(\mathbf{y}, \mathbf{z}) = 0$ .

Любые два вектора одного и того же пространства удовлетворяют *неравенству Коши – Шварца – Буняковского*

$$|(\mathbf{y}, \mathbf{z})| \leq \|\mathbf{y}\| \|\mathbf{z}\|.$$

К числу важнейших понятий линейной алгебры относится, в первую очередь, свойство *линейной зависимости* векторов.

**Определение 1.** Если ни одна нетривиальная линейная комбинация векторов  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k \in \mathbf{R}^n$  не равна нуль-вектору  $\theta_n = (0, 0, \dots, 0)^T$ ,  $\theta_n \in \mathbf{R}^n$ , т.е.

$$x_1 \mathbf{a}_1 + x_2 \mathbf{a}_2 + \dots + x_k \mathbf{a}_k = \theta_n \Rightarrow x_1 = x_2 = \dots = x_k = 0, \quad (4)$$

то эти векторы называются *линейно независимыми*. В противном случае эти векторы называются *линейно зависимыми*, и один из них является линейной комбинацией остальных.

Между линейной независимостью и ортогональностью существует простая связь.

**Теорема 1.** Если ненулевые векторы  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k$  попарно ортогональны (каждый вектор ортогонален любому другому), то они линейно независимы.

Совокупность линейно независимых векторов  $\mathbf{y}_i, i=1, 2, \dots, l$ , образует базис подпространства  $\mathbf{R}^l \subset \mathbf{R}^n$ , размерность которого равна  $l$ .

*Произвольный набор линейно независимых векторов может быть преобразован в набор попарно ортогональных векторов при помощи алгоритма Грама – Шмидта.*

Максимальное число  $r$  линейно независимых вектор-столбцов матрицы  $\mathbf{A}$  называется ее *рангом*. Векторы  $\mathbf{y} = \mathbf{A}\mathbf{x}$ ,  $\mathbf{x} \in \mathbf{R}^k$ , образуют множество значений матрицы  $\mathbf{A}$ , и это множество называется ее *образом* в пространстве  $\mathbf{R}^n$ . Поскольку в матрице  $\mathbf{A}$  имеются  $r \leq k$  линейно независимых вектор-столбцов, то ее образом в пространстве  $\mathbf{R}^n$  будет подпространство, размерность которого также равна  $r$ . Это подпространство называют *ранговым пространством* (пространством столбцов) матрицы  $\mathbf{A}$  и обозначают символом  $\mathbf{R}(\mathbf{A})$ ,  $\mathbf{R}(\mathbf{A}) \subset \mathbf{R}^n$ .

В матрично-векторном виде соотношение (4) можно записать в виде  $\mathbf{Ax} = \theta_n$ , т.е. векторы  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k$  линейно зависимы тогда и только тогда, когда существует нетривиальное решение системы

$$\mathbf{Ax} = \theta_n \quad (5)$$

и, следовательно, ранг  $r < k$ . Множество векторов  $\mathbf{x} \in \mathbf{R}^k$ , удовлетворяющих системе уравнений (5), образует так называемое *нуль-пространство*  $\mathbf{N}(\mathbf{A})$  матрицы  $\mathbf{A}$ . Размерность нуль-пространства  $\dim \mathbf{N}(\mathbf{A}) = k - r$ . Нуль-пространство матрицы называется также *ядром* этой матрицы, а его размерность называется *дефектом* матрицы. Обозначим дефект матрицы символом  $\delta(\mathbf{A})$ .

Так как ранг транспонированной матрицы  $\mathbf{A}^T$  также равен  $r$ , то ее *ранговое пространство*  $\mathbf{R}(\mathbf{A}^T)$ ,  $\mathbf{R}(\mathbf{A}^T) \subset \mathbf{R}^k$ , также имеет размерность  $r$ . Это подпространство называют еще *пространством строк* матрицы  $\mathbf{A}$ . Соответственно, *нуль-пространство*  $\mathbf{N}(\mathbf{A}^T)$  имеет размерность  $n - r$ .

Справедлива одна из наиболее важных теорем линейной алгебры:

**Теорема 2.**  $\dim \mathbf{R}(\mathbf{A}) = \dim \mathbf{R}(\mathbf{A}^T) = r,$   
 $\delta(\mathbf{A}) = k - r, \delta(\mathbf{A}^T) = n - r.$  (6)

Для произведения  $\mathbf{AB}$  прямоугольных матриц  $\mathbf{A}$ , размера  $k \times l$ , и  $\mathbf{B}$ , размера  $l \times n$ , – справедлива следующая теорема.

**Теорема 3.**  $rank(\mathbf{AB}) \leq rank(\mathbf{A}), rank(\mathbf{AB}) \leq r(\mathbf{B}), \delta(\mathbf{AB}) \geq \delta(\mathbf{B}).$

В частности, если столбцы матрицы  $\mathbf{A}$  линейно независимы, то для квадратной симметричной матрицы  $\mathbf{G} = \mathbf{A}^T \mathbf{A}$  имеем

$$rank(\mathbf{G}) = rank(\mathbf{A}) = k, \quad (7)$$

и матрица  $\mathbf{A}^T \mathbf{A}$  является *обратимой* матрицей.

Перейдем к следующему важному понятию *ортогональности* двух подпространств.

**Определение 2.** Два подпространства  $\mathbf{U}$  и  $\mathbf{V}$  одного и того же пространства  $\mathbf{R}^n$  называются *ортогональными*, если каждый вектор  $\mathbf{u} \in \mathbf{U}$  ортогонален каждому вектору  $\mathbf{v} \in \mathbf{V}$ , т.е. скалярное произведение  $(\mathbf{u}, \mathbf{v}) = 0$  для всех  $\mathbf{u} \in \mathbf{U}$  и  $\mathbf{v} \in \mathbf{V}$ . (Факт

ортогональности подпространств будем обозначать символом  $U \perp V$ ).

**Пример 1.** Пусть  $U$  является плоскостью, порожденной векторами  $\mathbf{u}_1 = (1, 0, 0, 0)^T$  и  $\mathbf{u}_2 = (1, 2, 0, 0)^T$ , а  $V$  – прямая, порожденная вектором  $\mathbf{v}_1 = (0, 0, 3, 4)^T$ . Тогда, так как вектор  $\mathbf{v}_1$  ортогонален как к вектору  $\mathbf{u}_1$ , так и к вектору  $\mathbf{u}_2$ , прямая  $V$  будет ортогональна к плоскости  $U$ , т.е.  $V \perp U$ .

*Прямой суммой*  $W = U \oplus V$  подпространств  $U$  и  $V$  называется подпространство  $W = \{\mathbf{w} : \mathbf{w} = \mathbf{u} + \mathbf{v}, \mathbf{u} \in U, \mathbf{v} \in V\}$ .

**Определение 3.** Пусть подпространство  $U \subset \mathbf{R}^n$ . Тогда подпространство всех  $n$ -мерных векторов, ортогональных к подпространству  $U$ , называется *ортогональным дополнением* к  $U$  и обозначается  $U^\perp$ ,  $U \oplus U^\perp = \mathbf{R}^n$ .

Для любого подпространства  $U \in \mathbf{R}^n$  существует ортогональное дополнение  $U^\perp \in \mathbf{R}^n$  и если  $\mathbf{y} \in \mathbf{R}^n$ , то найдутся единственные векторы  $\mathbf{u} \in U$  и  $\mathbf{u}_\perp \in U^\perp$  такие, что  $\mathbf{y} = \mathbf{u} + \mathbf{u}_\perp$ . Для этих векторов выполняется теорема Пифагора:

$$\|\mathbf{y}\|^2 = \|\mathbf{u}\|^2 + \|\mathbf{u}_\perp\|^2.$$

**Теорема 4.** Для любой матрицы  $A$  размера  $n \times k$  ранговое пространство  $\mathbf{R}(A^T)$  и нуль-пространство  $\mathbf{N}(A)$  являются ортогональными подпространствами в  $\mathbf{R}^k$  и выполняется условие

$$\dim \mathbf{R}(A^T) + \dim \mathbf{N}(A) = r + (k - r) = k,$$

т.е.  $\mathbf{R}(A^T)$  – *ортогональное дополнение* к  $\mathbf{N}(A)$ . Аналогично ранговое пространство  $\mathbf{R}(A)$  и нуль-пространство  $\mathbf{N}(A^T)$  являются ортогональными подпространствами в  $\mathbf{R}^n$ :

$$\dim \mathbf{R}(A) + \dim \mathbf{N}(A^T) = r + (n - r) = n,$$

$\mathbf{R}(A)$  – *ортогональное дополнение* к  $\mathbf{N}(A^T)$ .

Доказательство. Предположим, что  $\mathbf{u}$  – произвольный вектор из нуль-пространства  $\mathbf{N}(A)$ . Тогда  $A\mathbf{u} = \theta_n$ , т.е. вектор  $\mathbf{u}$  ортогонален каждой вектор-строке матрицы  $A$  или, что то же самое, к каждому

вектор-столбцу матрицы  $\mathbf{A}^T$ . Следовательно, он ортогонален ко всему пространству  $\mathbf{R}(\mathbf{A}^T)$ , порожденному этими столбцами. Это верно для любого вектора  $\mathbf{u}$  из нуль-пространства  $\mathbf{N}(\mathbf{A})$ , и, следовательно,  $\mathbf{R}(\mathbf{A}^T) \perp \mathbf{N}(\mathbf{A})$ . Проводя аналогичные рассуждения для матрицы  $\mathbf{A}^T$ , получаем  $\mathbf{R}(\mathbf{A}) \perp \mathbf{N}(\mathbf{A}^T)$ .

Так как  $\mathbf{R}(\mathbf{A}^T) \oplus \mathbf{N}(\mathbf{A}) = \mathbf{R}^k$ , то любой вектор  $\mathbf{x} \in \mathbf{R}^k$  можно однозначно разложить в сумму  $\mathbf{x} = \mathbf{x}_r + \mathbf{x}_0$ , где  $\mathbf{x}_r \in \mathbf{R}(\mathbf{A}^T)$ ,  $\mathbf{x}_0 \in \mathbf{N}(\mathbf{A})$ .

Так как  $\mathbf{R}(\mathbf{A}) \oplus \mathbf{N}(\mathbf{A}^T) = \mathbf{R}^n$ , то любой вектор  $\mathbf{y} \in \mathbf{R}^n$  можно однозначно разложить в сумму  $\mathbf{y} = \mathbf{y}_r + \mathbf{y}_0$ , где  $\mathbf{y}_r \in \mathbf{R}(\mathbf{A})$ ,  $\mathbf{y}_0 \in \mathbf{N}(\mathbf{A}^T)$ .

Таким образом, каждая матрица  $\mathbf{A}$  вместе с транспонированной к ней матрицей  $\mathbf{A}^T$  выделяют, соответственно, в пространствах  $\mathbf{R}^k$  и  $\mathbf{R}^n$  две пары подпространств  $\{\mathbf{R}(\mathbf{A}^T), \mathbf{N}(\mathbf{A})\}$  и  $\{\mathbf{R}(\mathbf{A}), \mathbf{N}(\mathbf{A}^T)\}$ . Теоремы 2 и 4 устанавливают чрезвычайно важные соотношения между этими основными подпространствами. Объединяя их вместе можно сформулировать теорему:

**Основная теорема линейной алгебры.** Для любой матрицы  $\mathbf{A}$  размера  $n \times k$  справедливо:

$$\begin{aligned} \mathbf{R}(\mathbf{A}^T) \perp \mathbf{N}(\mathbf{A}), & & \mathbf{R}(\mathbf{A}) \perp \mathbf{N}(\mathbf{A}^T), \\ \mathbf{R}(\mathbf{A}^T) \oplus \mathbf{N}(\mathbf{A}) = \mathbf{R}^k, & & \mathbf{R}(\mathbf{A}) \oplus \mathbf{N}(\mathbf{A}^T) = \mathbf{R}^n, \\ \dim \mathbf{R}(\mathbf{A}^T) = \dim \mathbf{R}(\mathbf{A}) = r, & & \delta(\mathbf{A}) = k - r, \delta(\mathbf{A}^T) = n - r. \end{aligned}$$

Действие произвольной матрицы  $\mathbf{A}$  показано (очень схематично) на рис. 1. Произвольный вектор  $\mathbf{x}$  однозначно разлагается в сумму  $\mathbf{x} = \mathbf{x}_r + \mathbf{x}_0$ . Матрица  $\mathbf{A}$  преобразует компоненту  $\mathbf{x}_r \in \mathbf{R}(\mathbf{A}^T)$  в вектор  $\mathbf{A}\mathbf{x}_r = \mathbf{y}_r \in \mathbf{R}(\mathbf{A})$ , в то время как компонента  $\mathbf{x}_0 \in \mathbf{N}(\mathbf{A})$  преобразуется в нулевой вектор  $\theta_n$ .

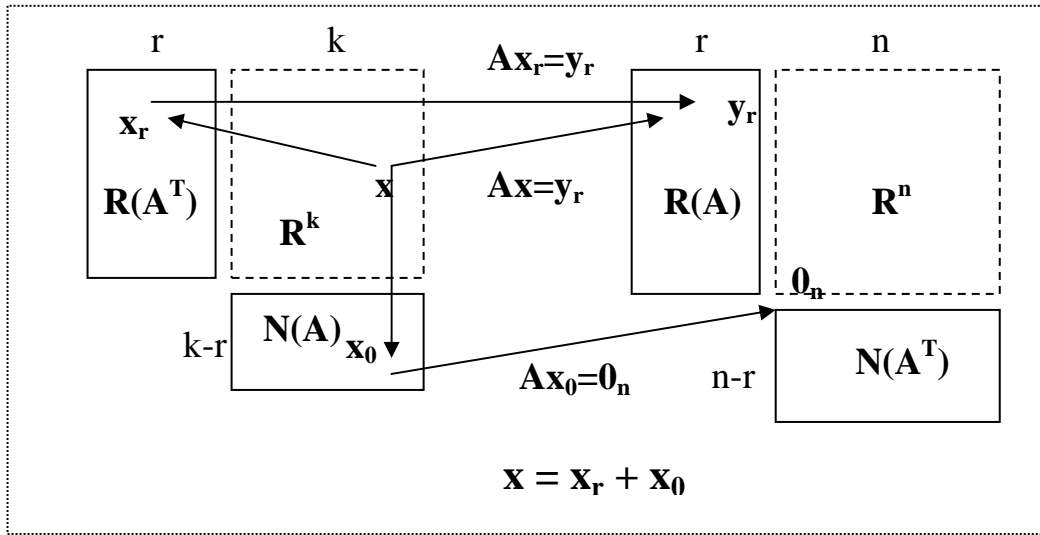


Рис. 1. Условное изображение расщепления пространств  $\mathbf{R}^k$  и  $\mathbf{R}^n$ .

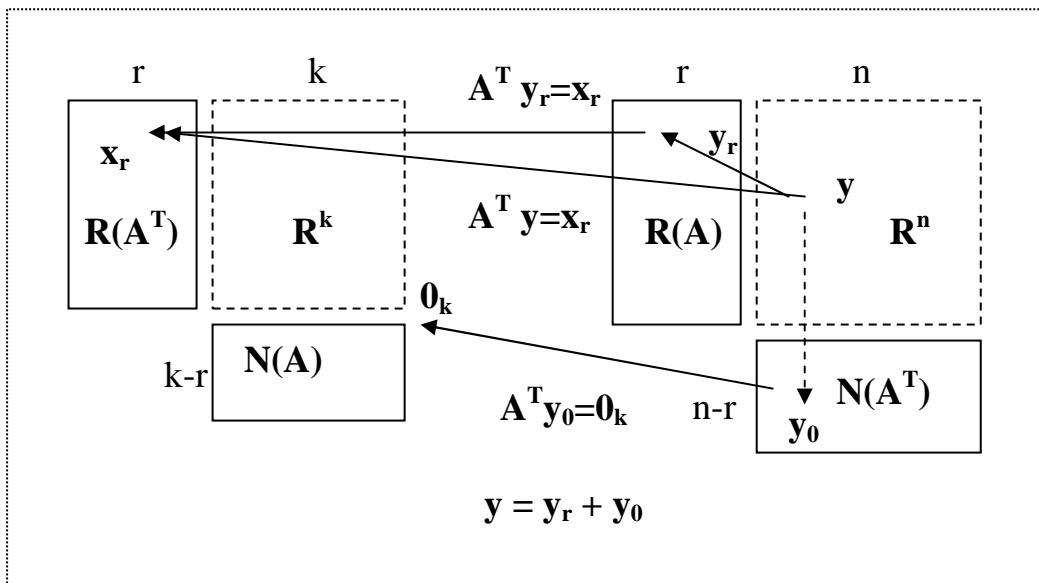


Рис. 2. Условное изображение действия матрицы  $\mathbf{A}^T$  размера  $k \times n$ ,  $\mathbf{R}(\mathbf{A}^T) \perp \mathbf{N}(\mathbf{A})$ ,  $\mathbf{R}(\mathbf{A}) \perp \mathbf{N}(\mathbf{A}^T)$ .

**Теорема 5.** Отображение пространства строк  $\mathbf{R}(\mathbf{A}^T)$  на пространство столбцов  $\mathbf{R}(\mathbf{A})$  с помощью матрицы  $\mathbf{A}$  невырожденное, т.е. обратимо: каждый вектор  $\mathbf{y}$  из пространства столбцов является образом одного и только одного вектора  $\mathbf{x}$ , принадлежащего пространству строк.

Доказательство. Если вектор  $\mathbf{y}$  принадлежит пространству столбцов, то он является некоторой линейной комбинацией  $\mathbf{A}\mathbf{x}$  столбцов матрицы  $\mathbf{A}$ . Раскладывая вектор  $\mathbf{x}$  в сумму  $\mathbf{x}_r + \mathbf{x}_0$ , где

$\mathbf{x}_r \in \mathbf{R}(\mathbf{A}^T)$ , а  $\mathbf{x}_0 \in \mathbf{N}(\mathbf{A})$ , получаем что  $\mathbf{Ax} = \mathbf{Ax}_r + \mathbf{Ax}_0 = \mathbf{Ax}_r = \mathbf{y}$ . Таким образом, мы нашли в пространстве строк подходящий вектор  $\mathbf{x}_r$ . Если бы там был другой вектор  $\mathbf{z}_r$ , такой, что  $\mathbf{Az}_r = \mathbf{y}$ , то  $\mathbf{A}(\mathbf{x}_r - \mathbf{z}_r) = \mathbf{y} - \mathbf{y} = \mathbf{0}_n$ . Поэтому вектор  $\mathbf{x}_r - \mathbf{z}_r$  принадлежит как пространству строк  $\mathbf{R}(\mathbf{A}^T)$ , так и нуль-пространству  $\mathbf{N}(\mathbf{A})$ . Следовательно, он ортогонален сам себе, а это возможно только для нулевого вектора, т.е.  $\mathbf{x}_r = \mathbf{z}_r$ .

Аналогично доказывается, что матрица  $\mathbf{A}^T$  является обратимым отображением в противоположном направлении из  $\mathbf{R}(\mathbf{A})$  на  $\mathbf{R}(\mathbf{A}^T)$ . Это не значит, что матрица  $\mathbf{A}^T$  является обратной к  $\mathbf{A}$ , просто для каждого  $\mathbf{x} \in \mathbf{R}(\mathbf{A}^T)$  существует один и только один элемент  $\mathbf{y} \in \mathbf{R}(\mathbf{A})$  такой, что  $\mathbf{x} = \mathbf{A}^T \mathbf{y}$ . Действие матрицы  $\mathbf{A}^T$  показано схематично на рис. 2.

*Каждая матрица  $\mathbf{A}$  обратима, если ее соответствующим образом воспринимать как отображение некоторого  $r$ - мерного подпространства на другое  $r$ - мерное подпространство, а именно  $\mathbf{R}(\mathbf{A}^T) \Leftrightarrow \mathbf{R}(\mathbf{A})$ .*

Основную теорему линейной алгебры можно сформулировать в виде *альтернативы Фредгольма*: для любых  $\mathbf{A}$  и  $\mathbf{y}$  одна и только одна из следующих задач:

$$(1) \quad \mathbf{Ax} = \mathbf{y}, \quad (2) \quad \mathbf{A}^T \mathbf{w} = \mathbf{0}, \quad (\mathbf{w}, \mathbf{y}) \neq 0,$$

имеет решение. Иначе говоря, либо вектор  $\mathbf{y} \in \mathbf{R}(\mathbf{A})$ , либо существует такой вектор  $\mathbf{w} \in \mathbf{N}(\mathbf{A}^T)$ , что  $(\mathbf{w}, \mathbf{y}) \neq 0$ . Получить подходящий вектор  $\mathbf{w} = \mathbf{y}_0$  можно, раскладывая вектор  $\mathbf{y}$  в сумму  $\mathbf{y} = \mathbf{y}_r + \mathbf{y}_0$ , где  $\mathbf{y}_r \in \mathbf{R}(\mathbf{A})$ ,  $\mathbf{y}_0 \in \mathbf{N}(\mathbf{A}^T)$ .

**Пример 2.** Рассмотрим отображение  $\mathbf{R}^3 \Rightarrow \mathbf{R}^2$  с помощью матрицы

$$\mathbf{A} = \begin{pmatrix} 1 & 0 & 2 \\ 1 & 1 & 4 \end{pmatrix}.$$

Одномерное нуль-пространство  $\mathbf{N}(\mathbf{A})$  порождается вектором  $\mathbf{x}_0 = (2, 2, -1)^T$  и этот вектор ортогонален пространству строк  $\mathbf{R}(\mathbf{A}^T)$ ,



порождаемому векторами  $\mathbf{x}_{r,1} = (1, 0, 2)^T$  и  $\mathbf{x}_{r,2} = (1, 1, 4)^T$ . Ранговое пространство  $\mathbf{R}(\mathbf{A})$  порождается векторами  $\mathbf{y}_{r,1} = (1, 1)^T$  и  $\mathbf{y}_{r,2} = (0, 1)^T$ , нуль-пространство  $\mathbf{N}(\mathbf{A}^T)$  состоит из единственного вектора  $\mathbf{y}_0 = (0, 0)^T$ . Из основной теоремы линейной алгебры заключаем что, например, вектор  $\mathbf{x} = (1, 1, 1)^T$  однозначно разлагается в сумму  $\mathbf{x} = \mathbf{x}_r + \mathbf{x}_0$ , где:  $\mathbf{x}_r = (-1, -1, 2)^T$ ,  $\mathbf{x}_r \in \mathbf{R}(\mathbf{A}^T)$ ,  $\mathbf{x}_0 = (2, 2, -1)^T$ ,  $\mathbf{x}_0 \in \mathbf{N}(\mathbf{A})$ , а система уравнений

$$\begin{pmatrix} 1 & 0 & 2 \\ 1 & 1 & 4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 3 \\ 6 \end{pmatrix}$$

имеет бесконечно много решений  $\mathbf{x} = \mathbf{x}_r \oplus \mathbf{N}(\mathbf{A})$  и общее решение

имеет вид:

$$\mathbf{x} = \mathbf{x}_r + \mathbf{x}_0 = \begin{pmatrix} -1 \\ -1 \\ 2 \end{pmatrix} + t \begin{pmatrix} 2 \\ 2 \\ -1 \end{pmatrix}, t \in \mathbf{R}.$$

Основная теорема линейной алгебры дает четкие ответы на вопросы существования и единственности решений систем уравнений произвольного порядка

$$\mathbf{Ax} = \mathbf{y}. \quad (8)$$

Рассмотрим наиболее часто встречающиеся в практических приложениях системы с вертикально-прямоугольными матрицами ( $n \geq k$ ).

Пусть вертикально-прямоугольная матрица  $\mathbf{A}$ , размера  $n \times k$ ,  $n \geq k$ , имеет *полный ранг*  $r = k$ . В этом случае,  $\dim \mathbf{N}(\mathbf{A}) = 0$ ,  $\dim \mathbf{R}(\mathbf{A}) = k$ ,  $\dim \mathbf{R}(\mathbf{A}^T) = k$  и  $\dim \mathbf{N}(\mathbf{A}^T) = n - k$ . Тогда имеем взаимно однозначное отображение  $\mathbf{R}^k \Leftrightarrow \mathbf{R}(\mathbf{A})$  и обратное отображение  $\mathbf{R}(\mathbf{A}) \Rightarrow \mathbf{R}^k$  задается матрицей

$$\mathbf{A}_L^{-1} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T, \quad (9)$$

т.к.  $\mathbf{A}_L^{-1}(\mathbf{Ax}) = (\mathbf{A}^T \mathbf{A})^{-1}(\mathbf{A}^T \mathbf{A})\mathbf{x} = \mathbf{x}$ . Матрица  $\mathbf{A}_L^{-1}$  называется *левой обратной* к матрице  $\mathbf{A}$ . Заметим, что матрица  $\mathbf{A}_L^{-1}$  размера  $k \times n$  существует, т.к. по теореме 3 квадратная матрица  $\mathbf{G} = \mathbf{A}^T \mathbf{A}$  – обратима. В частности, если матрица  $\mathbf{A}$  – полуортогональная (вектор-столбцы являются ортонормированными), то  $\mathbf{A}_L^{-1} = \mathbf{A}^T$ .

При решении переопределенных систем уравнений  $\mathbf{Ax} = \mathbf{y}$  полного ранга возможны три ситуации:

а) Вектор  $\mathbf{y} \in \mathbf{R}(\mathbf{A})$ , т.е.  $\mathbf{y} = \mathbf{y}_r \neq \theta_n$ ,  $\mathbf{y}_0 = \theta_n$ . Тогда система (8) имеет единственное *нетривиальное* решение

$$\mathbf{x} = \mathbf{A}_L^{-1} \mathbf{y} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}; \quad (10)$$

Эквивалентные условия *единственности* решения:

- 1) вектор-столбцы матрицы  $\mathbf{A}$  линейно независимы;
- 2) нуль-пространство  $\mathbf{N}(\mathbf{A})$  содержит только нулевой вектор;
- 3)  $\text{rank}(\mathbf{A}) = k$  ;
- 4) квадратная матрица  $\mathbf{A}^T \mathbf{A}$  является обратимой.

б) Ненулевой вектор  $\mathbf{y} \in \mathbf{N}(\mathbf{A}^T)$ . Тогда система (8) является *несовместной*, т.к. для любого  $\mathbf{x} \neq \theta_k$ , вектор  $\mathbf{y} = \mathbf{Ax} \notin \mathbf{N}(\mathbf{A}^T)$ ;

в) Вектор  $\mathbf{y} = \mathbf{y}_r + \mathbf{y}_0$ ,  $\mathbf{y}_r \neq \theta_n \in \mathbf{R}(\mathbf{A})$ ,  $\mathbf{y}_0 \neq \theta_n \in \mathbf{N}(\mathbf{A}^T)$ . В этом случае система (8) также является *несовместной*, т.к. для любого  $\mathbf{x} \in \mathbf{R}^k$  имеем  $\mathbf{Ax} = \mathbf{y} \in \mathbf{R}(\mathbf{A})$ .

Для несовместных систем полного ранга  $\mathbf{Ax} = \mathbf{y}_r + \mathbf{y}_0$  можно искать *решение наилучшего приближения*  $\mathbf{x}_{opt}$  такое, что для всех  $\mathbf{x} \in \mathbf{R}^k$  выполняется неравенство

$$\|\mathbf{Ax}_{opt} - \mathbf{y}\|^2 \leq \|\mathbf{Ax} - \mathbf{y}\|^2. \quad (11)$$

**Теорема 6.** Значение функции  $\Phi(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{y}\|^2$  достигает своего минимума в точке  $\mathbf{x}_{opt}$ , являющейся решением системы уравнений

$$\mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{y} = \mathbf{A}^T (\mathbf{y}_r + \mathbf{y}_0). \quad (12)$$

Доказательство.  $\Phi(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{y}\|^2 = \mathbf{x}^T \mathbf{A}^T \mathbf{Ax} - 2\mathbf{x}^T \mathbf{A}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}$ .

Тогда  $\mathbf{grad} \Phi(\mathbf{x}) = 2\mathbf{A}^T \mathbf{Ax} - 2\mathbf{A}^T \mathbf{y}$ . Так как  $\Phi(\mathbf{x}) > 0$  для всех  $\mathbf{x}$ , то из необходимого условия минимума квадратичной функции  $\Phi(\mathbf{x})$ , а именно  $\mathbf{grad} \Phi(\mathbf{x}) = \mathbf{0}$ , получаем так называемую *систему нормальных уравнений*

$$\mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{y}. \quad (13)$$

Система (13) имеет единственное решение  $\mathbf{x}_{opt} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$  и это решение называется *оптимальным*, потому что вектор

$$\mathbf{y}_r = \mathbf{Ax}_{opt} = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}, \quad \mathbf{y}_r \in \mathbf{R}(\mathbf{A}),$$

является ближайшим к исходному вектору  $\mathbf{y} = \mathbf{y}_r + \mathbf{y}_0$ .

Поскольку для произвольного вектора  $\mathbf{y}$  вектор

$$\mathbf{y}_r = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y} = \mathbf{P}_r \mathbf{y}$$

является компонентой этого вектора в пространстве столбцов матрицы, то матрица

$$\mathbf{P}_r = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T = \mathbf{AA}_L^{-1} \quad (14)$$

суть *матрица проектирования* на пространство столбцов  $\mathbf{R}(\mathbf{A})$ . Далее мы увидим, что в случае матриц неполного ранга матрица проектирования на пространство столбцов имеет аналогичный вид  $\mathbf{P}_r = \mathbf{AA}^+$ , где  $\mathbf{A}^+$  – так называемая обобщенная (*псевдообратная*) матрица.

Вектор  $\mathbf{y} - \mathbf{P}_r \mathbf{y}$  является компонентой в ортогональном дополнении  $\mathbf{N}(\mathbf{A}^T)$ . Поэтому матрица

$$\mathbf{P}_0 = \mathbf{E}_n - \mathbf{P}_r = \mathbf{E}_n - \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$$

является *матрицей проектирования* на ортогональное дополнение к  $\mathbf{R}(\mathbf{A})$  (т.е. нуль-пространство  $\mathbf{N}(\mathbf{A}^T)$ ). Таким образом, имеем матричную формулу для разбиения вектора на две взаимно перпендикулярные составляющие:

$$\mathbf{y} = \mathbf{P}_r \mathbf{y} + \mathbf{P}_0 \mathbf{y} = \mathbf{y}_r + \mathbf{y}_0, \quad \mathbf{y}_r \in \mathbf{R}(\mathbf{A}), \quad \mathbf{y}_0 \in \mathbf{N}(\mathbf{A}^T).$$

Матрицы проектирования обладают специфическими свойствами:

**Теорема 7.** Матрица проектирования  $\mathbf{P} = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$  обладает двумя основными свойствами – она является:

1) *идемпотентной*, т.е.  $\mathbf{P}\mathbf{P} = \mathbf{P}$ ;

2) *симметричной*, т.е.  $\mathbf{P} = \mathbf{P}^T$ .

И обратно, любая матрица с этими двумя свойствами представляет собой матрицу проектирования на свое пространство столбцов.

Доказательство.

$$\mathbf{P}\mathbf{P} = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T = \mathbf{P},$$

т.е. матрица идемпотентна. Для доказательства симметричности матрицы  $\mathbf{P}$  рассмотрим два произвольных вектора  $\mathbf{y}$  и  $\mathbf{z}$ . Тогда векторы

$$\mathbf{y}_r = \mathbf{P}\mathbf{y} \quad \text{и} \quad \mathbf{z}_0 = (\mathbf{E}_n - \mathbf{P})\mathbf{z}$$

являются взаимно ортогональными:

$$(\mathbf{P}\mathbf{y})^T (\mathbf{E}_n - \mathbf{P})\mathbf{z} = \mathbf{y}^T \mathbf{P}^T (\mathbf{E}_n - \mathbf{P})\mathbf{z} = 0.$$

Поскольку это справедливо для любых векторов  $\mathbf{y}$  и  $\mathbf{z}$ , получаем, что

$\mathbf{P}^T (\mathbf{E}_n - \mathbf{P}) = (\mathbf{0}_{n \times n})$ , т.е.  $\mathbf{P}^T = \mathbf{P}^T \mathbf{P}$ , или  $\mathbf{P} = \mathbf{P}^T \mathbf{P}$ . Следовательно,  $\mathbf{P}^T = \mathbf{P}$ . (Символ  $(\mathbf{0}_{n \times n})$  обозначает нуль-матрицу размера  $n \times n$ ).

Для доказательства обратного необходимо, используя свойства 1) и 2), показать, что  $\mathbf{P}$  является матрицей проектирования на свое пространство столбцов. Это пространство состоит из всех линейных комбинаций столбцов матрицы  $\mathbf{P}$ . Для любого вектора  $\mathbf{y}$  вектор  $\mathbf{P}\mathbf{y}$  обязательно лежит в этом пространстве. Кроме того, вектор  $\mathbf{y} - \mathbf{P}\mathbf{y}$  ортогонален этому пространству: для любого вектора  $\mathbf{P}\mathbf{z}$  из пространства столбцов свойства 1) и 2) дают

$$(\mathbf{y} - \mathbf{P}\mathbf{y})^T \mathbf{P}\mathbf{z} = \mathbf{y}^T (\mathbf{E}_n - \mathbf{P})^T \mathbf{P}\mathbf{z} = \mathbf{y}^T (\mathbf{P} - \mathbf{P}\mathbf{P})\mathbf{z} = 0.$$

Поскольку вектор  $\mathbf{y} - \mathbf{P}\mathbf{y}$  ортогонален пространству столбцов, он и является искомым перпендикуляром, так что  $\mathbf{P}$  оказывается матрицей проектирования. Матрица  $\mathbf{E}_n - \mathbf{P}$  также обладает этими двумя основными свойствами.

**Пример 3.** Пусть в подпространстве  $\mathbf{U} \in \mathbf{R}^3$  задан базис  $\mathbf{u}_1 = (1, 1, 0)^T$ ,  $\mathbf{u}_2 = (1, 0, 1)^T$  и вектор  $\mathbf{y} = (0, 2, 1)^T \notin \mathbf{U}$ . Тогда матрица  $\mathbf{A}$  со столбцами  $\mathbf{u}_1$  и  $\mathbf{u}_2$  имеет вид:

$$\mathbf{A} = \begin{pmatrix} 1 & 1 \\ 1 & 0 \\ 0 & 1 \end{pmatrix},$$

а матрица проектирования на подпространство  $\mathbf{R}(\mathbf{A})$ :

$$\mathbf{P}_r = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T = \frac{1}{3} \begin{pmatrix} 2 & 1 & 1 \\ 1 & 2 & -1 \\ 1 & -1 & 2 \end{pmatrix}.$$

Проекция вектора  $\mathbf{y}$  на подпространство  $\mathbf{R}(\mathbf{A})$  –  $\mathbf{y}_r = \mathbf{P}_r \mathbf{y} = (1, 1, 0)^T$ ; проекцией на ортогональное дополнение  $\mathbf{N}(\mathbf{A}^T)$  есть вектор  $\mathbf{y}_0 = (\mathbf{E}_3 - \mathbf{P}_r) \mathbf{y} = (-1, 1, 1)^T$ .

Для полуортогональных матриц  $\mathbf{A}$  размера  $n \times k$ ,  $k \leq n$  имеем  $\mathbf{A}^T \mathbf{A} = \mathbf{E}_k$ , и тогда матрицы проектирования на подпространства  $\mathbf{R}(\mathbf{A})$  и  $\mathbf{N}(\mathbf{A}^T)$  вычисляются очень легко, что видно из следующего примера.

**Пример 4.**

$$\mathbf{A} = \begin{pmatrix} 1/2 & 1/2 \\ 1/2 & -1/2 \\ 1/2 & 1/2 \\ 1/2 & -1/2 \end{pmatrix}, \quad \mathbf{P}_r = \mathbf{A} \mathbf{A}^T = \begin{pmatrix} 1/2 & 0 & 1/2 & 0 \\ 0 & 1/2 & 0 & 1/2 \\ 1/2 & 0 & 1/2 & 0 \\ 0 & 1/2 & 0 & 1/2 \end{pmatrix},$$

$$\mathbf{P}_0 = \mathbf{E}_4 - \mathbf{P}_r = \begin{pmatrix} 1/2 & 0 & -1/2 & 0 \\ 0 & 1/2 & 0 & -1/2 \\ -1/2 & 0 & 1/2 & 0 \\ 0 & -1/2 & 0 & 1/2 \end{pmatrix}.$$

Пусть для матрицы полного ранга мы решаем *несовместную* переопределенную систему уравнений

$$\mathbf{A} \mathbf{x} = \mathbf{y} = \mathbf{y}_r + \mathbf{y}_0 = \mathbf{P}_r \mathbf{y} + \mathbf{P}_0 \mathbf{y}. \quad (15)$$

Тогда оптимальное решение есть решение системы нормальных уравнений

$$\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{y} = \mathbf{A}^T (\mathbf{y}_r + \mathbf{y}_0) = \mathbf{A}^T \mathbf{P}_r \mathbf{y} + \mathbf{A}^T \mathbf{P}_0 \mathbf{y}.$$

Так как  $\mathbf{A}^T \mathbf{P}_0 \mathbf{y} = \mathbf{A}^T (\mathbf{E}_n - \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T) \mathbf{y} \equiv \theta_n$  для всех  $\mathbf{y}$ , то оптимальное решение системы есть вектор

$$\mathbf{x}_{opt} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}_r, \quad (16)$$

и, следовательно, вектор  $\mathbf{y}_0$  не оказывает никакого влияния на величину оптимального решения, а только определяет длину вектора оптимальной невязки

$$\|\mathbf{A} \mathbf{x}_{opt} - \mathbf{y}\| = \|\mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y} - (\mathbf{y}_r + \mathbf{y}_0)\| = \|\mathbf{y}_r - (\mathbf{y}_r + \mathbf{y}_0)\| = \|\mathbf{y}_0\|.$$

Следствие. В силу выше сказанного ясно, что при решении системы уравнений

$$\mathbf{A}(\mathbf{x} + \delta \mathbf{x}) = \mathbf{y}_r + \delta \mathbf{y},$$

содержащей погрешности  $\delta \mathbf{y} = \delta \mathbf{y}_r + \delta \mathbf{y}_0$ ,  $\delta \mathbf{y}_r \in \mathbf{R}(\mathbf{A})$ ,  $\delta \mathbf{y}_0 \in \mathbf{N}(\mathbf{A}^T)$ , вектор  $\delta \mathbf{y}_0$  не влияет на значение приближенного решения. Поэтому длина вектора оптимальной невязки  $\|\delta \mathbf{y}_0\|$  не может служить характеристикой близости приближенного решения к истинному (см. рис. 3). Возмущение решения  $\delta \mathbf{x} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \delta \mathbf{y}_r$  зависит только от компоненты  $\delta \mathbf{y}_r$ . Однако, длина  $\|\delta \mathbf{x}\|$  зависит не только от величины  $\|\delta \mathbf{y}_r\|$ , но также от направления вектора  $\delta \mathbf{y}_r$ , что будет подробно обсуждаться ниже. В частности, при фиксированном значении  $\|\delta \mathbf{y}\|$ , вектор  $\delta \mathbf{x}$  изменяется в диапазоне:

$$\begin{cases} \delta \mathbf{x}_{\min} = \theta_k, & \delta \mathbf{y} \in \mathbf{N}(\mathbf{A}^T) \\ \delta \mathbf{x}_{\max} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \delta \mathbf{y}, & \delta \mathbf{y} \in \mathbf{R}(\mathbf{A}) \end{cases}.$$

Если погрешность  $\delta \mathbf{y}$  по своей природе имеет случайный характер, т.е. направление вектора погрешности непредсказуемо, то длина вектора возмущения решения будет удовлетворять неравенству

$$\begin{aligned} 0 \leq \|\delta \mathbf{x}\| &\leq \|\delta \mathbf{x}_{\max}\| = \|(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T\| \|\delta \mathbf{y}_r\| = \\ &= \|(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T\| \sqrt{\|\delta \mathbf{y}\|^2 - \|\delta \mathbf{y}_0\|^2} \end{aligned} \quad (17)$$

Величина

$$\rho = \|(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T\| \sqrt{\|\delta \mathbf{y}\|^2 - \|\delta \mathbf{y}_0\|^2} \quad (18)$$

определяет максимально возможное расстояние между истинным и приближенным решениями и это расстояние определяет *радиус доверительной области* возмущенного решения.

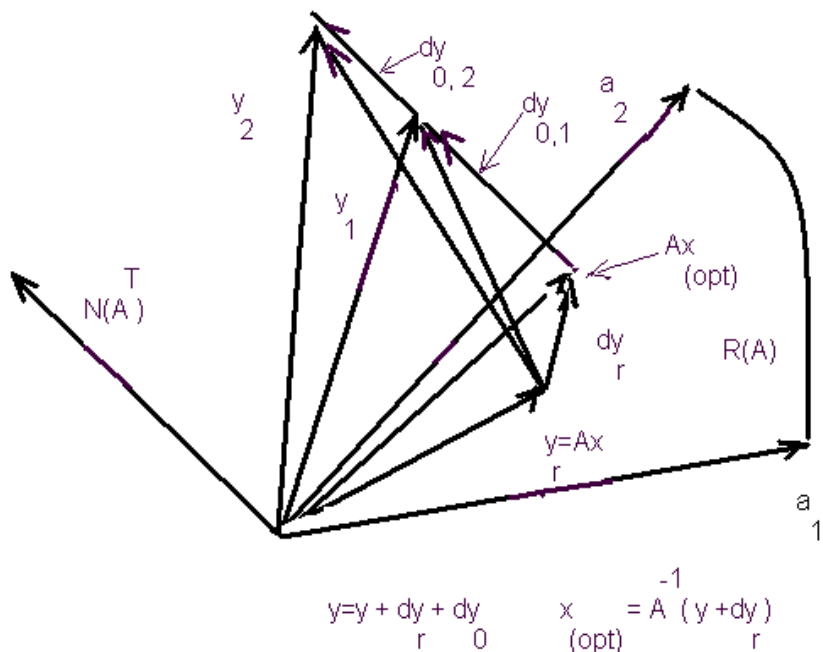


Рис. 3. Оптимальное решение не зависит от компоненты  $dy_0 \in N(A^T)$ .

Рассмотрение вопросов решения систем уравнений с матрицами неполного ранга мы проведем после изложения *метода сингулярного разложения матриц*.

## §2. Сингулярное разложение матриц

*Метод решения хорош, если с самого начала мы можем предвидеть – и далее подтвердить это, – что, следуя этому методу, мы достигнем цели.*

*Лейбниц*

Сингулярное разложение матриц (SVD – Singular Value Decomposition) является своеобразным “*томографом высокого разрешения*”. Этот метод позволяет определить ранг матрицы  $\mathbf{A}$ , меру и тип обусловленности, найти ортогональные базисы подпространств  $\mathbf{R}(\mathbf{A})$ ,  $\mathbf{R}(\mathbf{A}^T)$ ,  $\mathbf{N}(\mathbf{A})$ ,  $\mathbf{N}(\mathbf{A}^T)$  и, самое главное, вычислить обобщенную обратную матрицу (псевдообратную), найти оптимальное решение несовместной системы уравнений с матрицами неполного ранга, приближенное решение систем уравнений с плохо обусловленной матрицей полного ранга, радиус доверительной области приближенных решений и ... многое другое! Великолепный набор качеств этого метода заслуживает его детального рассмотрения. Алгоритм SVD входит в десятку наиболее важных алгоритмических достижений 20-го века.

Напомним, что нижеследующие условия единственности решения системы уравнений  $\mathbf{Ax} = \mathbf{y}$  являются эквивалентными:

- (1) столбцы матрицы  $\mathbf{A}$  размера  $n \times k$ ,  $n \geq k$  – линейно независимы;
- (2) нуль-пространство  $\mathbf{N}(\mathbf{A})$  содержит только нулевой вектор;
- (3)  $\text{rank}(\mathbf{A}) = k$ ;
- (4) квадратная симметричная матрица  $\mathbf{A}^T \mathbf{A}$  является обратимой.

Для этого случая система  $\mathbf{Ax} = \mathbf{y}$  имеет единственное решение, которое вычисляется по сравнительно простой формуле:

$$\mathbf{x} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y} = \mathbf{A}_L^{-1} \mathbf{y}.$$

Но если условия 1)÷4) не выполняются и вектор  $\mathbf{x}$  определяется из системы  $\mathbf{Ax} = \mathbf{y}$  неединственным образом или система несовместна, то необходимо определить *псевдообратную* матрицу  $\mathbf{A}^+$  такую, что



$\mathbf{x} = \mathbf{A}^+ \mathbf{y}$  будет оптимальным решением по методу наименьших квадратов системы  $\mathbf{Ax} = \mathbf{y}$ .

**Определение.** *Оптимальным среди всех решений системы  $\mathbf{Ax} = \mathbf{y}$  по методу наименьших квадратов является решение с минимальной длиной (евклидовой нормой).*

Пусть рассматриваемая система уравнений с матрицей неполного ранга является совместной. Для отыскания оптимального решения вспомним, что пространство строк  $\mathbf{R}(\mathbf{A}^T)$  и нуль-пространство  $\mathbf{N}(\mathbf{A})$  представляют собой ортогональные дополнения друг к другу в пространстве  $\mathbf{R}^k$ . Это означает, что любой вектор может быть разложен на две перпендикулярные составляющие – проекцию на пространство строк и проекцию на нуль-пространство. Применим это разложение к одному из решений (обозначим его через  $\mathbf{x}_s$ ) системы  $\mathbf{Ax} = \mathbf{y}$ . Тогда  $\mathbf{x}_s = \mathbf{x}_r + \mathbf{x}_0$ , где  $\mathbf{x}_r \in \mathbf{R}(\mathbf{A}^T)$ ,  $\mathbf{x}_0 \in \mathbf{N}(\mathbf{A})$ . Здесь возникают три существенных момента:

- 1) Компонента  $\mathbf{x}_r$  сама является решением системы  $\mathbf{Ax} = \mathbf{y}$ : поскольку  $\mathbf{Ax}_0 = \mathbf{0}$ , имеем  $\mathbf{Ax}_s = \mathbf{A}(\mathbf{x}_r + \mathbf{x}_0) = \mathbf{Ax}_r = \mathbf{y}$ ;
- 2) Все решения системы  $\mathbf{Ax} = \mathbf{y}$  имеют одну и ту же компоненту  $\mathbf{x}_r$  в пространстве строк  $\mathbf{R}(\mathbf{A}^T)$  и отличаются лишь компонентой  $\mathbf{x}_0$  из нуль-пространства  $\mathbf{N}(\mathbf{A})$ , т.к. общее решение равняется сумме частного решения (в данном случае это  $\mathbf{x}_r$ ) и произвольного решения  $\mathbf{x}_0 \in \mathbf{N}(\mathbf{A})$  однородного уравнения  $\mathbf{Ax}_0 = \mathbf{0}$ .
- 3) Длина такого решения  $\mathbf{x}_r + \mathbf{x}_0$  удовлетворяет теореме Пифагора, поскольку две компоненты взаимно ортогональны:

$$\|\mathbf{x}_r + \mathbf{x}_0\|^2 = \|\mathbf{x}_r\|^2 + \|\mathbf{x}_0\|^2.$$

Отсюда мы заключаем: *решением с минимальной длиной является вектор  $\mathbf{x}_r$ .*

*Оптимальным решением по методу наименьших квадратов произвольной системы  $\mathbf{Ax} = \mathbf{y}$  является вектор  $\mathbf{x}_r$ , удовлетворяющий двум условиям:*

- 1) *вектор  $\mathbf{Ax}_r$  равняется проекции вектора  $\mathbf{y}$  на пространство столбцов  $\mathbf{R}(\mathbf{A})$ ;*

2) вектор  $\mathbf{x}_r$  принадлежит пространству строк  $\mathbf{R}(\mathbf{A}^T)$ .

**Определение.** Матрица  $\mathbf{A}^+$  называется *псевдообратной* к  $\mathbf{A}$ , если для  $\mathbf{x}_r = \mathbf{A}^+ \mathbf{y}$  выполняются вышеприведенные условия 1) и 2).

Смысл псевдообратной матрицы можно легко понять из геометрических соображений, если обратиться к четырем основным подпространствам матрицы  $\mathbf{A}$ . Матрица  $\mathbf{A}^+$  должна сочетать в себе выполнение двух отдельных шагов:

- 1) проектировать вектор  $\mathbf{y}$  на  $\mathbf{R}(\mathbf{A})$ , т.е. вычислять  $\mathbf{y}_r = \mathbf{P}_r \mathbf{y}$ ;
- 2) выбирать единственный вектор  $\mathbf{x}$  из пространства строк  $\mathbf{R}(\mathbf{A}^T)$ , являющийся решением системы  $\mathbf{A}\mathbf{x} = \mathbf{y}_r$ .

Один крайний случай возникает, когда вектор  $\mathbf{y}$  перпендикулярен  $\mathbf{R}(\mathbf{A})$ , т.е. когда вектор  $\mathbf{y} \in \mathbf{N}(\mathbf{A}^T)$ . Тогда  $\mathbf{y}_r = \mathbf{0}_n$ ,  $\mathbf{x} = \mathbf{0}_k$  и матрица  $\mathbf{A}^+$  должна переводить все векторы из  $\mathbf{N}(\mathbf{A}^T)$  в нуль-вектор  $\mathbf{0}_k$ . Другая крайность состоит в том, что вектор  $\mathbf{y}$  может лежать целиком в пространстве столбцов  $\mathbf{R}(\mathbf{A})$  и тогда решение должно находиться путем “обращения” матрицы  $\mathbf{A}$ . (Мы уже отмечали в §1, что матрица  $\mathbf{A}$  обратима, если рассматривать ее просто как отображение из ее пространства строк в пространство столбцов и обратной должна быть  $\mathbf{A}^+$ ). Из этого описания, подкрепленного иллюстрацией на рис. 4, можно сформулировать некоторые основные свойства псевдообратной матрицы:

- 1) Матрица  $\mathbf{A}^+$  имеет размер  $k \times n$ , т.к. она применяется к вектору  $\mathbf{y} \in \mathbf{R}^n$  и дает  $\mathbf{x} \in \mathbf{R}^k$ ;
- 2) Пространство столбцов матрицы  $\mathbf{A}^+$  совпадает с пространством строк матрицы  $\mathbf{A}$ , а пространство строк матрицы  $\mathbf{A}^+$  совпадает с пространством столбцов  $\mathbf{R}(\mathbf{A})$ . Заметим, что в качестве следствия получаем  $\text{rank}(\mathbf{A}^+) = \text{rank}(\mathbf{A})$ ;
- 3) Псевдообратная к матрице  $\mathbf{A}^+$  есть матрица  $\mathbf{A}$ ;

4) В общем случае  $\mathbf{A}\mathbf{A}^+ \neq \mathbf{E}_n$ , поскольку матрица  $\mathbf{A}$  может не иметь правой обратной, но матрица  $\mathbf{A}\mathbf{A}^+$  всегда совпадает с матрицей  $\mathbf{P}_r$ , осуществляющей проектирование на пространство столбцов  $\mathbf{R}(\mathbf{A})$ :

$$\mathbf{A}\mathbf{A}^+ \mathbf{y} = \mathbf{A}\mathbf{x}_r = \mathbf{y}_r = \mathbf{P}_r \mathbf{y}, \quad \text{т.е.} \quad \mathbf{A}\mathbf{A}^+ = \mathbf{P}_r.$$

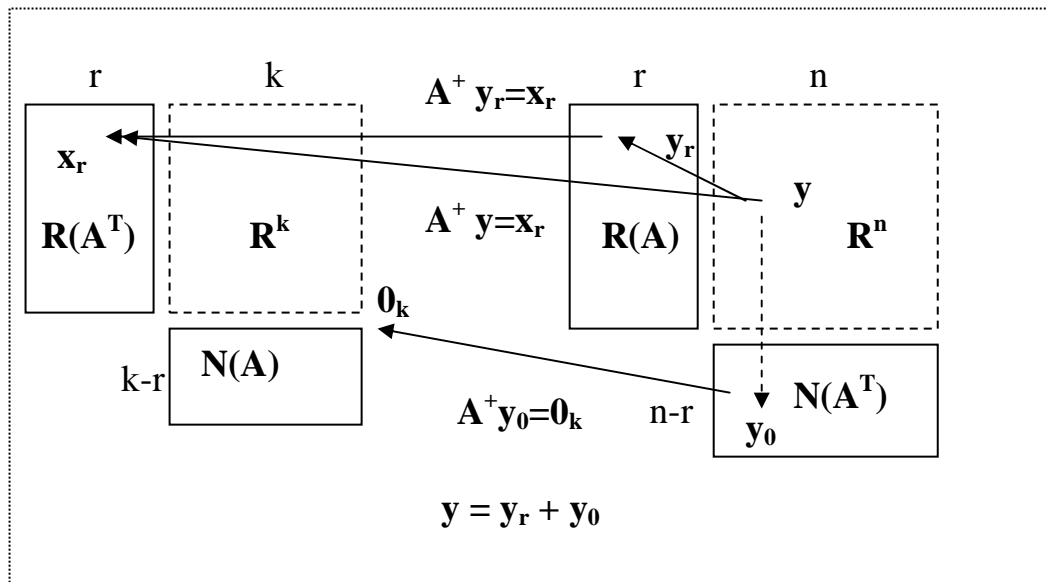


Рис. 4. Проекционные действия псевдообратной матрицы  $\mathbf{A}^+$ .

**Пример 5.** Пусть

$$\mathbf{A} = \begin{pmatrix} 10 & 0 & 0 \\ 0 & 0.5 & 0 \\ 0 & 0 & 0 \end{pmatrix}.$$

Базисами пространства столбцов и пространства строк являются векторы  $\mathbf{e}_1 = (1, 0, 0)^T$  и  $\mathbf{e}_2 = (0, 1, 0)^T$ , вектор  $\mathbf{e}_3 = (0, 0, 1)^T$  порождает нуль-пространства  $\mathbf{N}(\mathbf{A})$  и  $\mathbf{N}(\mathbf{A}^T)$ . Решение системы  $\mathbf{A}\mathbf{x} = \mathbf{y}$  начинается с проектирования  $\mathbf{y}$  на пространство столбцов  $\mathbf{R}(\mathbf{A})$ :

$$\text{если } \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ y_3 \end{pmatrix}, \quad \text{то } \mathbf{y}_r = \mathbf{P}_r \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ 0 \end{pmatrix}.$$

Затем решается система  $\mathbf{Ax} = \mathbf{y}_r$ :

$$\begin{pmatrix} 10 & 0 & 0 \\ 0 & 0.5 & 0 \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \\ 0 \end{pmatrix}.$$

Решение выписывается легко:  $x_1 = 0.1y_1$ ,  $x_2 = 2y_2$  и  $x_3 = t$ ,  $t \in R$ . Из этого бесконечного семейства решений выбираем одно, имеющее минимальную длину, полагая  $x_3 = 0$ . Таким образом, мы получили решение  $\mathbf{x}_{opt} = (0.1y_1, 2y_2, 0)^T$ , которое лежит в пространстве строк. Это решение можно получить с помощью псевдообратной матрицы следующим образом. Если  $\mathbf{S} = \mathbf{A}$ , тогда

$$\mathbf{A}^+ = \mathbf{S}^+ = \begin{pmatrix} 0.1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \end{pmatrix}, \quad \mathbf{x}_{opt} = \mathbf{A}^+ \mathbf{y}, \quad \mathbf{SS}^+ = \mathbf{P}_r = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}.$$

Этот пример является типичным для всего семейства матриц  $\mathbf{S}$  специального вида, имеющих в качестве первых  $r$  элементов главной диагонали положительные числа  $\sigma_1, \sigma_2, \dots, \sigma_r$  и остальные элементы равны нулю. Псевдообратная  $\mathbf{S}^+$  к матрице  $\mathbf{S}$  вычисляется легко: первые  $r$  элементов на ее главной диагонали будут числа  $1/\sigma_i$ , остальные равны нулю. Если матрица  $\mathbf{S}$  имеет размер  $n \times k$ , то  $\mathbf{S}^+$  будет иметь размер  $k \times n$ . Легко убедиться, что ранг матрицы  $\mathbf{S}^+$  равен рангу  $\mathbf{S}$ , дефект матрицы  $\mathbf{S}$  равен  $k - r$ , дефект матрицы  $\mathbf{S}^+$  равен  $n - r$  и  $(\mathbf{S}^+)^+ = \mathbf{S}$ , т.е. псевдообратная к  $\mathbf{S}^+$  совпадает с  $\mathbf{S}$ .

В общем случае основой вычисления псевдообратных матриц служит теорема о **сингулярном разложении матриц**. Доказательство этой теоремы существенным образом опирается на известные теоремы линейной алгебры:

**Теорема 8.** Если  $\mathbf{A}$  – симметричная матрица размера  $n \times n$  с действительными элементами, то её собственные значения – действительные числа.

Доказательство. Если  $(\mathbf{A} - \lambda \mathbf{E})\mathbf{x} = \theta_n$ , то и  $(\mathbf{A} - \lambda^* \mathbf{E})\mathbf{x} = \theta_n$ , где  $\lambda^*$  – комплексно сопряженная величина по отношению к  $\lambda$ . Следовательно,  $\lambda = \lambda^*$ .

**Теорема 9.** Для матрицы  $\mathbf{A}$  произвольного размера собственные значения симметричной матрицы  $\mathbf{G} = \mathbf{A}^T \mathbf{A}$  – неотрицательны, т.е.  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$ .

**Теорема 10.** (О диагонализации матриц). Если  $\mathbf{A}$  – симметричная матрица размера  $n \times n$  с действительными элементами и собственными числами  $\lambda_1, \lambda_2, \dots, \lambda_n$ , то существует ортогональная матрица  $\mathbf{V}$  такая, что

$$\mathbf{V}^T \mathbf{A} \mathbf{V} = \Lambda,$$

где  $\Lambda$  – диагональная матрица с диагональными элементами  $\lambda_j$  собственных значений. Матрица  $\mathbf{V}$  – ортогональная, столбцы матрицы  $\mathbf{V}$  – собственные векторы матрицы  $\mathbf{A}$ . При этом  $j$ -й столбец  $\mathbf{v}_j$  матрицы  $\mathbf{V}$  соответствует  $j$ -му собственному вектору, т.е.  $\mathbf{A} \mathbf{v}_j = \lambda_j \mathbf{v}_j$ .

**Теорема 11.** (Сингулярное разложение матриц). Произвольная матрица  $\mathbf{A}$  размера  $n \times k$ ,  $n \geq k$  может быть представлена в виде

$$\mathbf{A} = \mathbf{U} \mathbf{S} \mathbf{V}^T, \quad (19)$$

где  $\mathbf{U}$  – ортогональная матрица размера  $n \times n$ ,  $\mathbf{V}$  – ортогональная матрица размера  $k \times k$ , а матрица  $\mathbf{S}$  размера  $n \times k$  имеет специальную диагональную форму: на главной диагонали находятся числа  $\sigma_i > 0$ ,  $i = 1, 2, \dots, r$ ,  $r \leq k$  и  $\sigma_i = 0$ ,  $i = r + 1, \dots, k$ . Числа  $\sigma_i$  называются *сингулярными числами* матрицы.

Доказательство. Доказательство справедливости утверждения (19) основано на использовании теорем 8 ÷ 10: для квадратной симметричной матрицы  $\mathbf{G} = \mathbf{A}^T \mathbf{A}$ , размера  $k \times k$ , существует ортонормированный набор собственных векторов  $\mathbf{v}_i$ ,  $i = 1, 2, \dots, k$ ,  $\mathbf{v}_i \in \mathbf{R}^k$  (которые образуют столбцы матрицы  $\mathbf{V}$ ),

$$\text{т.е. } \mathbf{A}^T \mathbf{A} \mathbf{v}_i = \lambda_i \mathbf{v}_i, \text{ причем } \mathbf{v}_i^T \mathbf{v}_i = 1 \text{ и } \mathbf{v}_i^T \mathbf{v}_j = 0 \text{ для } i \neq j. \quad (20)$$

Умножая скалярно на  $\mathbf{v}_i$ , получаем, что  $\lambda_i \geq 0$ :

$$\mathbf{v}_i^T \mathbf{A}^T \mathbf{A} \mathbf{v}_i = \lambda_i \mathbf{v}_i^T \mathbf{v}_i, \text{ или } \|\mathbf{A} \mathbf{v}_i\|^2 = \lambda_i. \quad (21)$$

Предположим, что числа  $\lambda_1, \lambda_2, \dots, \lambda_r$  положительны, а остальные  $k-r$  векторов  $\mathbf{A} \mathbf{v}_i$  и чисел  $\lambda_i$  равны нулю. Для  $\lambda_i > 0$  положим  $\sigma_i = \sqrt{\lambda_i}$  и  $\mathbf{u}_i = \mathbf{A} \mathbf{v}_i / \sigma_i$ ,  $\mathbf{u}_i \in \mathbf{R}^n$ . Из (21) следует, что  $\mathbf{u}_i$  являются единичными векторами, причем ортогональными (что вытекает из (20)):

$$(\mathbf{u}_i, \mathbf{u}_j) = \frac{\mathbf{v}_i^T \mathbf{A}^T \mathbf{A} \mathbf{v}_j}{\sigma_i \sigma_j} = \frac{\lambda_j \mathbf{v}_i^T \mathbf{v}_j}{\sigma_i \sigma_j} = \begin{cases} 1, & i = j \\ 0, & i \neq j \end{cases}.$$

Следовательно, мы имеем систему  $r$  ортонормированных векторов, которая может быть расширена до полного ортонормированного базиса  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r, \dots, \mathbf{u}_n$ , образующего столбцы матрицы  $\mathbf{U}$ , при помощи алгоритма ортогонализации Грама-Шмидта. Тогда элементами матрицы  $\mathbf{U}^T \mathbf{A} \mathbf{V}$  будут числа  $\mathbf{u}_i^T \mathbf{A} \mathbf{v}_j$ , равные нулю при  $j > r$

(поскольку в этом случае  $\mathbf{A} \mathbf{v}_j = \theta_n$ ) и равные  $\mathbf{u}_i^T \sigma_j \mathbf{u}_j$  при  $j \leq r$ . Последние равняются нулю при  $i \neq j$  и  $\sigma_j$  при  $i = j$ . Другими словами, матрица  $\mathbf{U}^T \mathbf{A} \mathbf{V}$  в точности совпадает со специальной матрицей  $\mathbf{S}$ , у которой на главной диагонали расположены числа  $\sigma_i > 0$ ,  $i = 1, 2, \dots, r$  и  $\sigma_i = 0$  для  $i = r+1, \dots, n$ . Следовательно,

$$\mathbf{A} = \mathbf{U} \mathbf{S} \mathbf{V}^T.$$

Псевдообратная матрица вычисляется по формуле:

$$\mathbf{A}^+ = \mathbf{V} \mathbf{S}^+ \mathbf{U}^T. \quad (22)$$

Эта формула доказывается непосредственно из принципа наименьших квадратов. Минимизируемая невязка  $\|\mathbf{A} \mathbf{x} - \mathbf{y}\|$  с учетом разложения (19) имеет вид:

$$\Phi_1(\mathbf{x}) = \|\mathbf{A} \mathbf{x} - \mathbf{y}\| = \|\mathbf{U} \mathbf{S} \mathbf{V}^T \mathbf{x} - \mathbf{y}\|. \quad (23)$$

Так как умножение на ортогональную матрицу  $\mathbf{U}^T$  не изменяет длину вектора  $\mathbf{A} \mathbf{x} - \mathbf{y}$ , то

$$\Phi_1(\mathbf{x}) = \left\| \mathbf{S}\mathbf{V}^T \mathbf{x} - \mathbf{U}^T \mathbf{y} \right\|.$$

Введем новый неизвестный вектор  $\mathbf{x}' = \mathbf{V}^T \mathbf{x} = \mathbf{V}^{-1} \mathbf{x}$ , длина которого совпадает с длиной вектора  $\mathbf{x}$ . Тогда будем минимизировать величину  $\left\| \mathbf{S}\mathbf{x}' - \mathbf{U}^T \mathbf{y} \right\|$ , и оптимальное решение  $\mathbf{x}'_{opt}$  (т.е. решение с минимальной длиной) имеет вид  $\mathbf{x}'_{opt} = \mathbf{S}^+ \mathbf{U}^T \mathbf{y}$ . Следовательно, оптимальное решение системы  $\mathbf{A}\mathbf{x} = \mathbf{y}$  есть

$$\mathbf{x}_{opt} = \mathbf{V}\mathbf{S}^+ \mathbf{U}^T \mathbf{y} = \mathbf{A}^+ \mathbf{y}, \quad \text{т.е.} \quad \mathbf{A}^+ = \mathbf{V}\mathbf{S}^+ \mathbf{U}^T. \quad (24)$$

Так как  $\mathbf{x}_{opt} = \mathbf{A}^+ \mathbf{y} \in \mathbf{R}(\mathbf{A}^T)$ , то  $\mathbf{A}\mathbf{x}_{opt} = \mathbf{A}\mathbf{A}^+ \mathbf{y} \in \mathbf{R}(\mathbf{A})$ . Следовательно,  $\mathbf{P}_r = \mathbf{A}\mathbf{A}^+$  является матрицей проектирования на подпространство  $\mathbf{R}(\mathbf{A})$ .

Матрица  $\mathbf{A}^+$  называется *обратной матрицей Мура – Пенроуза*, поскольку именно они ввели это понятие, или *обобщенной обратной* к  $\mathbf{A}$ . Следует отметить, что название “*обобщенная обратная*” применяется к другим матрицам, которые обладают не всеми свойствами матрицы  $\mathbf{A}^+$ . Именно поэтому чаще пользуются термином “*псевдообратная*”.

*Спектральной нормой* матрицы  $\mathbf{A}$  называется число, определяемое соотношением

$$\|\mathbf{A}\| = \max_{\mathbf{x} \neq \theta} (\|\mathbf{A}\mathbf{x}\| / \|\mathbf{x}\|) = \max_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|, \quad (25)$$

т.е. норма матрицы  $\mathbf{A}$  измеряет *коэффициент наибольшего изменения* длины векторов  $\mathbf{x}$  при умножении  $\mathbf{x}$  на эту матрицу.

Из определения нормы непосредственно следует, что:

- 1)  $\|c\mathbf{A}\| = |c|\|\mathbf{A}\|$  для всех вещественных  $c$  и всех  $\mathbf{A}$ ;
- 2)  $\|(\theta)\| = 0$  и  $\|\mathbf{A}\| > 0$ , если  $\mathbf{A} \neq (\theta)$  (здесь  $(\theta)$  означает нулевую матрицу).

Из определения нормы (25) также вытекает, что

$$\|\mathbf{A}\mathbf{x}\| \leq \|\mathbf{A}\| \|\mathbf{x}\| \quad (26)$$

для всех  $\mathbf{x}$  и  $\mathbf{A}$ , причем равенство достигается, по крайней мере, на одном ненулевом векторе.

Для того чтобы вычислить коэффициент наибольшего изменения в общем случае, возведем в квадрат обе части определяющего норму соотношения:

$$\|\mathbf{A}\|^2 = \max_{\|\mathbf{x}\|=1} \|\mathbf{Ax}\|^2 = \max_{\|\mathbf{x}\|=1} (\mathbf{x}^T \mathbf{A}^T \mathbf{Ax}) = \lambda_{\max}(\mathbf{A}^T \mathbf{A}),$$

где  $\lambda_{\max}(\mathbf{A}^T \mathbf{A})$  есть наибольшее собственное значение матрицы  $\mathbf{G} = \mathbf{A}^T \mathbf{A}$ . Таким образом

$$\|\mathbf{A}\| = \sqrt{\lambda_{\max}(\mathbf{A}^T \mathbf{A})}. \quad (27)$$

Поскольку квадраты сингулярных чисел  $\sigma_i$  матрицы  $\mathbf{A}$  являются (по построению) собственными числами матрицы  $\mathbf{A}^T \mathbf{A}$ , то другой формулой для нормы будет  $\|\mathbf{A}\| = \sigma_{\max}$ . Действительно, в соотношении  $\|\mathbf{Ax}\| = \|\mathbf{USV}^T \mathbf{x}\|$  ортогональные матрицы  $\mathbf{U}$  и  $\mathbf{V}$  не изменяют длины векторов и наибольший увеличивающий множитель равен наибольшему сингулярному числу  $\sigma_{\max}$ .

Вектор  $\mathbf{v}$ ,  $\|\mathbf{v}\|=1$ , который при отображении матрицей  $\mathbf{A}$  увеличивается максимальным образом, является собственным вектором матрицы  $\mathbf{A}^T \mathbf{A}$ , соответствующим наибольшему собственному значению:

$$\frac{\mathbf{v}^T \mathbf{A}^T \mathbf{A} \mathbf{v}}{\mathbf{v}^T \mathbf{v}} = \frac{\mathbf{v}^T \lambda_{\max} \mathbf{v}}{\mathbf{v}^T \mathbf{v}} = \lambda_{\max} = \|\mathbf{A}\|^2.$$

Норма Фробениуса (называется еще евклидовой матричной нормой – аналог евклидовой нормы векторов) матрицы  $\mathbf{A}$  определяется формулой

$$\|\mathbf{A}\|_F = \left( \sum_{i=1}^n \sum_{j=1}^k a_{ij}^2 \right)^{1/2}.$$

Спектральная норма и норма Фробениуса удовлетворяют соотношениям

$$\max_{i,j} |a_{ij}| \leq \|\mathbf{A}\| \leq \|\mathbf{A}\|_F \leq m^{1/2} \|\mathbf{A}\|, \quad (28)$$

где  $\mathbf{A}$  – матрица размера  $n \times k$ ,  $m = \min(n, k)$ .

Для произведения матриц справедливо неравенство:

$$\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|. \quad (29)$$



Действительно, по определению нормы имеем  $\|\mathbf{A}\mathbf{B}\mathbf{x}\| \leq \|\mathbf{A}\| \|\mathbf{B}\mathbf{x}\| \leq \|\mathbf{A}\| \|\mathbf{B}\| \|\mathbf{x}\|$ . Разделив это неравенство на  $\|\mathbf{x}\|$  и максимизируя, приходим к неравенству (29).

*Неравенство треугольника* для матричных норм: из неравенства треугольника для векторов имеем  $\|\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{x}\| \leq \|\mathbf{A}\mathbf{x}\| + \|\mathbf{B}\mathbf{x}\|$ . Разделив это неравенство на  $\|\mathbf{x}\|$  и максимизируя, получаем

$$\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|. \quad (30)$$

Обе названные нормы удовлетворяют мультипликативным неравенствам:

$$\|\mathbf{A}\mathbf{B}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|, \quad \|\mathbf{A}\mathbf{B}\|_F \leq \|\mathbf{A}\|_F \|\mathbf{B}\|_F.$$

Спектральная норма псевдообратной матрицы  $\mathbf{A}^+$  есть число

$$\|\mathbf{A}^+\| = 1/\sigma_{\min}, \quad \sigma_{\min} > 0.$$

*Числом обусловленности* (condition number) матрицы  $\mathbf{A}$  называется величина

$$c = \text{cond}(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^+\| = \sigma_{\max} / \sigma_{\min}. \quad (31)$$

Это число фигурирует в неравенстве для оценки относительной погрешности решения системы  $\mathbf{A}(\mathbf{x} + \delta\mathbf{x}) = \mathbf{y} + \delta\mathbf{y}$ , а именно

$$\|\delta\mathbf{x}\|/\|\mathbf{x}\| \leq c \|\delta\mathbf{y}\|/\|\mathbf{y}\|. \quad (32)$$

Действительно, так как  $\mathbf{A}\mathbf{x} = \mathbf{y}$ ,  $\mathbf{A}\delta\mathbf{x} = \delta\mathbf{y}$ , то непосредственно из (26) получаем, что

$$\|\mathbf{y}\| = \|\mathbf{A}\mathbf{x}\| \leq \|\mathbf{A}\| \|\mathbf{x}\| \quad \text{и} \quad \|\delta\mathbf{x}\| \leq \|\mathbf{A}^+\| \|\delta\mathbf{y}\|,$$

т.е.

$$\|\delta\mathbf{x}\|/\|\mathbf{x}\| \leq \|\mathbf{A}\| \|\mathbf{A}^+\| \|\delta\mathbf{y}\|/\|\mathbf{y}\|. \quad (33)$$

Из неравенства (33) следует, что для матриц с большим числом обусловленности относительная погрешность решения системы уравнений  $\mathbf{A}(\mathbf{x} + \delta\mathbf{x}) = (\mathbf{y} + \delta\mathbf{y})$  может значительно превышать относительную погрешность в правой части. Однако практическая ценность этого неравенства мала, поскольку неизвестными величинами являются числа  $\|\delta\mathbf{x}\|$ ,  $\|\mathbf{x}\|$ ,  $\|\delta\mathbf{y}\|$  и  $\|\mathbf{y}\|$ . Более содержательные результаты можно получить при вычислении оценок абсолютной погрешности решения. Именно оценки абсолютных погрешностей решений более всего интересуют специалистов-

прикладников. Эти вопросы подробно рассматриваются в следующем параграфе.

Сформулируем общий итог теоремы о сингулярном разложении матриц:

1.  $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^T$ ,  $\mathbf{U}^T\mathbf{A}\mathbf{V} = \mathbf{S}$ ,  $\mathbf{A}^+ = \mathbf{V}\mathbf{S}^+\mathbf{U}^T$ ,  $\mathbf{x}_{opt} = \mathbf{A}^+\mathbf{y}$ ,  $\mathbf{P}_r = \mathbf{A}\mathbf{A}^+$ ;
2. Количество отличных от нуля сингулярных чисел определяет ранг матрицы  $\mathbf{A}$ ;
3.  $\|\mathbf{A}\| = \sigma_{\max}$ ,  $\|\mathbf{A}^+\| = 1/\sigma_{\min}$ ,  $cond(\mathbf{A}) = \sigma_{\max} / \sigma_{\min}$ ;
4. Алгоритмы сингулярного разложения построены таким образом, что сингулярные числа вычисляются в невозрастающем порядке. Поэтому первые  $r$  столбцов матрицы  $\mathbf{U}$  образуют ортогональный базис подпространства  $\mathbf{R}(\mathbf{A})$ , остальные  $n - r$  столбцов образуют ортогональный базис подпространства  $\mathbf{N}(\mathbf{A}^T)$ ;
5. Первые  $r$  столбцов матрицы  $\mathbf{V}$  образуют ортогональный базис подпространства  $\mathbf{R}(\mathbf{A}^T)$ , остальные  $k - r$  столбцов образуют ортогональный базис подпространства  $\mathbf{N}(\mathbf{A})$ ;
6. Для матриц полного ранга наличие очень малых сингулярных чисел указывает на близость матрицы к вырождению;
7. Определитель невырожденной матрицы удовлетворяет условию  $|\det(\mathbf{A})| = |\det(\mathbf{U}) \det(\mathbf{S}) \det(\mathbf{V}^T)| = \sigma_1 \sigma_2 \dots \sigma_n$ , так как определители ортогональных матриц  $\mathbf{U}$  и  $\mathbf{V}$  равны  $\pm 1$ .

### §3. Спектральная классификация матриц произвольного размера и критерии плохой обусловленности

*Свет мой, зеркальце! Скажи  
Да всю правду доложи:  
Я ль на свете всех милее,  
Всех румяней и белее? ...*

*Пушкин А.С. Сказка о мертвой царевне  
и о семи богатырях.*

Характерной особенностью операций отображения конечномерных пространств с помощью матрицы  $\mathbf{A}$  является простота вычислений результатов отображений и возможность геометрической интерпретации этих результатов (для больших размерностей – на интуитивном уровне по аналогии с двух- и трехмерными пространствами).

Операция конечномерного отображения является *линейной* в том смысле, что

$$\mathbf{A}(t_1\mathbf{x}_1 + t_2\mathbf{x}_2) = t_1\mathbf{A}\mathbf{x}_1 + t_2\mathbf{A}\mathbf{x}_2 = t_1\mathbf{y}_1 + t_2\mathbf{y}_2, \quad t_1, t_2 \in \mathbf{R}. \quad (34)$$

В частности, прямолинейный отрезок

$$\mathbf{l}_x : (1-t)\mathbf{x}_1 + t\mathbf{x}_2, \quad t \in [0, 1], \quad (35)$$

отображается в прямолинейный отрезок

$$\mathbf{l}_y : \mathbf{A}\mathbf{l}_x = \mathbf{A}((1-t)\mathbf{x}_1 + t\mathbf{x}_2) = (1-t)\mathbf{A}\mathbf{x}_1 + t\mathbf{A}\mathbf{x}_2 = (1-t)\mathbf{y}_1 + t\mathbf{y}_2, \quad (36)$$

однако длины отрезков  $|\mathbf{l}_x|$  и  $|\mathbf{l}_y|$ , в общем случае, различны.

Деформация различных множеств или даже всего пространства при отображениях зависит от корпоративной структуры совокупности вектор-столбцов  $\mathbf{a}_i$ ,  $i = 1, 2, \dots, k$ . Выяснение особенностей картины деформаций при отображениях будет предметом нашего особого внимания.

**Пример 6.** Рассмотрим упоминавшуюся во введении матрицу

$$\mathbf{A}_t = \begin{pmatrix} 3 & 5 \\ 2 & 4 \end{pmatrix}.$$

Эта невырожденная матрица задает отображение  $\mathbf{R}^2 \Rightarrow \mathbf{R}^2$ :

$$\mathbf{y} = \mathbf{A}_t \mathbf{x} = x_1 \mathbf{a}_1 + x_2 \mathbf{a}_2 = x_1 \begin{pmatrix} 3 \\ 2 \end{pmatrix} + x_2 \begin{pmatrix} 5 \\ 4 \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}.$$

Так как  $\mathbf{A}_t \mathbf{e}_1 = \mathbf{a}_1$ ,  $\mathbf{A}_t \mathbf{e}_2 = \mathbf{a}_2$ , где  $\mathbf{e}_1$ ,  $\mathbf{e}_2$  – орты стандартного базиса, то на основании (36) получаем, что прямолинейный отрезок

$$\mathbf{l}_x : (1-t)\mathbf{e}_1 + t\mathbf{e}_2, \quad t \in [0, 1]$$

отображается в отрезок

$$\mathbf{l}_y : \mathbf{A}_t((1-t)\mathbf{e}_1 + t\mathbf{e}_2) = (1-t)\mathbf{a}_1 + t\mathbf{a}_2.$$

Соответственно, множество точек на сторонах квадрата  $Q_x = \{\mathbf{x} : |x_1| + |x_2| = 1\}$  отображается на стороны параллелограмма, вершины которого задаются векторами  $\{\mathbf{a}_1, \mathbf{a}_2, -\mathbf{a}_1, -\mathbf{a}_2\}$ . Таким образом, мы установили, что матрица  $\mathbf{A}_t$  при отображении деформирует квадрат  $Q_x$  в параллелограмм и осуществляет его поворот. Очевидно, что степень деформации и угол поворота зависят от корпоративной структуры вектор-столбцов матриц, а именно: *чем меньше углы между вектор-столбцами и чем больше различие в их длинах, тем больше эффект деформации. Эффект деформации отсутствует, когда  $\mathbf{A}$  – тождественное отображение.*

Топология эффекта деформации более содержательна, если рассмотреть отображение окружности  $\Omega_x = \{\mathbf{x} : \|\mathbf{x}\| = 1\}$ . Матрица  $\mathbf{A}_t$  отображает окружность  $\Omega_x$  в сильно вытянутый эллипс, длины полуосей  $l_i$  которого равны значениям сингулярных чисел матрицы  $\mathbf{A}_t$ , а именно:  $l_1 = \sigma_1 = 7.34$ ,  $l_2 = \sigma_2 = 0.27$ . Соответственно, матрица  $\mathbf{A}_t^+ = \mathbf{A}_t^{-1}$  отображает единичную окружность  $\Omega_y$  в эллипс, длины осей которого  $l_1 = 1/\sigma_1 \approx 0.14$  и  $l_2 = 1/\sigma_2 \approx 3.7$  (см. рис. 5б, стр. 37). Большое значение  $cond(\mathbf{A}_t) = \sigma_1/\sigma_2 = 26.96$  сигнализирует о большой деформации окружности при ее отображении матрицей  $\mathbf{A}_t$  и возможной плохой обусловленности системы уравнений  $\mathbf{A}_t(\mathbf{x} + \delta\mathbf{x}) = \mathbf{y} + \delta\mathbf{y}$ . Действительно, так как  $\sigma_2 < 1$  и  $\delta\mathbf{x} = \mathbf{A}_t^+ \delta\mathbf{y}$ , то  $\|\delta\mathbf{x}\| \leq \|\mathbf{A}_t^+\| \|\delta\mathbf{y}\| = (1/\sigma_2) \|\delta\mathbf{y}\|$  и, следовательно, длина погрешности

решения может увеличиваться в  $1/\sigma_2$  раз по сравнению с длиной вектора погрешности  $\delta\mathbf{y}$  при наиболее неблагоприятной его ориентации. На рис. 5а и 5б приводится геометрическое объяснение “парадоксального” результата, когда погрешность  $\|\delta\mathbf{y}_2\| < \|\delta\mathbf{y}_1\|$  приводит к погрешности  $\|\delta\mathbf{x}_2\| \gg \|\delta\mathbf{x}_1\|$ . Это обусловлено тем, что матрица  $\mathbf{A}_t^{-1}$  отображает единичную окружность в сплюснуто-вытянутый эллипс большой анизотропии ( $\text{cond}(\mathbf{A}_t) = 26.96$ ). Рассматриваемая система уравнений действительно является плохо обусловленной. В зависимости от ориентации вектора  $\delta\mathbf{y}$  (фиксированной длины), длина вектора  $\delta\mathbf{x} = \mathbf{A}_t^{-1}\delta\mathbf{y}$  может достигать максимальной величины  $\|\delta\mathbf{x}\|_{\max} = \|\mathbf{A}_t^+\| \|\delta\mathbf{y}\| = (1/\sigma_{\min}) \|\delta\mathbf{y}\| \gg \|\delta\mathbf{y}\|$ . В нашем случае вектор  $\mathbf{y} + \delta\mathbf{y}_2$  отображается с помощью матрицы  $\mathbf{A}^{-1}$  в точку  $\mathbf{x} + \delta\mathbf{x}_2$ , которая расположена относительно точки  $\mathbf{x} = (1, 1)^T$  значительно дальше, чем точка  $\mathbf{x} + \delta\mathbf{x}_1 = \mathbf{A}_t^{-1}(\mathbf{y} + \delta\mathbf{y}_1)$ . Такой поразительный эффект вызван именно тем, что направление  $\delta\mathbf{y}_1$  – наиболее благоприятное и совпадает с направлением первого собственного вектора ( $\lambda_1 \cong 53.88$ ) матрицы  $\mathbf{A}_t^T \mathbf{A}_t$ , а направление  $\delta\mathbf{y}_2$  – наиболее неблагоприятное и совпадает с направлением второго собственного вектора ( $\lambda_2 \cong 0.073$ ).

Следует заметить, что большое значение числа обусловленности матрицы не является признаком плохой обусловленности системы относительно абсолютной погрешности. Например, для матрицы  $\mathbf{D} = \text{diag}(1000, 10)$ ,  $\text{cond}(\mathbf{D}) = 100$ , однако  $\|\delta\mathbf{x}\| \leq \|\mathbf{D}^{-1}\| \|\delta\mathbf{y}\| = 0.1 \|\delta\mathbf{y}\|$ , т.е. система уравнений хорошо обусловлена.

**Определение.** Система алгебраических уравнений  $\mathbf{A}(\mathbf{x} + \delta\mathbf{x}) = \mathbf{y} + \delta\mathbf{y}$  *хорошо обусловлена*, если  $\|\delta\mathbf{x}\| \leq \|\delta\mathbf{y}_r\|$  и *плохо обусловлена*, если  $\|\delta\mathbf{x}\| > \|\delta\mathbf{y}_r\|$ , где  $\delta\mathbf{y}_r$  – проекция вектора  $\delta\mathbf{y}$  на ранговое пространство  $\mathbf{R}(\mathbf{A})$ .

Так как  $\|\delta\mathbf{x}\| \leq \|\mathbf{A}^+\| \|\delta\mathbf{y}\| = (1/\sigma_{\min}) \|\delta\mathbf{y}_r\|$ , то можно сформулировать эквивалентные критерии обусловленности:

$\sigma_{\min} \geq 1$  – система уравнений хорошо обусловлена;

$0 < \sigma_{\min} < 1$  – система уравнений плохо обусловлена.

Далее мы покажем, что эти критерии справедливы для матриц произвольного размера и ранга.

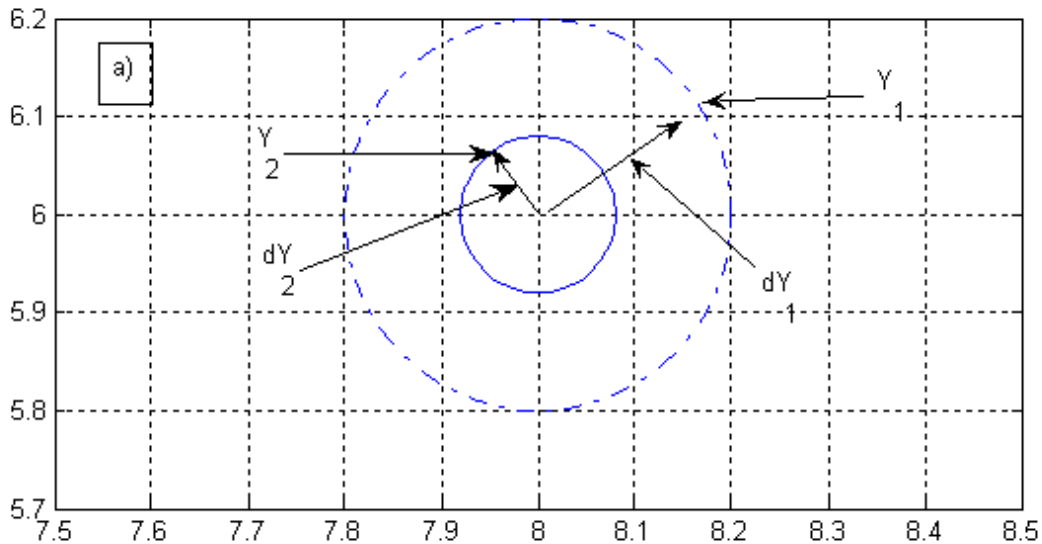


Рис. 5а. Погрешности  $\delta y_1$  и  $\delta y_2$  коллинеарны собственным векторам матрицы  $\mathbf{A}_t^T \mathbf{A}_t$ .

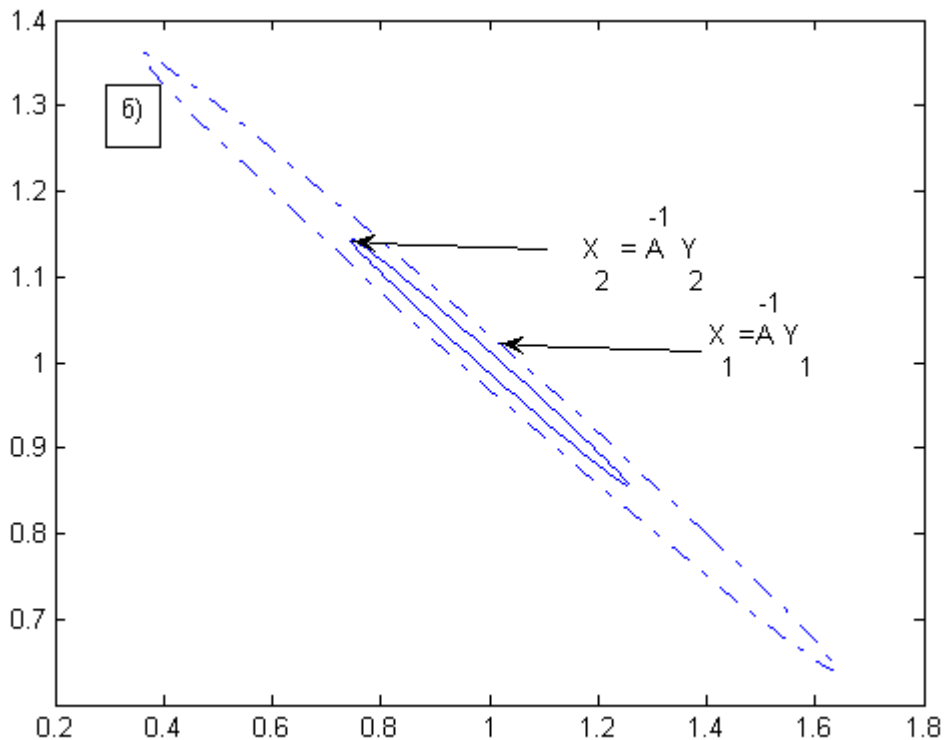


Рис. 5б. Сильная анизотропия при отображении обратной матрицей

приводит к существенному различию возмущенных решений.

Сингулярное разложение матриц позволяет получить геометрическую картину топологии отображений  $\mathbf{R}^k \Leftrightarrow \mathbf{R}^n$ .

**Теорема 12.** Матрица  $\mathbf{A}$ , произвольного размера  $n \times k$  и ранга  $r$ , отображает гиперсферу  $\Omega_{\mathbf{x}} = \{\mathbf{x} : \|\mathbf{x}\| = 1, \mathbf{x} \in \mathbf{R}^k\}$  в гиперэллипсоид  $\mathbf{E}_{\mathbf{y}} \in \mathbf{R}(\mathbf{A})$ , размерность которого равна  $r$ . Аналогично, матрица  $\mathbf{A}^+$  отображает гиперсферу  $\Omega_{\mathbf{y}} = \{\mathbf{y} : \|\mathbf{y}\| = 1, \mathbf{y} \in \mathbf{R}^n\}$  в гиперэллипсоид  $\mathbf{E}_{\mathbf{x}} \in \mathbf{R}(\mathbf{A}^T)$  размерность которого также равна  $r$ .

Доказательство. Рассмотрим *ортогональное* преобразование  $\mathbf{V}$  системы координат в пространстве  $\mathbf{R}^k = \mathbf{X}$ , в результате которого вектор  $\mathbf{x}$  получит новое представление  $\mathbf{x}'$ , где  $\mathbf{x} = \mathbf{V}\mathbf{x}'$ . Таким же образом, применяя другое ортогональное преобразование  $\mathbf{U}$  координат в  $\mathbf{R}^n = \mathbf{Y}$ , мы получим новые координаты  $\mathbf{y}$ , именно  $\mathbf{y}'$ , где  $\mathbf{y} = \mathbf{U}\mathbf{y}'$ . Если  $\mathbf{U}$  и  $\mathbf{V}$  – матрицы из сингулярного разложения матрицы  $\mathbf{A}$ , тогда

$$\mathbf{y}' = \mathbf{U}^T \mathbf{y} = \mathbf{U}^T \mathbf{A} \mathbf{x} = \mathbf{U}^T \mathbf{A} (\mathbf{V} \mathbf{x}') = (\mathbf{U}^T \mathbf{A} \mathbf{V}) \mathbf{x}' = \mathbf{S} \mathbf{x}'.$$

В новой ортогональной системе координат преобразование с помощью матрицы  $\mathbf{A}$  имеет очень простое представление:

$$\left\{ \begin{array}{l} y'_1 = \sigma_1 x'_1 \\ y'_2 = \sigma_2 x'_2 \\ \dots \\ y'_r = \sigma_r x'_r \\ y'_{r+1} = 0 \\ \dots \\ y'_n = 0 \end{array} \right. \quad (37)$$

Преобразование  $\mathbf{y}' = \mathbf{S} \mathbf{x}'$  просто отображает первую координатную ось пространства  $\mathbf{X}'$  в первую координатную ось пространства  $\mathbf{Y}'$  с коэффициентом растяжения  $\sigma_1 > 0$ . То же самое проделывается со 2-й, 3-й, ...,  $r$ -ой координатными осями пространства  $\mathbf{X}'$ , причем  $\sigma_2, \sigma_3, \dots, \sigma_r$  играют роль соответствующих коэффициентов растяже-

ния. Последующие  $(r+1)$ -я, ...,  $k$ -я координатные оси  $\mathbf{X}'$  отображаются в нулевой вектор  $\theta_n$  пространства  $\mathbf{Y}'$ .

Из соотношений (37) следует, что матрица  $\mathbf{S}$  отображает единичную гиперсферу  $\Omega_{\mathbf{x}} = \{\mathbf{x}' : \|\mathbf{x}'\| = 1\}$ ,  $\Omega_{\mathbf{x}} \in \mathbf{R}^k$  в  $r$ -мерный гиперэллипсоид  $\mathbf{E}_{\mathbf{y}}$ ,  $\mathbf{E}_{\mathbf{y}} \in \mathbf{R}(\mathbf{A})$  векторов  $\mathbf{y}'$  таких, что

$$\frac{(y'_1)^2}{\sigma_1^2} + \frac{(y'_2)^2}{\sigma_2^2} + \dots + \frac{(y'_r)^2}{\sigma_r^2} \leq 1 \quad \text{и} \quad y'_{r+1} = \dots = y'_n = 0, \quad (38)$$

т.е. длины полуосей определяются значениями сингулярных чисел  $\sigma_i > 0$ . Следует отметить, что в (38) стоит именно знак нестрого равенства  $\leq$ , так как проекцией гиперсферы на подпространство меньшей размерности является гипершар.

Как уже отмечалось ранее, численные алгоритмы сингулярного разложения позволяют вычислять сингулярные числа в *невозрастающем* порядке (т.е.  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ ). Поэтому в  $\mathbf{E}_{\mathbf{y}}$  наиболее удаленная точка имеет координаты  $\mathbf{y}'_{\max} = (\sigma_1, 0, \dots, 0)^T$ , а наименее удаленная точка имеет координаты  $\mathbf{y}'_{\min} = (0, \dots, \sigma_r, 0, \dots, 0)^T$ . Аналогично, матрица  $\mathbf{S}^+$  отображает единичную гиперсферу  $\Omega_{\mathbf{y}} = \{\mathbf{y} : \|\mathbf{y}\| = 1\}$ ,  $\mathbf{y} \in \mathbf{R}^n$  в  $r$ -мерный гиперэллипсоид  $\mathbf{E}_{\mathbf{x}} \in \mathbf{R}(\mathbf{A}^T)$ , длины полуосей которого равны  $1/\sigma_i$ ,  $\sigma_i > 0$ .

Свойство линейности операции отображения с помощью матрицы  $\mathbf{A}$  своеобразно проявляется в следующем. Семейство концентрических гиперсфер в  $\mathbf{R}^k$  с центром  $\mathbf{x}_c$  отображается в виде концентрических гиперэллипсоидов в  $\mathbf{R}^n$  с центром  $\mathbf{y}_c = \mathbf{A}\mathbf{x}_c$ . Поэтому действие матрицы  $\mathbf{A}$  можно интерпретировать как действие своеобразной “линзы”. Любопытно отметить, что если  $\mathbf{x}_c \neq \theta_k \in \mathbf{N}(\mathbf{A})$ , то “линза”  $\mathbf{A}$  концентрирует гиперэллипсоиды в  $\mathbf{R}^n$  вокруг начала координат. Соответствующим своеобразным свойством обладает и псевдообратная матрица  $\mathbf{A}^+$ , которая отображает концентрические гиперсферы с центром  $\mathbf{y}_c$  из  $\mathbf{R}^n$  в гиперэллипсоиды в  $\mathbf{R}(\mathbf{A}^T)$  с центром  $\mathbf{x}_c = \mathbf{A}^+\mathbf{y}_c$ . В обоих случаях



размерности этих гиперэллипсоидов определяются количеством ненулевых сингулярных чисел.

Важную информацию можно извлечь из анализа спектра сингулярных чисел. Всевозможные матрицы целесообразно разбить на 4 типа относительно распределения сингулярных чисел:

$$\left. \begin{array}{l} 1) \quad \sigma_i > 1, i = 1, 2, \dots, r, \\ 2) \quad \sigma_i = 1, i = 1, 2, \dots, r, \\ 3) \quad 0 < \sigma_{\min} < 1 < \sigma_{\max}, \\ 4) \quad 0 < \sigma_i < 1, i = 1, 2, \dots, r. \end{array} \right\} \quad (39)$$

Деформации единичных гиперсфер матрицами разных типов качественно существенно различаются.

Матрицы первого типа отображают гиперсферу  $\Omega_{\mathbf{x}} = \{\mathbf{x} : \|\mathbf{x}\| = 1\}$  в гиперэллипсоид, у которого длины каждой из полуосей больше единицы. Эти матрицы представляют оператор растяжения, и, соответственно,  $\mathbf{A}^+$  – оператор сжатия. Длина вектора погрешности решения для систем уравнений с такими матрицами удовлетворяет неравенству

$$\|\delta\mathbf{x}\| \leq \|\mathbf{A}^+\| \|\delta\mathbf{y}\| = (1/\sigma_{\min}) \|\delta\mathbf{y}_r\| < \|\delta\mathbf{y}_r\|$$

и, следовательно, системы уравнений с матрицами такого типа являются *хорошо-обусловленными*.

Очевидно, что ко второму типу относятся все ортогональные и полуортогональные матрицы, которые при отображении не изменяют длину векторов, так как  $\|\mathbf{Ax}\|^2 = (\mathbf{x}^T \mathbf{A}^T \mathbf{Ax}) = \|\mathbf{x}\|^2$ . Для таких матриц  $\|\mathbf{A}\| = \|\mathbf{A}^+\| = 1$ , поэтому справедливо равенство

$$\|\delta\mathbf{x}\| = \|\delta\mathbf{y}\|,$$

т.е. системы уравнений с матрицами такого типа – *нейтрально-обусловлены*.

Матрицы третьего типа в определенных направлениях осуществляют растяжение, а в других – сжатие. Для таких матриц справедливо неравенство

$$(1/\sigma_{\max}) \|\delta\mathbf{y}_r\| \leq \|\delta\mathbf{x}\| \leq (1/\sigma_{\min}) \|\delta\mathbf{y}_r\|.$$

При случайном направлении вектора погрешности  $\delta\mathbf{y}$  (фиксированной длины) существует такая неблагоприятная его ориентация,

при которой длина вектора погрешности решения может достигать значения

$$\|\delta \mathbf{x}\| = (1/\sigma_{\min}) \|\delta \mathbf{y}_r\| > \|\delta \mathbf{y}_r\|,$$

и поэтому следует признать, что системы уравнений с матрицами третьего типа являются *плохо-обусловленными*.

И, наконец, матрицы четвертого типа осуществляют сжатие (вообще говоря, с различными коэффициентами сжатия  $0 < \sigma_i < 1$ ). Соответственно  $\mathbf{A}^+$  осуществляет растяжение. Для таких матриц справедливо неравенство

$$(1/\sigma_{\max}) \|\delta \mathbf{y}_r\| \leq \|\delta \mathbf{x}\| \leq (1/\sigma_{\min}) \|\delta \mathbf{y}_r\|,$$

т.е. для всевозможных ориентаций вектора  $\delta \mathbf{y}$  имеем  $\|\delta \mathbf{x}\| > \|\delta \mathbf{y}_r\|$ , а это означает – матрицы четвертого типа заведомо *плохо-обусловлены*.

**Общий вывод: системы уравнений с матрицами первого и второго типов – хорошо-обусловлены, а с матрицами третьего и четвертого типов – плохо-обусловлены.**

На рис. ба ÷ бг показаны действия  $(2 \times 2)$ -матриц и обратных к ним всех типов, где отчетливо видна взаимно ортогональная ориентация эллипсов прямого и обратного отображений окружности.

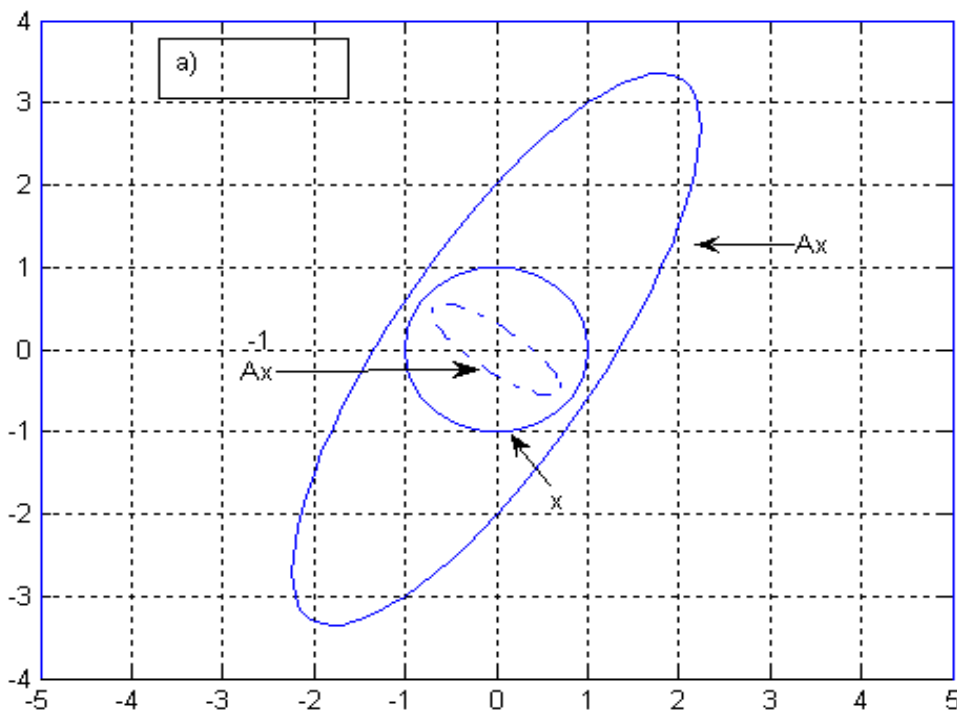


Рис. ба. Отображение окружности матрицей первого типа:  
 $\sigma_{\min} > 1$ ,  $cond(\mathbf{A}) = 3.3$ .

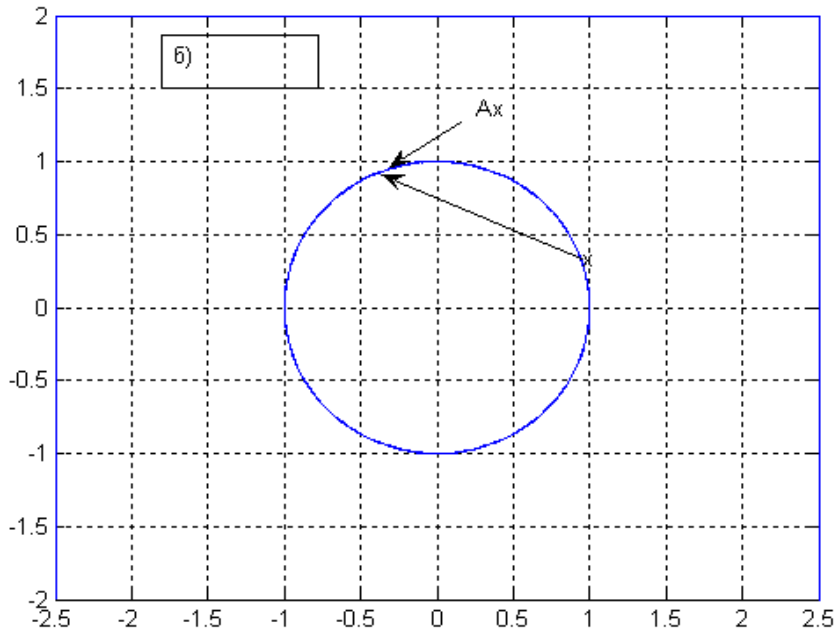


Рис. 6б. Отображение окружности ортогональной матрицей:  
 $\sigma_{\min} = \sigma_{\max} = 1, \text{cond}(\mathbf{A}) = 1.$

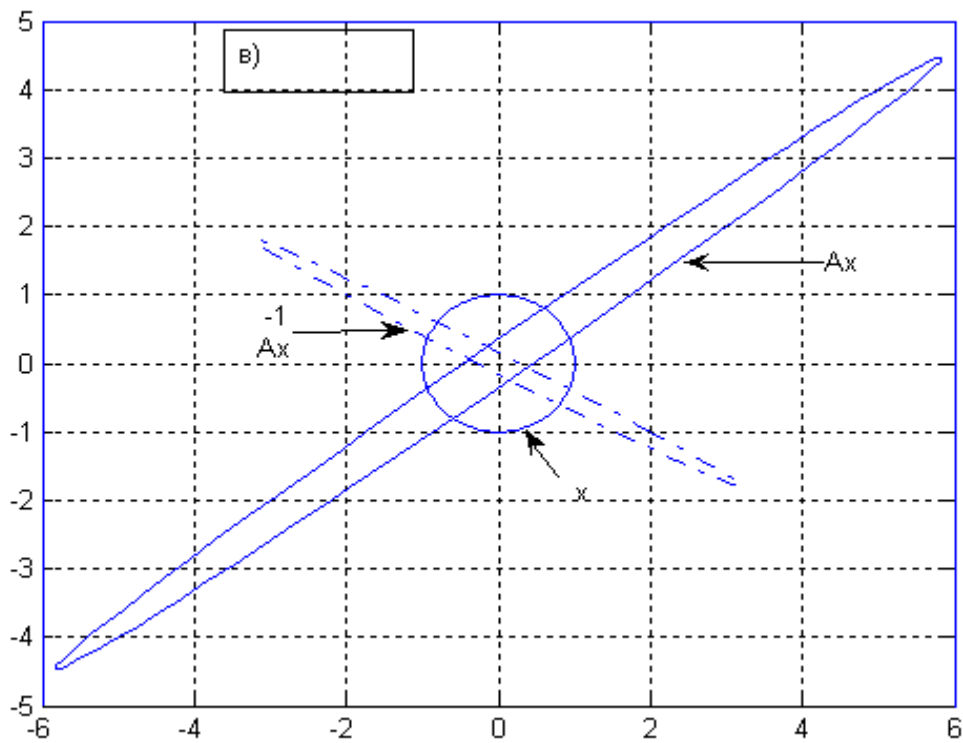


Рис. 6в. Отображение окружности матрицей третьего типа:  
 $0 < \sigma_{\min} < 1 < \sigma_{\max}, \text{cond}(\mathbf{A}) = 26.96.$

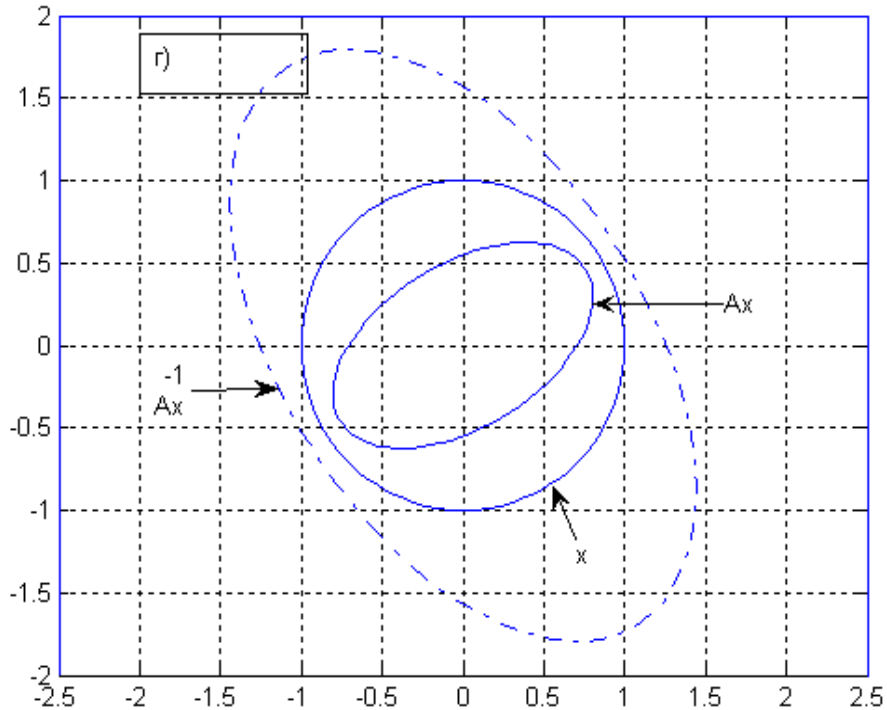


Рис. бг. Отображение окружности матрицей четвертого типа:  
 $0 < \sigma_{\min} < \sigma_{\max} < 1$ ,  $\text{cond}(\mathbf{A}) = 1.78$ .

При решении систем линейных уравнений  $\mathbf{A}(\mathbf{x} + \delta\mathbf{x}) = \mathbf{y} + \delta\mathbf{y}$ , содержащих погрешности в правой части, длина вектора погрешности решения удовлетворяет неравенству

$$\|\delta\mathbf{x}\| \leq \|\mathbf{A}^+\| \|\delta\mathbf{y}\| = (1/\sigma_{\min}) \|\delta\mathbf{y}_r\|, \quad (40)$$

при этом из (37) следует, что равенство достигается тогда, когда в новой ортогональной системе координат  $\delta\mathbf{y}' = (0, \dots, 0, \delta y_r, 0, \dots, 0)^T$ . Из неравенства (40) и предложенной классификации матриц (39) следует, что:

- а) для матриц первого и второго типов длина вектора погрешности решения никогда не превышает длины компоненты  $\delta y_r$  вектора погрешности правой части;
- б) для матриц третьего типа длина вектора погрешности решения может превышать или не превышать  $\|\delta\mathbf{y}_r\|$  (в зависимости от ориентации  $\delta\mathbf{y}$ );
- в) для матриц четвертого типа длина вектора погрешности решения всегда превышает  $\|\delta\mathbf{y}_r\|$ .

Таким образом, матрицы третьего и четвертого типов естественно называть *плохо-обусловленными*. Число  $q = 1/\sigma_{\min}$  – коэффициент максимального усиления длины вектора погрешности решения, а число  $cond(\mathbf{A}) = \sigma_{\max}/\sigma_{\min}$  является мерой *анизотропии* конечномерного отображения с помощью матрицы  $\mathbf{A}$ .

Вектор  $\delta\mathbf{y}$  в системе уравнений  $\mathbf{A}(\mathbf{x} + \delta\mathbf{x}) = \mathbf{y} + \delta\mathbf{y}$  является зачастую результатом погрешностей экспериментальных измерений и по своей природе является случайной величиной. Направление  $\delta\mathbf{y}$  непредсказуемо и, в общем случае,  $\delta\mathbf{y} = \delta\mathbf{y}_r + \delta\mathbf{y}_0$ ,  $\delta\mathbf{y}_r \in \mathbf{R}(\mathbf{A})$ ,  $\delta\mathbf{y}_0 \in \mathbf{N}(\mathbf{A}^T)$ . Так как  $\mathbf{A}^+\delta\mathbf{y}_0 \equiv \theta_n$ , то *радиусом доверительной области* (максимальное уклонение возмущенного решения  $\mathbf{x} + \delta\mathbf{x}$  от истинного) является величина

$$\rho = \|\delta\mathbf{x}\|_{\max} = (1/\sigma_{\min}) \|\delta\mathbf{y}_r\|. \quad (41)$$

В принципе, величину  $\|\delta\mathbf{y}\|$  можно оценить достаточно точно из условий проведения эксперимента, а  $\|\delta\mathbf{y}_0\| = \|\mathbf{A}\mathbf{x}_{opt} - (\mathbf{y} + \delta\mathbf{y})\|$  есть значение оптимальной невязки в задаче МНК. Так как  $\delta\mathbf{y}_r \perp \delta\mathbf{y}_0$ , то  $\|\delta\mathbf{y}_r\| = \sqrt{\|\delta\mathbf{y}\|^2 - \|\delta\mathbf{y}_0\|^2}$  и радиус доверительной области можно оценить более точно:

$$\rho = (1/\sigma_{\min}) \sqrt{\|\delta\mathbf{y}\|^2 - \|\delta\mathbf{y}_0\|^2}. \quad (42)$$

Специального рассмотрения заслуживает случай, когда в спектре распределения сингулярных чисел наблюдаются быстрый спад или пороговый нырок к нулю. В этом случае матрица  $\mathbf{A}$  является близкой к вырожденной и заведомо плохо-обусловленной. Наличие вычислительных погрешностей в значениях малых сингулярных чисел может привести к огромным погрешностям при вычислении псевдообратной матрицы  $\mathbf{A}^+$  и, соответственно, оптимальное решение  $\mathbf{x}_{opt} = \mathbf{A}^+\mathbf{y}$  будет содержать большие искажения. Разумным подходом является вычисление приближенного оптимального решения, применяя способ замены малых сингулярных чисел на некоторое пороговое значение  $\sigma_l$ . Замена истинной матрицы  $\mathbf{A}$  в системе уравнений  $\mathbf{A}\mathbf{x} = \mathbf{USV}^T \mathbf{x} = \mathbf{y}$  на приближенную  $\hat{\mathbf{A}} = \mathbf{U}\hat{\mathbf{S}}\mathbf{V}^T$ ,

где  $\hat{\mathbf{S}} = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_l, \sigma_l, \sigma_l, \dots, \sigma_l)$ , приведет к тому, что приближенное решение  $\mathbf{x}_{opt,app} = \hat{\mathbf{A}}^+ \mathbf{y}$  будет отличаться от  $\mathbf{x}_{opt} = \mathbf{A}^+ \mathbf{y}$ . Относительную погрешность решения  $\|\Delta \mathbf{x}\| / \|\mathbf{x}\|$  можно оценить, используя матрицы

$$\begin{cases} \hat{\mathbf{A}} = \mathbf{U} \hat{\mathbf{S}} \mathbf{V}^T, & \hat{\mathbf{S}} = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_l, \sigma_l, \sigma_l, \dots, \sigma_l) \\ \mathbf{A}_{0,l} = \mathbf{U} \mathbf{S}_{0,l} \mathbf{V}^T, & \mathbf{S}_{0,l} = \text{diag}(0, 0, \dots, 0, \sigma_l, \sigma_l, \dots, \sigma_l) \end{cases}. \quad (43)$$

Тогда

$$\|\Delta \mathbf{x}\| / \|\mathbf{x}\| \approx \|\mathbf{A}_{0,l}^+ \mathbf{y}\| / \|\hat{\mathbf{A}}^+ \mathbf{y}\|.$$

## Часть 2

### §4. Линейные математические модели: оптимальные параметры и радиус доверительной области

*... Я почувствовал, что уже увяз по пояс,  
а коня моего почти не было видно, торчали одни уши.  
Я крепко сжал бока лошади своими ногами,  
схватился рукой за свой собственный чуб и  
...представьте, вытащил себя вместе с конем  
из этого топкого болота!*

*Рудольф Распэ. Приключения барона Мюнхаузена*

При обработке и интерпретации экспериментальных данных возникает задача установления функциональной зависимости между группами переменных  $\mathbf{t} = (t_1, t_2, \dots, t_n)^T$  и  $\mathbf{f} = (f_1, f_2, \dots, f_n)^T$ . Мы будем считать независимыми переменными компоненты вектора  $\mathbf{t}$ , влияющими на значения компонент вектора  $\mathbf{f}$  по некоторому закону  $f = F(t, \mathbf{p}, \mathbf{q})$ . Функциональная зависимость  $F(t, \mathbf{p}, \mathbf{q})$ , принадлежащая некоторому параметрическому семейству функций относительно векторов параметров  $\mathbf{p}$  и  $\mathbf{q}$ , называется *математической моделью* взаимосвязи между наблюдаемыми величинами.

В самой общей постановке, задача установления функциональной взаимосвязи включает в себе две задачи:

1. *Выбор математической модели, которая адекватно описывает изучаемое явление;*
2. *Отыскание оптимальных параметров  $\mathbf{p}_{opt}$  и  $\mathbf{q}_{opt}$ , при которых данная математическая модель наилучшим образом соответствует опытными данным.*

Каких либо теоретических указаний на способ выбора адекватных моделей не существует, поэтому при выборе конкретной математической модели руководствуются априорными предположениями об изучаемом явлении.

Приведем конкретный пример. Пусть мы работаем с радиоактивным материалом. Тогда наблюдаемыми величинами  $f(t)$  являются показания счетчика Гейгера в различные моменты времени  $t$ . Пусть известно, что изучаемый материал представляет собой смесь двух

радиоактивных веществ с неизвестными периодами полураспада. Требуется определить периоды полураспадов, и в какой пропорции эти вещества смешаны. Тогда адекватная математическая модель, описывающая показания счетчика Гейгера имеет вид

$$f(t, p_1, p_2, q_1, q_2) = p_1 \exp(-q_1 t) + p_2 \exp(-q_2 t), \quad (44)$$

где параметры  $p_1, p_2, q_1, q_2$  подлежат определению на основании серии опытных измерений.

На практике радиоактивность измеряется в дискретные моменты времени и показания счетчика не будут точно соответствовать формуле (44), так как всегда присутствуют экспериментальные погрешности измерений. Вместо этого мы имеем серию показаний  $f_1, f_2, \dots, f_n$  в различные моменты времени  $t_1, t_2, \dots, t_n$  и соотношение (44) выполняются лишь приближенно:

$$\begin{aligned} p_1 \exp(-q_1 t_1) + p_2 \exp(-q_2 t_1) &\approx f_1 \\ p_1 \exp(-q_1 t_2) + p_2 \exp(-q_2 t_2) &\approx f_2 \\ &\dots\dots\dots \\ p_1 \exp(-q_1 t_n) + p_2 \exp(-q_2 t_n) &\approx f_n. \end{aligned} \quad (45)$$

В общем случае система уравнений (45) является несовместной и тогда, используя метод наименьших квадратов, возникает задача поиска оптимальных значений параметров  $p_1, p_2, q_1, q_2$ , при которых достигается минимальное расхождение между теоретическими и наблюдаемыми значениями  $\mathbf{f}$ , т.е. решается задача

$$\Phi(p_1, p_2, q_1, q_2) = \sum_{i=1}^n \left[ f_i^{teor}(p_1, p_2, q_1, q_2, t_i) - f_i^{exp} \right]^2 \rightarrow \min. \quad (46)$$

Параметры  $q_1$  и  $q_2$  в математической модели (44) входят нелинейным образом. Задача поиска глобального минимума нелинейного функционала (46) связана с определенными вычислительными трудностями и для ее решения требуется применение специальных вычислительных методов.

Мы остановимся на рассмотрении отдельного класса задач МНК, когда математические модели выбираются в виде линейной комбинации некоторого количества базисных функций, а именно

$$f(t) = \sum_{j=1}^k p_j \varphi_j(t), \quad (47)$$

где  $p_j$  – компоненты вектора искомых параметров.

Математические модели вида (47) называются *линейными математическими моделями*. Для таких моделей уравнения (45)



записываются в виде системы линейных алгебраических уравнений  $\mathbf{Fp} = \mathbf{f}$  :

$$\left\{ \begin{array}{l} p_1\varphi_1(t_1) + p_2\varphi_2(t_1) + \dots + p_k\varphi_k(t_1) = f(t_1) \\ p_1\varphi_1(t_2) + p_2\varphi_2(t_2) + \dots + p_k\varphi_k(t_2) = f(t_2) \\ \dots\dots\dots\dots\dots\dots\dots\dots\dots\dots\dots\dots\dots\dots\dots\dots\dots\dots \\ p_1\varphi_1(t_n) + p_2\varphi_2(t_n) + \dots + p_k\varphi_k(t_n) = f(t_n) \end{array} \right. . \quad (48)$$

Матрица  $\mathbf{F}$ , размера  $n \times k$ ,  $n > k$ , имеет вид

$$\mathbf{F} = \begin{pmatrix} \varphi_1(t_1) & \varphi_2(t_1) & \dots & \varphi_k(t_1) \\ \varphi_1(t_2) & \varphi_2(t_2) & \dots & \varphi_k(t_2) \\ \dots & \dots & \dots & \dots \\ \varphi_1(t_n) & \varphi_2(t_n) & \dots & \varphi_k(t_n) \end{pmatrix}$$

и называется *матрицей плана*. Такое название связано с тем, что численные значения элементов матрицы  $\mathbf{F}$  определяются запланированным выбором базисных функций и точек измерений  $t_i$  наблюдаемой зависимости  $f(t_i) = f_i$ .

В общем случае, когда наблюдаемые значения  $f_i$  содержат погрешности, переопределенная система уравнений (48) является несовместной. Это означает, что вектор  $\mathbf{f}$  не принадлежит ранговому пространству, т.е.  $\mathbf{f} = \mathbf{f}_r + \mathbf{f}_0$ ,  $\mathbf{f}_r \in \mathbf{R}(\mathbf{F})$ ,  $\mathbf{f}_0 \in \mathbf{N}(\mathbf{F}^T)$ ,  $\mathbf{f}_0 \neq \theta_n$ . Тогда оптимальное решение в смысле наименьших квадратов системы  $\mathbf{Fp} = \mathbf{f}$  вычисляется по формуле

$$\mathbf{p}_{opt} = \mathbf{F}^+ \mathbf{f},$$

где  $\mathbf{F}^+$  – псевдообратная матрица Мура – Пенроуза, а радиус доверительной области искомым параметров оценивается по формуле

$$\rho = \|\delta \mathbf{p}\|_{\max} = (1 / \sigma_{\min}) \sqrt{\|\delta \mathbf{f}\|^2 - \|\delta \mathbf{f}_0\|^2},$$

где:  $\|\delta \mathbf{f}\|$  – длина вектора погрешности, оцениваемая из условий эксперимента,  $\|\delta \mathbf{f}_0\| = \|\mathbf{Fp}_{opt} - \mathbf{f}\|$  – значение оптимальной невязки в задаче МНК.

Рассмотрим несколько популярных математических моделей.

1. Математическая модель вида  $f(t) = p_1 + p_2 t$  является простейшей полиномиальной моделью. Поскольку базисными функциями является набор  $\{\varphi_1(t) = 1, \varphi_2(t) = t\}$ , то для такой модели матрица плана имеет вид:

$$\mathbf{F} = \begin{pmatrix} 1 & t_1 \\ 1 & t_2 \\ \dots & \dots \\ 1 & t_n \end{pmatrix}.$$

Квадраты сингулярных чисел матрицы  $\mathbf{F}$  есть собственные значения матрицы нормальных уравнений

$$\mathbf{G} = \mathbf{F}^T \mathbf{F} = \begin{pmatrix} n & \sum_{i=1}^n t_i \\ \sum_{i=1}^n t_i & \sum_{i=1}^n t_i^2 \end{pmatrix}.$$

Докажем, что для рассматриваемой математической модели можно применить специальное преобразование, при котором матрица  $\mathbf{F}$  будет хорошо обусловленной. Для этого рассмотрим математическую модель  $f_c(\tau) = q_1 + q_2 \tau$ , где  $\tau = w(t - t_c)$ ,  $t_c = (1/n) \sum_{i=1}^n t_i$  – геометри-

ческий центр точек  $\{t_1, t_2, \dots, t_n\}$ . Так как  $\sum_{i=1}^n \tau_i = \sum_{i=1}^n w(t_i - t_c) = 0$ ,

то рассматриваемая математическая модель генерирует матрицу  $\mathbf{F}_c$  такую, что матрица соответствующей системы нормальных уравнений принимает диагональный вид

$$\mathbf{G}_c = \mathbf{F}_c^T \mathbf{F}_c = \begin{pmatrix} n & 0 \\ 0 & \sum_{i=1}^n \tau_i^2 \end{pmatrix}.$$

Масштабирующий множитель  $w$  необходимо выбрать так, чтобы выполнялось неравенство  $\sum_{i=1}^n \tau_i^2 > 1$ , и тогда собственные значения матрицы и, следовательно, сингулярные числа матрицы  $\mathbf{F}_c$  будут больше единицы, т.е. матрица  $\mathbf{F}_c$  принадлежит к классу хорошо-

обусловленных матриц. Легко убедиться, что между параметрами  $p_1$ ,  $p_2$  и  $q_1$ ,  $q_2$  существует простая взаимосвязь:

$$p_1 = q_1 - q_2 w t_c, \quad p_2 = q_2 w. \quad (49)$$

2) Математическая модель  $f(t) = p_1 \exp(p_2 t)$  редуцируется к простейшей линейной модели  $g(t) = \ln f(t) = q_1 + q_2 t$ , где  $q_1 = \ln p_1$ ,  $q_2 = p_2$ .

3) Математическая модель  $f(t) = p_1 \exp(-p_2(t_c - t)^2)$  редуцируется к параболической модели  $g(\tau) = \ln f(\tau) = q_1 + q_2 \tau^2$ , где  $\tau = t_c - t + \tau_s$ ,  $q_1 = \ln p_1$ ,  $q_2 = -p_2$ . Матрица плана такой модели имеет вид:

$$\mathbf{F} = \begin{pmatrix} 1 & \tau_1^2 \\ 1 & \tau_2^2 \\ \dots & \dots \\ 1 & \tau_n^2 \end{pmatrix},$$

соответствующая ей матрица нормальных уравнений

$$\mathbf{G} = \mathbf{F}^T \mathbf{F} = \begin{pmatrix} n & \sum_{i=1}^n \tau_i^2 \\ \sum_{i=1}^n \tau_i^2 & \sum_{i=1}^n \tau_i^4 \end{pmatrix}.$$

Для заданной серии точек  $\{t_i\}$  можно выбрать такое  $\tau_s$ , чтобы добиться хорошей обусловленности матрицы  $\mathbf{F}$ . Другими словами, необходимо выбирать такое значение  $\tau_s$ , чтобы для собственных значений матрицы  $\mathbf{G}$  выполнялось условие  $\lambda_{\min} > 1$ .

4) Математическая модель в виде полинома степени  $k$

$$f(t) = \sum_{j=0}^k p_j t^j \quad (50)$$

генерирует матрицу плана:

$$\mathbf{F} = \begin{pmatrix} 1 & t_1 & t_1^2 & \dots & t_1^k \\ 1 & t_2 & t_2^2 & \dots & t_2^k \\ \cdot & \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \cdot & \dots & \cdot \\ 1 & t_n & t_n^2 & \dots & t_n^k \end{pmatrix}.$$

Спектр сингулярных чисел матрицы  $\mathbf{F}$  существенным образом зависит от распределения точек  $\{t_i\}$ , поэтому матрица  $\mathbf{F}$  может с большой вероятностью принадлежать к классу плохо-обусловленных матриц. Это связано с тем, что для базисных функций  $\{t^i, i = 0, 1, \dots, k, k \geq 3\}$  не существует конечного интервала, на котором эти функции были бы попарно ортогональными. Именно поэтому аппроксимация функций с помощью алгебраического полинома высокого порядка не применяется в практических приложениях.

Примером задачи с плохой обусловленностью является следующая задача.

5) Полиномиальная аппроксимация функций по значениям алгебраических моментов.

Пусть на отрезке  $[0, 1]$  задана функция  $f(t)$  и известно множество моментов

$$m_i = \int_0^1 t^i f(t) dt, \quad i = 0, 1, 2, \dots, k.$$

(Значения  $m_i$ , полученные путем численного интегрирования, возможно, содержат некоторые погрешности). Ставится задача аппроксимации функции  $f(t)$  алгебраическим полиномом  $P(t)$  степени  $k$ :

$$p_0 + p_1 t + p_2 t^2 + \dots + p_k t^k = f(t). \quad (51)$$

Умножая (51) последовательно на  $t^j$ ,  $j = 0, 1, 2, \dots, k$ , и интегрируя на  $[0, 1]$ , получаем систему уравнений  $\mathbf{Gp} = \mathbf{m}$  с матрицей Гильберта (элементы матрицы вычисляются по формуле  $g_{ij} = 1/(i + j - 1)$ ):

$$\mathbf{G} = \begin{pmatrix} 1 & 1/2 & 1/3 & \dots & 1/(k+1) \\ 1/2 & 1/3 & 1/4 & \dots & 1/(k+2) \\ 1/3 & 1/4 & 1/5 & \dots & 1/(k+3) \\ \cdot & \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \cdot & \dots & \cdot \\ 1/(k+1) & 1/(k+2) & 1/(k+3) & \dots & 1/(2k+1) \end{pmatrix}.$$

Поскольку углы между вектор-столбцами матрицы  $\mathbf{G}$  малы, то следует ожидать, что матрица плохо-обусловлена. На самом деле, матрица Гильберта относится к числу плохо-обусловленных матриц третьего типа. Характерная особенность этой матрицы в том, что при возрастании ее порядка минимальные собственные числа (сингулярные числа) очень быстро стремятся к нулю. Матрица Гильберта является плохо-обусловленной и уже для  $k > 10$  ничтожно малые погрешности в правой части системы  $\mathbf{G}(\mathbf{p} + \delta\mathbf{p}) = \mathbf{m} + \delta\mathbf{m}$  могут вызвать гигантское увеличение  $\|\delta\mathbf{p}\|$ . Глубинная причина этого эффекта в том же: для базисных функций  $\{t^i, i = 0, 1, \dots, k\}$  не существует интервала, на котором даже первые три функции были бы ортогональны. Именно поэтому в практических приложениях не используется аппроксимация функций алгебраическим многочленом порядка  $k > 3$ . Существенно предпочтительнее использовать ортогональные многочлены (например, многочлены Лежандра, Чебышева и др.).

Приведем пример использования ортогональных функций.

б) Аппроксимация функций тригонометрическим многочленом.

Пусть на отрезке  $[0, 1]$  задана функция  $f(t)$  и известно множество коэффициентов Фурье

$$c_l = \int_0^1 \exp(i2\pi lt) f(t) dt, \quad l = 0, 1, 2, \dots, k.$$

Численные значения  $c_l$  возможно содержат некоторые погрешности. Ставится задача аппроксимации функции  $f(t)$  тригонометрическим многочленом  $T_k(t)$  степени  $k$ :

$$p_0 + p_1 \exp(i2\pi t) + p_2 \exp(i2\pi 2t) + \dots + p_k \exp(i2\pi kt) = f(t). \quad (52)$$

Функции  $\{\exp(i2\pi lt), l = 0, 1, 2, \dots, k\}$  – ортогональны на интервале  $[0, 1]$ , т.е.  $\int_0^1 \exp(i2\pi(j-l)t) dt = \delta_{jl}$ . Умножая (52) последовательно

на  $\exp(-i2\pi lt)$ ,  $l = 0, 1, 2, \dots, k$ , и интегрируя на  $[0, 1]$ , получаем систему уравнений

$$\mathbf{E}\mathbf{p} = \mathbf{c}^*, \quad (53)$$

где:  $\mathbf{E}$  – единичная матрица,  $\mathbf{c}^*$  – вектор комплексно сопряженных коэффициентов Фурье. Из (53) следует, что  $\mathbf{p} \equiv \mathbf{c}^*$  и  $\delta\mathbf{p} \equiv \delta\mathbf{c}^*$ , т.е. задача аппроксимации функции тригонометрическим многочленом является *нейтрально-обусловленной*.

В общем случае, доверительный радиус оптимальных параметров выбранной (возможно адекватной) линейной математической модели зависит от обусловленности соответствующей матрицы плана  $\mathbf{F}$  и, конечно, от точности экспериментальных данных. Отметим, что по визуальному анализу графического изображения дискретного набора  $\{t_i, f(t_i)\}$  невозможно судить о качестве эксперимента. Действительно, если вектор погрешностей  $\delta\mathbf{f} = \delta\mathbf{f}_r + \delta\mathbf{f}_0$ , где  $\delta\mathbf{f}_r$  принадлежит ранговому пространству  $\mathbf{R}(\mathbf{F})$ , а  $\delta\mathbf{f}_0$  лежит в нуль-пространстве  $\mathbf{N}(\mathbf{F}^T)$  и тогда  $\mathbf{F}^+ \delta\mathbf{f}_0 \equiv \theta_k$ , то для двух серий измерений  $\{t_i, f_1(t_i)\}$  и  $\{t_i, f_2(t_i)\}$  таких, что  $\delta_{1,r}\mathbf{f} = \delta_{2,r}\mathbf{f}$ , а  $\delta_{1,0}\mathbf{f} \neq \delta_{2,0}\mathbf{f}$  имеем  $\|\mathbf{f} + \delta_1\mathbf{f}\| \neq \|\mathbf{f} + \delta_2\mathbf{f}\|$ , однако оценки оптимальных параметров в обоих случаях тождественно совпадают. Такие серии экспериментальных измерений являются *эквивалентными* в смысле оценки оптимальных параметров по МНК. Приведенные рассуждения можно подтвердить на конкретном примере параболической зависимости  $f(t) = 7 - 4t + t^2$ . Две серии погрешностей “экспериментальных” данных моделировались с помощью датчика случайных чисел, распределенных по нормальному закону с различными дисперсиями. Затем, после разложения векторов погрешностей на ортогональные составляющие, образованы новые векторы погрешностей  $\delta_1\mathbf{f}$  и  $\delta_2\mathbf{f}$  такие, что  $\delta_{1,r}\mathbf{f} = \delta_{2,r}\mathbf{f}$ , а  $\delta_{1,0}\mathbf{f} \neq \delta_{2,0}\mathbf{f}$ . Численные значения “экспериментальных” данных изображены графически на рис. 7 и приведены в таб. 1. Обе серии “экспериментальных” значений дают одинаковые оптимальные параметры

$$\mathbf{p}_{1,opt} = \mathbf{p}_{2,opt} = (7.888, -4.539, 1.018)^T,$$

хотя визуальный анализ графического изображения дает предпочтение первой серии измерений (на рис. 7 эти данные изображены

(—\*—) линией). В принципе возможна ситуация, когда  $\|\delta_1 \mathbf{f}\| \gg \|\delta_2 \mathbf{f}\|$ , но  $\|\delta_{1,r} \mathbf{f}\| \ll \|\delta_{2,r} \mathbf{f}\|$ . В этом случае оценка оптимальных параметров по данным первой серии измерений будет более точной. Именно поэтому мерой качества эксперимента является *радиус доверительной области* оптимальных параметров.

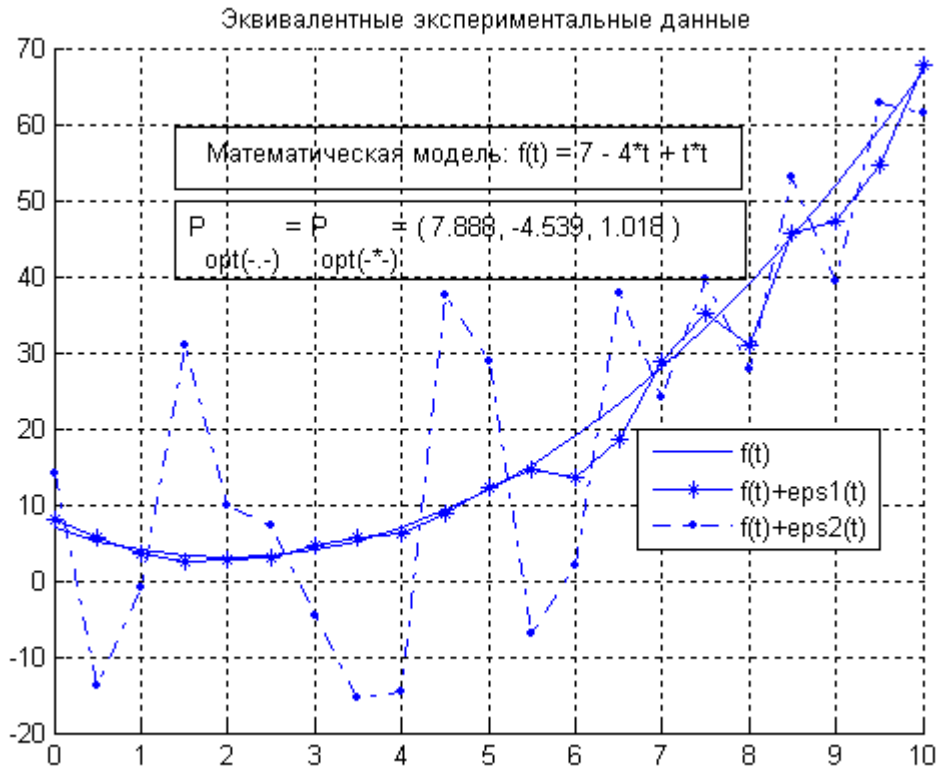


Рис. 7. Две серии эквивалентных “экспериментальных” данных с различными уровнями погрешностей.

Таб. 1. Эквивалентные “экспериментальные”  
 данные для параболической модели  $f(t) = p_1 + p_2t + p_3t^2$

$i$	$t_i$	$f_i + \delta_1 f_i$	$f_i + \delta_2 f_i$
1	0.0000	8.0013	14.0312
2	0.5000	5.7449	-13.6938
3	1.0000	3.5893	-0.9193
4	1.5000	2.4487	30.8792
5	2.0000	2.8870	9.7729
6	2.5000	3.0828	7.2042
7	3.0000	4.5480	-4.7272
8	3.5000	5.7406	-15.4413
9	4.0000	6.0883	-14.6353
10	4.5000	8.7071	37.3732
11	5.0000	12.3312	28.9211
12	5.5000	14.5882	-6.9125
13	6.0000	13.4434	2.0876
14	6.5000	18.4239	37.8237
15	7.0000	28.8513	23.9817
16	7.5000	35.1849	39.6465
17	8.0000	30.8896	27.8764
18	8.5000	45.6048	53.0226
19	9.0000	47.1660	39.3771
20	9.5000	54.6112	62.7062
21	10.0000	67.8583	61.4169



## §5. Алгоритмы численного решения интегрального уравнения типа свертки

*Только Истина есть непреложный закон,  
Этот мир Бытия лишь ему подчинен.  
Все что есть – быть должно в Бытии и Вселенной,  
То, что быть не должно, то мираж или сон ...*

*Омар Хайям. Рубаи*

Определим линейный оператор, осуществляющий свертывание функции  $h(t)$  с сигналом произвольной формы  $f(t)$ . Мы будем предполагать, что  $h(t)$  описывает реакцию некоторого приемного устройства на единичное импульсное воздействие,  $f(t)$  – входной сигнал. В результате воздействия приемного устройства на входной сигнал получается выходной сигнал  $g(t)$ , который называется *сверткой* и определяется следующим образом.

Пусть даны две функции  $f(t)$  и  $h(t)$  с областью определения  $-\infty < t < \infty$  и такие, что существует интеграл свертки

$$g(t) = \int_{-\infty}^{\infty} f(t - \tau)h(\tau)d\tau. \quad (54)$$

Функцию  $g(t)$  называют еще *конволюционным произведением*, а операция свертки кратко обозначается  $g = f \otimes h$ . Легко убедиться, что операция свертки обладает свойством коммутативности, т.е.

$$g(t) = \int_{-\infty}^{\infty} f(\tau)h(t - \tau)d\tau. \quad (55)$$

Как следует из определения свертки, вычисление  $g(t)$  производится следующим образом. Функцию  $f(-\tau)$  необходимо сместить на величину  $t$ . Площадь фигуры, полученная в результате перемножения функций  $f(t - \tau)$  и  $h(\tau)$ , дает значение функции  $g(t)$ , определяемое интегралом (54). Заметим, что использование в формуле (54) инвертированной функции  $f(-\tau)$  отражает тот факт, что на приемное устройство в момент времени  $t$  поступил сигнал  $f(t - \tau)$ .

Рассмотрим некоторые свойства свертки функций, предполагая осуществимыми все операции, которые мы будем производить.

Определим:

$$S_f = \int_{-\infty}^{\infty} f(t)dt \quad , \quad c_f = \frac{1}{S_f} \int_{-\infty}^{\infty} tf(t)dt \quad , \quad w_f^2 = \frac{1}{S_f} \int_{-\infty}^{\infty} (t - c_f)^2 f(t)dt .$$

Величину  $S_f$  называют *площадью*,  $c_f$  – *центром*, а  $w_f$  – *шириной* функции  $f(t)$ . Рассмотрим, как изменяются эти величины при операциях свертки.

Свойство 1. При свертке функций их площади *перемножаются*. Действительно,

$$S_g = \int_{-\infty}^{\infty} g(t)dt = \int_{-\infty}^{\infty} dt \int_{-\infty}^{\infty} f(\tau)h(t-\tau)d\tau = \int_{-\infty}^{\infty} f(\tau) \left[ \int_{-\infty}^{\infty} h(t-\tau)dt \right] d\tau = S_f S_h .$$

Свойство 2. При свертке функций их центры *складываются*:

$$\begin{aligned} c_g &= \frac{1}{S_f S_h} \int_{-\infty}^{\infty} tg(t)dt = \frac{1}{S_f S_h} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} tf(\tau)h(t-\tau)dtd\tau = \\ &= \frac{1}{S_f S_h} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \tau f(\tau)h(t-\tau)dtd\tau + \frac{1}{S_f S_h} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (t-\tau) f(\tau)h(t-\tau)dtd\tau = \\ &= c_f + c_h . \end{aligned}$$

Свойство 3. Ширина  $w_g$  свертки функций  $f(t)$  и  $h(t)$ :

$$w_g = (w_f^2 + w_h^2)^{1/2} .$$

Из определения свертки (54) непосредственно следует, что если сигналы  $f(t)$  и  $h(t)$  имеют финитные длительности  $T_1$  и  $T_2$ , то длительность выходного сигнала  $g(t)$  равна  $T_1 + T_2$ . В общем случае для дефинитных функций также можно говорить об их ширине. Тогда, в предположении о существовании соответствующих интегралов, имеем

$$\begin{aligned} w_g^2 &= \frac{1}{S_f S_h} \int_{-\infty}^{\infty} (t - c_g)^2 g(t)dt = \frac{1}{S_f S_h} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (t - c_f - c_h)^2 f(\tau)h(t-\tau)dtd\tau = \\ &= \frac{1}{S_f S_h} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} [(\tau - c_f) + (t - \tau - c_h)]^2 f(\tau)h(t-\tau)dtd\tau = w_f^2 + w_h^2 . \end{aligned}$$

Здесь использованы равенства:

$$\int_{-\infty}^{\infty} (\tau - c_f) f(\tau) d\tau = 0 \quad \text{и} \quad \int_{-\infty}^{\infty} (t - \tau - c_h) h(\tau) d\tau = 0.$$

Итак,

$$w_g = (w_f^2 + w_h^2)^{1/2},$$

т.е. в результате свертки происходит неизбежное *уширение* сигнала.

**Пример.** Свертка функции с прямоугольным импульсом

$$\Pi_a(t) = \begin{cases} 1, & |t| \leq a \\ 0, & |t| > a \end{cases}.$$

Легко убедиться, что

$$g(t) = \frac{1}{2a} \Pi_a(t) \otimes f(t) = \frac{1}{2a} \int_{t-a}^{t+a} f(\tau) d\tau$$

представляет собой *скользящее среднее значение* функции  $f(t)$ . Скользящим усреднением достигается определенный эффект сглаживания осцилляций функции, но, в силу Свойства 3, ценой неизбежного увеличения ее ширины.

Множественная свертка прямоугольного импульса

$$\Pi_{n,a}(t) = \underbrace{\Pi_a(t) \otimes \Pi_a(t) \otimes \dots \otimes \Pi_a(t)}_n$$

является гладкой функцией с непрерывными производными до  $(n-1)$ -го порядка. Так, например,  $\Pi_{3,a}(t)$  состоит из трех гладко сшитых параболических отрезков:

$$\Pi_{3,a}(t) = \begin{cases} 0.5(3a+t)^2, & -3a \leq t \leq -a \\ 3a^2 - t^2, & |t| \leq a \\ 0.5(3a-t)^2, & a \leq t \leq 3a \end{cases}.$$

Таким образом, многократное применение скользящего усреднения с прямоугольным окном фиксированной ширины эквивалентно однократному скользящему усреднению с более гладким окном. Этот способ может применяться для сглаживания экспериментальных данных, измеренных на равномерной, но достаточно плотной сетке значений  $\{t_i\}$ , и содержащих малые случайные погрешности.

Для дефинитных функций наиболее эффективный способ решения интегрального уравнения (55) сводится к применению преобразова-

ния Фурье. В тоже время, для финитных функций можно проводить численные решения, применяя дискретизацию на равномерной сетке.

Для функций, заданных на дискретной и равномерной сетке, корректное определение свертки вводится следующим образом. Пусть  $\{f_1, f_2, \dots, f_m\}$  и  $\{h_1, h_2, \dots, h_n\}$  – дискретные выборки с постоянным шагом значений функций  $f(t)$  и  $h(t)$ . Объем этих выборок может быть различным. Образуем векторы размерности  $k = m + n$  следующим образом:

$$\mathbf{f}_{+0} = (f_1, f_2, \dots, f_m, \underbrace{0, 0, \dots, 0}_n)^T \text{ и } \mathbf{h}_{+0} = (h_1, h_2, \dots, h_n, \underbrace{0, 0, \dots, 0}_m)^T.$$

Необходимость дополнения нулевыми компонентами вызвана тем, что при свертке результирующий объем выборки  $k = m + n$ .

Вектор  $\mathbf{g} = \mathbf{h} \otimes \mathbf{f}$  называется сверткой векторов, если его компоненты вычисляются по формуле

$$g_k = \sum_{j=1}^k h_{k-j+1} f_j, \quad k = 1, 2, \dots, (m+n). \quad (56)$$

Выражения (56) можно записать в матрично-векторном виде

$$\mathbf{H}_{+0} \mathbf{f}_{+0} = \mathbf{g}, \quad (57)$$

где квадратная ленточная матрица  $\mathbf{H}_{+0}$ , размера  $(m+n) \times (m+n)$  и шириной ленты  $n$ , имеет вид:

$$\mathbf{H}_{+0} = \begin{pmatrix} h_1 & 0 & 0 & \cdot & \cdot & \cdot & \cdot & \cdot & 0 \\ h_2 & h_1 & 0 & \cdot & \cdot & \cdot & \cdot & \cdot & 0 \\ h_3 & h_2 & h_1 & 0 & \cdot & \cdot & \cdot & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ h_n & h_{n-1} & h_{n-2} & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & h_n & h_{n-1} & 0 & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & h_2 & h_1 & 0 \\ 0 & \cdot & \cdot & \cdot & \cdot & \cdot & h_3 & h_2 & h_1 \end{pmatrix}.$$

Таким образом, с помощью выражений (56) можно вычислить свертку дискретных функций, а решая систему уравнений (57) можно определить вектор неизвестных  $\mathbf{f}$ .

Из ленточной структуры матрицы  $\mathbf{H}_{+0}$  видно, что дополнение вектора  $\mathbf{f}$  нулевыми компонентами не изменяет значений компонент

вектора свертки  $\mathbf{g}$  и, кроме того,  $g_{m+n} \equiv 0$ . Поэтому уравнение (57) можно заменить эквивалентным уравнением

$$\mathbf{H}_1 \mathbf{f} = \mathbf{g},$$

где ленточная матрица  $\mathbf{H}_1$ , размера  $(m+n-1) \times m$  и шириной ленты  $n$ , имеет вид:

$$\mathbf{H}_1 = \begin{pmatrix} h_1 & 0 & \cdot & \cdot & \cdot & 0 \\ h_2 & h_1 & \cdot & \cdot & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ h_n & h_{n-1} & \cdot & \cdot & \cdot & 0 \\ 0 & h_n & \cdot & \cdot & \cdot & h_1 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & h_n & h_{n-1} \\ 0 & \cdot & \cdot & \cdot & 0 & h_n \end{pmatrix}.$$

Алгоритм построения матрицы  $\mathbf{H}_1$  простой: ее вектор-столбцы суть

$$\mathbf{h}_j = \mathbf{T}_j \mathbf{h}_{+0}, \quad j = 1, 2, \dots, m,$$

где  $\mathbf{T}_j$  – оператор циклического сдвига компонент вектора  $\mathbf{h}_{+0}$  вниз на  $(j-1)$  позиций.

Если  $h_1 \neq 0$ , то матрица  $\mathbf{H}_1$  имеет полный ранг  $r = m$ . Однако ленточные матрицы такой структуры принадлежат к третьему типу обусловленности, поэтому оптимальное решение предпочтительнее вычислять по формуле  $\mathbf{f}_{opt} = \mathbf{H}_1^+ \mathbf{g}$  нежели по формуле  $\mathbf{f}_{opt} = \mathbf{H}_L^{-1} \mathbf{g}$ . При больших объемах выборок  $m$  и  $n$  необходимо проводить анализ спектра сингулярных чисел и вычислять приближенное оптимальное решение  $\mathbf{f}_{opt,app}$ , применяя процедуры пороговой модификации сингулярных чисел.

## Часть 3

### §6. Задания специального вычислительного практикума

*Каждая решенная мною задача становилась образцом, который использовался впоследствии для решения других задач.*  
Рене Декарт. Рассуждения о методе

Задания специального вычислительного практикума предназначены для более глубокого усвоения теоретического материала, изложенного в предыдущих параграфах. Компьютерные вычисления рекомендуется выполнять в среде системы программирования научно-технических расчетов **MATLAB**. Отчеты о выполненных заданиях должны содержать постановки задач, комментарии к полученным результатам и их графические изображения. Следуя рекомендациям Декарта, необходимо проявлять инициативу самостоятельного любопытства и творческого поиска.

Ниже приводится краткий список тестовых вопросов. Если учащийся в состоянии дать убедительные ответы, то ему нет необходимости выполнять задачи этого практикума. В противном случае мы настоятельно рекомендуем выполнить этот интересный, и, безусловно, полезный вычислительный практикум. Заметим, что даже опытные специалисты по численному моделированию затрудняются дать вразумительные ответы на эти вопросы.

#### Тестовые вопросы

##### 1. Легко ли открыть “псевдоэффект” ?

Решаются системы линейных алгебраических уравнений (СЛАУ)  $A\vec{x} = \vec{b}_i$ , содержащие погрешности в компонентах векторов правой части. Пусть

$$A = \begin{pmatrix} 10 & 7 \\ 4 & 3 \end{pmatrix}, \det(A) = 2.$$

Рассмотрим решения СЛАУ для различных векторов  $\vec{b}_i$ .

Пусть  $\vec{b}_0 = \begin{pmatrix} 7 \\ 3 \end{pmatrix}$  – точные значения коэффициентов в правой части. Тогда  $\vec{x}_0 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$  – точное решение системы.

Если коэффициенты содержат погрешности измерений, например  $\vec{b}_1 = \begin{pmatrix} 7.20 \\ 3.08 \end{pmatrix}$ , то возмущенное решение будет  $\vec{x}_1 = \begin{pmatrix} 0.02 \\ 1.00 \end{pmatrix}$ .

Пусть при более тщательном измерении коэффициентов имеем  $\vec{b}_2 = \begin{pmatrix} 7.00 \\ 3.06 \end{pmatrix}$ . Тогда  $\vec{x}_2 = \begin{pmatrix} -0.21 \\ +1.30 \end{pmatrix}$ .

Этот результат может служить основой интерпретации **псевдоэффекта**.

Очевидно, что  $\|\vec{b}_0 - \vec{b}_2\| \ll \|\vec{b}_0 - \vec{b}_1\|$ , однако

$$\|\vec{x}_0 - \vec{x}_2\| \gg \|\vec{x}_0 - \vec{x}_1\| !!!$$

Как объяснить этот эффект: *меньшие* погрешности в исходных данных вызвали *большее* возмущение в решении, чем *большие* погрешности?

Для каких матриц будут иметь место аналогичные эффекты ?

Указать класс матриц таких, что если

$$\|\vec{b}_0 - \vec{b}_2\| < \|\vec{b}_0 - \vec{b}_1\|, \text{ то и } \|\vec{x}_0 - \vec{x}_2\| < \|\vec{x}_0 - \vec{x}_1\|.$$

## 2. Какие решения получают, решая *несовместные* системы уравнений ?

Почему для всякой *несовместной* СЛАУ  $A\vec{x} = \vec{b}$  система уравнений  $A^T A\vec{y} = A^T \vec{b}$  всегда *совместна* и имеет, возможно, бесконечно много решений ?

**Пример:** Система уравнений  $\begin{pmatrix} 1 & 3 \\ 2 & 6 \end{pmatrix} \vec{x} = \begin{pmatrix} 5 \\ 8 \end{pmatrix}$  – несовместна.

Однако система уравнений  $\begin{pmatrix} 1 & 2 \\ 3 & 6 \end{pmatrix} \begin{pmatrix} 1 & 3 \\ 2 & 6 \end{pmatrix} \vec{y} = \begin{pmatrix} 1 & 2 \\ 3 & 6 \end{pmatrix} \begin{pmatrix} 5 \\ 8 \end{pmatrix}$  имеет

бесконечно много решений:  $\vec{y} = 21/50 \begin{pmatrix} 1 \\ 3 \end{pmatrix} + \begin{pmatrix} -3t \\ t \end{pmatrix}, t \in R.$

### 3. “Полезные ископаемые” на различных сериях экспериментальных данных.

Для описания функциональной зависимости выбрана некоторая линейная математическая модель  $f(\vec{p}, t) = \sum_{i=1}^{n_p} p_i \varphi_i(t)$ . Оптимальные параметры  $p_i$  для выбранной модели определяются путем решения задачи МНК. Предположим, что имеются две серии экспериментальных значений  $f_j, j = 1, 2, \dots, n, n > n_p$ , измеренных в одних и тех же точках  $t_j$ , однако с существенно различными уровнями погрешностей измерений.

**Могут ли совпадать значения оптимальных параметров, вычисленные по этим сериям измерений ?**

**Могут ли две серии измерений с одинаковым уровнем погрешностей, дать существенно различные оценки оптимальных параметров ?**

### 4. Волшебное превращение Золушки в прекрасную Принцессу.

Пусть квадратная матрица  $A$  такая, что вектор-столбцы этой матрицы являются попарно ортогональными, т.е.  $A^T A = D$ ,  $D$  – диагональная матрица. Однако вектор-строки этой матрицы могут и не быть попарно ортогональными, то есть матрица  $AA^T$  не является диагональной. Почему матрица  $H = AD^+$ , где  $D^+$  – диагональная матрица с диагональными элементами  $d_{ii}^+ = 1/\sqrt{d_{ii}}$  является унитарной, то есть

$$H^T H = H H^T = E, E \text{ – единичная матрица, } H^{-1} = H^T.$$



Этот факт – своеобразный “бриллиант” линейной алгебры !

**Каковы его полезные применения?**

**Тема 1. Топология отображений конечномерных пространств матрицами различных типов обусловленности**

1. Сформировать матрицы  $\mathbf{A}_i$ ,  $i = 1, 2, 3, 4$ , размера  $2 \times 2$ , различающиеся типом обусловленности. Какого типа будет матрица

$$\mathbf{A} = \sum_{i=1}^4 \mathbf{A}_i ?$$

С помощью этих матриц выполнить отображение  $\mathbf{R}^2 = \mathbf{X} \Rightarrow \mathbf{R}^2 = \mathbf{Y}$  следующих множеств точек  $M(\mathbf{x})$ :

а)  $M_C(\mathbf{x}) = \left\{ \mathbf{x}(t) : \mathbf{x}(t) = \begin{pmatrix} \cos(t) \\ \sin(t) \end{pmatrix}, 0 \leq t \leq 1 \right\}$  – Miss Circle.

Отображение первой четверти окружности произвести с помощью матрицы первого типа, второй четверти – с помощью матрицы второго типа и т.д.;

б)  $M_S(\mathbf{x}) = \left\{ \mathbf{x}(t) : \mathbf{x}(t) = \begin{pmatrix} t \\ \sin(2\pi t) \end{pmatrix}, -a \leq t \leq a \right\}$  – Sine-горки;

в)  $M_{As}(\mathbf{x}) = \left\{ \mathbf{x}(t) : \mathbf{x}(t) = \begin{pmatrix} \sin(2\pi t) \\ t \end{pmatrix}, -a \leq t \leq a \right\}$  – змея Arcsin;

г)  $M_F(\mathbf{x}) = \left\{ \mathbf{x}(t) : \mathbf{x}(t) = \begin{pmatrix} t \\ 1/|t| \end{pmatrix}, -a \leq t \leq a \right\}$  – контуры Фудзи-ямы;

д)  $M_{Chebur}(\mathbf{x})$  – контуры графического изображения Чебурашки. Будут ли на множествах  $M_{Chebur,i}(\mathbf{y} = \mathbf{A}_i \mathbf{x})$  присутствовать изображения веселых улыбок или гримасы разочарования?

е) Существуют ли матрицы, которые производят инверсию графических изображений слов:

$$\begin{aligned} \text{ТОМ} &\Leftrightarrow \text{МОТ} \\ \text{КОТ} &\Leftrightarrow \text{ТОК}. \end{aligned}$$

ж)  $M_{VIP}(\mathbf{x})$  – задается в соответствии с эстетическими вкусами и уровнем фантазии исполнителя.

Указания:

а) Матрицы  $\mathbf{A}_j$  должны быть общего вида (т.е. не диагональные);

б) Значения параметра  $a$  выбираются так, чтобы отчетливые графические изображения каждой пары множеств  $M_j(\mathbf{x})$  и  $M_j(\mathbf{y})$  умещались на одной странице формата А4;

в) Формирование матриц различного типа обусловленности проще всего выполнить следующим образом. Очевидно, что матрица

$$\mathbf{A}_2 = \begin{pmatrix} \cos \varphi & -\sin \varphi \\ \sin \varphi & \cos \varphi \end{pmatrix}$$

является ортогональной, т.е. это матрица второго типа. Далее определим матрицу

$$\mathbf{A}(t_1, t_2) = \begin{pmatrix} t_1 \cos \varphi & -t_2 \sin \varphi \\ t_1 \sin \varphi & t_2 \cos \varphi \end{pmatrix}.$$

Тогда, задавая значения  $t_1 > 1$  и  $t_2 > 1$ ,  $t_1 \neq t_2$ , получаем матрицу первого типа. Задавая  $t_1 > 1$  и  $0 < t_2 < 1$  имеем матрицу третьего типа, а для  $0 < t_1 < 1$  и  $0 < t_2 < 1$ ,  $t_1 \neq t_2$  получаем матрицу четвертого типа.

г) Формирование множества  $M_{VIP}(\mathbf{x})$  является обязательным для каждого исполнителя!

2. Сформировать матрицу третьего типа  $\mathbf{A}_3$  размера  $2 \times 2$ , число обусловленности которой  $cond(\mathbf{A}_3) \geq 100$ . Задать вектор  $\mathbf{x}_1$  и вычислить  $\mathbf{y}_1 = \mathbf{A}_3 \mathbf{x}_1$ . Задать вектор погрешности  $\delta \mathbf{y}$ . Определить при каких ориентациях вектора  $\delta \mathbf{y}$  (фиксированной длины  $\|\delta \mathbf{y}\| \approx 0.03 * \|\mathbf{y}_1\|$ ) длины векторов  $\mathbf{x} + \delta \mathbf{x} = \mathbf{A}_3^{-1}(\mathbf{y}_1 + \delta \mathbf{y})$  будут достигать своего наименьшего и наибольшего значений?

Примечание (шутка). Сочинитель уникальной матрицы третьего типа, для которой длина возмущенного решения (при малых значениях погрешностей в правой части) может достигают катастрофи-

чески большого значения, будет рекомендован для регистрации в Книге рекордов России.

3. С помощью некоторой матрицы  $\mathbf{A}$ , размера  $5 \times 2$  и ранга 2, выполнить отображение  $\mathbf{R}^2 \Rightarrow \mathbf{R}^4$  множества точек  $\Omega_{\mathbf{x}} = \{\mathbf{x} : \|\mathbf{x}\| = 1\}$ . Чтобы выяснить, что представляет собой множество  $\Omega_{\mathbf{y}} = \{\mathbf{y} : \mathbf{y} = \mathbf{Ax}, \|\mathbf{x}\| = 1\}$  можно воспользоваться одним из следующих способов:

а) Построить специальный летательный аппарат, произвести полет над ранговым пространством  $\mathbf{R}(\mathbf{A})$  и визуально изучить множество  $\Omega_{\mathbf{y}}$ ;

б) Выполнить экспликацию множества точек  $\Omega_{\mathbf{y}}$ . Результат экспликации изобразить графически в полярной системе координат. Визуально убедиться в том, что полученная фигура представляет собой эллипс. Вычислить отношение длин осей этого эллипса. Сравнить его площадь с площадью единичного круга;

в) Произвести нормировку векторов из множества  $\Omega_{\mathbf{y}}$  и сформировать  $\Omega_{\mathbf{e}} = \{\mathbf{e} = \mathbf{Ax} / \|\mathbf{Ax}\|\}$ . Изобразить графически множество векторов  $\mathbf{x} = \mathbf{A}_L^{-1}\mathbf{e}$ ,  $\mathbf{e} \in \Omega_{\mathbf{e}}$  и объяснить полученный результат.

Комментарий: Ясно, что способ а) относится к области научной фантастики, способ б) – достаточно трудоемкий, а способ в) является интеллектуально наиболее элегантным. В прочем, не возбраняется воспользоваться всеми тремя способами.

## Тема 2. Сингулярное разложение матриц и его различные применения

1. Найти размерности и ортогональные базисы четырех основных подпространств  $\mathbf{R}(\mathbf{A})$ ,  $\mathbf{N}(\mathbf{A})$ ,  $\mathbf{R}(\mathbf{A}^T)$  и  $\mathbf{N}(\mathbf{A}^T)$  матрицы  $\mathbf{A}$ :

$$\mathbf{A} = \begin{pmatrix} 1 & 0 & 1 & 2 \\ 2 & 1 & 2 & 4 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 2 \\ 2 & 3 & 4 & 6 \\ 6 & 7 & 8 & 14 \end{pmatrix}.$$

Определить тип матрицы и число обусловленности. Разложить вектор  $\mathbf{y}_1 = (1, 2, 3, 4, 5, 6)^T$  на ортогональные составляющие  $\mathbf{y}_1 = \mathbf{y}_r + \mathbf{y}_0$ ,  $\mathbf{y}_r \in \mathbf{R}(\mathbf{A})$ ,  $\mathbf{y}_0 \in \mathbf{N}(\mathbf{A}^T)$ . Используя сингулярное разложение матрицы найти оптимальное решение и оптимальную невязку для системы  $\mathbf{A}\mathbf{x} = \mathbf{y}_1$ . Найти вектор  $\mathbf{y}_2$  такой, что  $\mathbf{x}_{opt} = \mathbf{A}^+ \mathbf{y}_1 = \mathbf{A}^+ \mathbf{y}_2$ . Найти наиболее и наименее благоприятные ориентации вектора погрешности фиксированной длины.

2. Для ленточной матрицы  $\mathbf{A}_{m,n}$  размера  $(m+n-1) \times n$ ,  $m=5$ ,  $n=4$ , вычислить спектр сингулярных чисел, определить тип и число обусловленности.

$$\mathbf{A}_{5,4} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

Как будут изменяться числа обусловленности матриц такой структуры при изменении их размеров?

3. Элементами матрицы  $\mathbf{A}$ , размера  $n \times n$ ,  $n=1000$ , являются случайные числа, равномерно распределенные на интервале  $[-1,1]$ . Выполнить сингулярное разложение матрицы и построить график спектра сингулярных чисел. Будет ли матрица  $\mathbf{A}$  близка к вырожденной? Как будут изменяться спектры сингулярных чисел при

изменении интервалов равномерных распределений. Произвести аналогичные вычисления для матриц, элементами которой являются нормально распределенные случайные числа с некоторым стандартным отклонением  $\delta \neq 1$ .

### Тема 3. Сжатие информации с применением метода сингулярного разложения матриц

1. На планете ZG5-V10 имеется горный массив, состоящий из 10 потухших вулканов. Линейные размеры массива –  $(20 \times 20)$  км. Рельеф поверхности этого массива описывается суперпозицией функций Гаусса

$$h(u, v) = \sum_{k=1}^{10} w_k \exp(-((u - u_k)^2 + (v - v_k)^2) / d_k^2),$$

где:  $w_k$  – максимальная высота  $k$ -го вулкана,  $(u_k, v_k)$  – координаты центров вулканов,  $d_k$  – радиусы вулканов. Фотографический снимок этого горного массива в оцифрованном виде задан в виде матрицы  $\mathbf{H}$  размера  $n \times n$ ,  $n = 1000$ . Объем информации для хранения или передачи такого оцифрованного снимка  $I(h(u_i, v_j)) = 10^6$  чисел. Цель задания – выяснить возможности сжатия объема информации без потери существенных деталей изображения горного массива.

#### Указания:

а) Задать набор чисел  $\{w_k, u_k, v_k, d_k, k = 1, 2, \dots, 10\}$ . Произвести дискретную выборку значений функции  $h(u_i, v_j)$ ,  $i, j = 1, 2, \dots, n$ ,  $n = 1000$ , и сформировать матрицу  $\mathbf{H}$ , элементы которой  $h_{ij} = h(u_i, v_j)$ ;

б) Выполнить сингулярное разложение матрицы  $\mathbf{H} = \mathbf{USV}^T$  и построить график спектра сингулярных чисел  $\sigma_i, i = 1, 2, \dots, n$ ;

в) Произвести обнуление сингулярных чисел  $\sigma_i = 0, i = k + 1, \dots, n$ , и вычислить матрицу  $\tilde{\mathbf{H}} = \tilde{\mathbf{U}}\tilde{\mathbf{S}}\tilde{\mathbf{V}}^T$ . Матрица  $\tilde{\mathbf{U}}$ , размера  $n \times k$ , такая, что  $k$  вектор-столбцов матрицы  $\tilde{\mathbf{U}}$  совпадают с первыми  $k$  вектор-столбцами матрицы  $\mathbf{U}$ .  $\tilde{\mathbf{S}}$  – матрица усеченных сингулярных чисел размера  $n \times k$ . Аналогично формируется матрица  $\tilde{\mathbf{V}}$ . На основе визуального анализа графического изображения матриц  $\tilde{\mathbf{H}}$  выяснить, при каких минимальных значениях  $k$  изображение горного массива

можно считать удовлетворительным (т.е. достаточно точно локализируются координаты всех вулканов, отсутствуют артефакты в виде ложных вулканов и т.д.). Сравнить нормы Фробениуса матриц  $\mathbf{H}$  и  $\tilde{\mathbf{H}}$  (при удовлетворительном сжатом изображении горного массива);

г) Вычислить коэффициент максимально допустимого сжатия информации

$$\alpha = \frac{I(\tilde{\mathbf{H}})}{I(\mathbf{H})} = \frac{2nk + k}{n^2},$$

где:  $I(\mathbf{H})$  – объем информации для хранения матрицы  $\mathbf{H}$ ,  $I(\tilde{\mathbf{H}})$  – объем информации для хранения матрицы  $\tilde{\mathbf{H}}$ .

д) Как изменится контрастность снимка, если

$$\hat{\mathbf{H}} = \mathbf{H} + \mathbf{W},$$

где  $\mathbf{W}$  – матрица “белого” шума, т.е. значениями элементов  $w_{ij}$  являются случайные числа с равномерным законом распределения?

#### Тема 4. Оптимальные параметры линейных математических моделей и оценка их доверительных радиусов

1) Убедиться, что две серии экспериментальных данных, приведенных в таб. 1, являются эквивалентными в смысле МНК для линейной математической модели

$$f(t) = p_1 + p_2 * t + p_3 * t^2.$$

Указания:

а) Сформировать матрицу плана  $\mathbf{F}$ , размера  $21 \times 3$ , следующего вида

$$\mathbf{F} = \begin{pmatrix} 1 & t_1 & t_1^2 \\ 1 & t_2 & t_2^2 \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ 1 & t_n & t_n^2 \end{pmatrix}.$$

Компоненты векторов  $\mathbf{t}$ ,  $\mathbf{f}_{e,1}$  и  $\mathbf{f}_{e,2}$  взять из таб. 1;

б) Вычислить и сравнить векторы оптимальных параметров для двух серий экспериментальных данных  $\mathbf{f}_{e,1}$  и  $\mathbf{f}_{e,2}$  по формуле

$$\mathbf{p}_{opt,i} = (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{f}_{e,i};$$

в) Оценить длины векторов погрешностей в двух сериях измерений и вычислить радиусы доверительной области векторов оптимальных параметров по формуле

$$\rho_i = \|\mathbf{F}^+\| \sqrt{\|\boldsymbol{\delta}_i\|^2 - \|\boldsymbol{\delta}_{0,i}\|^2},$$

где:  $\|\mathbf{F}^+\| = 1/\sigma_{\min}(\mathbf{F})$ ,  $\|\tilde{\boldsymbol{\delta}}_{0,i}\|^2 = \|\mathbf{F}\mathbf{p}_{opt,i} - \mathbf{f}_{e,i}\|^2$  – остаточные значения функционала МНК в оптимальных точках  $\mathbf{p}_{opt,i}$ .

2) Для экспериментальных данных, приведенных в таб. 1, произвести  $n$ -кратное скользящее усреднение данных ( $n=1, 2, 3$ ), вычислить оптимальные параметры модели после каждого усреднения и сравнить результаты. Скользящее усреднение проводить по двум соседним точкам.

3) Вычислить две серии эквивалентных “экспериментальных” данных для линейной математической модели

$$f(t) = p_1 + p_2 * t + p_3 * t^2 + p_4 * t^3 + p_5 * t^4.$$

Указания:

а) Задать векторы  $\mathbf{p} = (p_1, p_2, \dots, p_5)^T$  и  $\mathbf{t} = (t_1, t_2, \dots, t_n)^T$ ;

б) Сформировать матрицу плана  $\mathbf{F}$ , вычислить вектор значений  $\mathbf{f}$ ;

в) Сформировать два вектора погрешностей  $\delta_1 \mathbf{f}$  и  $\delta_2 \mathbf{f}$  с помощью датчика случайных чисел, распределенных по нормальному закону с различными дисперсиями. Разложить эти векторы на ортогональные составляющие

$$\delta_i \mathbf{f} = \delta_{i,r} \mathbf{f} + \delta_{i,0} \mathbf{f}, \quad i = 1, 2,$$

где

$$\delta_{i,r} \mathbf{f} = \mathbf{F}(\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \delta_i \mathbf{f},$$

$$\delta_{i,0} \mathbf{f} = (\mathbf{E}_n - \mathbf{F}(\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T) \delta_i \mathbf{f};$$

г) Сформировать два вектора:  $\tilde{\mathbf{f}}_1 = (\mathbf{f} + \delta_{1,r} \mathbf{f}) + \delta_{1,0} \mathbf{f}$ ,

$$\tilde{\mathbf{f}}_2 = (\mathbf{f} + \delta_{2,r} \mathbf{f}) + \delta_{2,0} \mathbf{f}.$$

Убедиться, что эти векторы являются эквивалентными в смысле определения оптимальных параметров по МНК. Можно ли на основе визуального анализа графиков отдать предпочтение одной из серий “экспериментальных” измерений –  $\{t_j, \tilde{\mathbf{f}}_1(t_j)\}$  или  $\{t_j, \tilde{\mathbf{f}}_2(t_j)\}$ ?

## Тема 5. Аппроксимация функций

1. С помощью многочленов степени  $k$ ,  $k = 5, 6$ , аппроксимировать функции:

а)  $f_1(x) = \exp(-x)$ ,  $0 \leq x \leq 2$ ;

б)  $f_2(x) = \exp(-x^2)$ ,  $0 \leq x \leq 2$ ;

в)  $f_3(x) = \frac{1}{1+x^2}$ ,  $0 \leq x \leq 2$ ;

г)  $f_4(x)$ ,  $f_5(x)$ , ... – выбрать самостоятельно.

2. Аппроксимировать выше упомянутые функции тригонометрическими многочленами  $T_k(x)$ ,  $k = 5, 6$  и сравнить результаты.

## Тема 6. Численное решение интегрального уравнения свертки

1. Для заданных функций  $f(t)$  и  $h(t)$ ,  $-a \leq t \leq a$ ,  $a = 1$ , вычислить свертку

$$g(t) = \int_{\max(t, -a)}^{\min(t, a)} f(t - \tau)h(\tau)d\tau.$$

Произвести дискретную выборку с некоторым постоянным шагом значений функций  $f(t)$  и  $h(t)$  (объем выборки  $n > 100$ ) и образовать векторы  $\mathbf{f}$  и  $\mathbf{h}$ . Сформировать матрицу дискретной свертки  $\mathbf{N}_1$  и вычислить свертку  $\mathbf{g} = \mathbf{N}_1\mathbf{f}$ . Вычислить  $\mathbf{f}_{sol} = \mathbf{N}_1^+\mathbf{g}$  и сравнить с вектором  $\mathbf{f}$  и функцией  $f(t)$ . Провести анализ спектра сингулярных чисел матрицы  $\mathbf{N}_1$  и применить процедуру пороговой модификации малых сингулярных чисел с целью получения приближенных решений  $\tilde{\mathbf{f}}_{sol} = \tilde{\mathbf{N}}_1^+\mathbf{g}$ .

Примечание. Функции  $f(t)$  и  $h(t)$  каждый исполнитель выбирает самостоятельно.



## Тема 7. Приближенное решение плохо-обусловленных систем линейных уравнений большого порядка

### Указания:

1. Для матрицы  $\mathbf{H}$ , рассматриваемой в Теме 3, провести анализ спектра сингулярных чисел.

2. Сформировать вектор  $\mathbf{v}$  с компонентами  $v_j = h_{ij}$ ,  $j = 1, 2, \dots, n$ .

Номер строки  $i$  должен соответствовать наиболее выразительному одномерному сечению функции  $h(u, v)$ . Вычислить:

$$\mathbf{u}_{1,r} = \mathbf{H}\mathbf{v} \in \mathbf{R}(\mathbf{H}), \quad \mathbf{v}_r = \mathbf{H}^T \mathbf{u}_{1,r} \in \mathbf{R}(\mathbf{H}^T) \quad \text{и} \quad \mathbf{u}_{2,r} = \mathbf{H}\mathbf{v}_r \in \mathbf{R}(\mathbf{H}).$$

Эти манипуляции необходимы для того, чтобы  $\mathbf{v}_r$  и  $\mathbf{u}_{2,r}$  принадлежали ранговым пространствам  $\mathbf{R}(\mathbf{H}^T)$  и  $\mathbf{R}(\mathbf{H})$ , соответственно. Тогда решением системы  $\mathbf{H}\mathbf{v} = \mathbf{u}_{2,r}$  будет вектор  $\mathbf{v} = \mathbf{v}_r = \mathbf{H}^+ \mathbf{u}_{2,r}$ . Однако, если в спектре сингулярных чисел наблюдается “нырок” к очень малым числам, то необходимо проводить пороговую модификацию сингулярных чисел и вычислять приближенные решения.

4. Проводя серию пороговой модификации спектра сингулярных чисел матрицы  $\mathbf{H}$ , вычислить

$$\mathbf{v}_{opt,app} = \tilde{\mathbf{H}}^+ \mathbf{u}_{2,r}, \quad \tilde{\mathbf{H}} = \mathbf{U} \tilde{\mathbf{S}} \mathbf{V}^T.$$

5. Сравнить графики дискретных функций  $\{j, v_{r,j}\}$  и  $\{j, v_{opt,app,j}\}$ , где  $j$  – номер компоненты векторов  $\mathbf{v}_r$  и  $\mathbf{v}_{opt,app}$ .

## Рекомендуемая литература

1. Форсайт Дж., Молер К. Численное решение систем линейных алгебраических уравнений. – М.: Мир, 1969.
2. Лоусон Ч., Хенсон Р. Численное решение задач метода наименьших квадратов. – М.: Наука, 1986.
3. Голуб Дж., Ван Лоун Ч. Матричные вычисления. – М.: Мир, 1999.
4. Деммель Дж. Вычислительная линейная алгебра. Теория и приложения. – М.: Мир, 2001.

## Содержание

<b>Введение</b> .....	3
<b>Часть 1</b> ..	8
<b>§1.</b> Основная теорема линейной алгебры .....	8
<b>§2.</b> Сингулярное разложение матриц .....	23
<b>§3.</b> Спектральная классификация матриц произвольного размера и критерии плохой обусловленности .....	34
<b>Часть 2</b> .....	46
<b>§4.</b> Линейные математические модели: оптимальные параметры и радиусы доверительной области .....	46
<b>§5.</b> Алгоритмы численного решения интегрального уравнения типа свертки .....	56
<b>Часть 3</b> .....	62
<b>§6.</b> Задания специального вычислительного практикума .....	62
Рекомендуемая литература .....	74

Учебное издание

АНДРУШЕВСКИЙ Николай Матвеевич

Анализ устойчивости решений  
систем линейных алгебраических уравнений

Методическое пособие  
специального вычислительного практикума

Издательский отдел  
Факультета вычислительной математики и кибернетики  
МГУ имени М.В. Ломоносова  
Лицензия ИД N 05899 о 24.09.01 г.

119992, ГСП-2, Москва, Ленинские горы, МГУ им. М.В. Ломоносова,  
2-й учебные корпус

Напечатано с готового оригинал-макета  
в издательстве ООО <<МАКС Пресс>>  
Лицензия ИД N 00510 от 01.12.99 г.  
Подписано к печати 12.11.2008  
Формат 60x90 1/16. Усл.печ.л. 4,38. Тираж 100 экз. Заказ 659.

119992, ГСП-2, Москва, Ленинские горы, МГУ им. М.В. Ломоносова,  
2-й учебный корпус, 627 к.  
Тел. 939- 3890. Тел./Факс 939-3891