Programming, numerics and optimization Lecture B-4: Linear systems II

Łukasz Jankowski ljank@ippt.pan.pl

Institute of Fundamental Technological Research Room 4.32, Phone +22.8261281 ext. 428

April 20, 2021¹

¹Current version is available at http://info.ippt.pan.pl/~ljank.

Outline

- Least-squares problem
- 2 Conditioning
- 3 Numerical regularization
- 4 Large Toeplitz systems
- 5 Further reading

6 Homework 8

Outline



- Normal equations
- Gauss-Markov theorem
- SVD-based solution
- Iterative methods

Least-squares problem

Least-squares problems originate from estimation of model parameter in problems of fitting noisy data.

Consider the following model

$$y(t)=f(t;\mathbf{x}),$$

where $\mathbf{x} \in \mathbb{R}^n$ is a vector of unknown model parameters, to be determined from a series of (noisy) measurements (y_i, t_i) , $i = 1, 2, \ldots, m$. To reduce the noise influence, m > n.

The related least-squares problem is in general a nonlinear optimization problem:

find **x**, which minimizes
$$\sum_{i=1}^{m} (y_i - f(t_i; \mathbf{x}))^2$$
.

Linear least-squares problem

If the dependence on parameters \mathbf{x} is linear,

$$y(t) = \sum_{j=1}^{n} x_j f_j(t),$$

a linear least-squares problem is obtained,

find **x**, which minimizes
$$\sum_{i=1}^{m} \left[y_i - \sum_{j=1}^{n} x_j f_j(t_i) \right]^2$$
,

which is equivalent to the following convex optimization problem:

find **x**, which minimizes $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2$, where $\mathbf{x} = [x_1, \dots, x_n]^T$, $\mathbf{y} = [y_1, \dots, y_m]^T$ and $\mathbf{A} \in \mathbb{R}^{m \times n}$ is an $m \times n$ matrix composed of $f_j(t_i)$.

Linear least-squares problem

Linear least-squares problem

```
find x, which minimizes \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2,
```

where $\mathbf{x} \in \mathbb{R}^{n}$, $\mathbf{y} \in \mathbb{R}^{m}$ and $\mathbf{A} \in \mathbb{R}^{m \times n}$.

The related linear equation $\mathbf{y} = \mathbf{A}\mathbf{x}$ can be in general inconsistent (m > n) and have no exact solutions. However, the least-squares solution \mathbf{x}_{LS} is the best solution in the sense of the *minimum* residual norm:

- Ax_{LS} is the closest point to y out of all points from range A.
- Ax_{LS} is the orthogonal projection of y onto range A.

The least-squares solution is unique, if rank $\mathbf{A} = n$ (full column rank), that is if \mathbf{A} corresponds to an injection.

Normal equations

Normal equations

The vector **x** minimizes $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2$ if and only if it satisfies the *normal equations*,

$$\mathbf{A}^{\mathsf{T}}\mathbf{A}\mathbf{x} = \mathbf{A}^{\mathsf{T}}\mathbf{y},$$

or equivalently $\mathbf{A}^{\mathsf{T}}(\mathbf{y} - \mathbf{A}\mathbf{x}) = \mathbf{0}$.

- If m >> n, the matrix $\mathbf{A}^{\mathsf{T}} \mathbf{A} \in \mathbb{R}^{n \times n}$ is much smaller than the initial matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$.
- In the general case, **A**^T**A** is (square) symmetric positive semidefinite.
- If **A** is full column rank, then **A**^T**A** is full-rank positive definite and normal equations have a unique solution.
- A disadvantage of the normal equations is squaring of the condition number (if A is ill-conditioned, then conditioning A^TA is even worse).

Gauss-Markov theorem

The Gauss-Markov theorem

- recasts the problem of estimating a vector of parameters in a (noisy) linear model using the statistical terminology, and
- states that the least-squares solution is the optimum estimator of the parameter vector in the statistical sense, that it is
 - unbiased (mean estimation error is zero) and
 - of minimum variance.

Gauss-Markov theorem

Conditioning

Least-squares problem

000000000000

Outline

Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ be a full column rank matrix of real numbers, $\mathbf{y} \in \mathbb{R}^m$ be a vector of (noisy) measurements (or observations) and $\mathbf{x} \in \mathbb{R}^n$ be an unknown parameter vector. Assume the linear model

Regularization

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \boldsymbol{\epsilon},$$
 (*

Large Toeplitz systems

Reading

where ϵ is a vector of uncorrelated random measurement errors with zero mean and the same variance,

$$\mathsf{E}[\boldsymbol{\epsilon}] = \mathbf{0}, \qquad \qquad \mathsf{V}[\boldsymbol{\epsilon}] = \sigma^2 \mathbf{I}.$$

The best unbiased estimator $\hat{\mathbf{x}}$ of \mathbf{x} is the solution to the least-squares problem associated with (*):

$$\hat{\mathbf{x}} = (\mathbf{A}^{\mathsf{T}}\mathbf{A})^{-1}\mathbf{A}^{\mathsf{T}}\mathbf{y}, \qquad \qquad \mathbf{V}[\hat{\mathbf{x}}] = \sigma^{2}(\mathbf{A}^{\mathsf{T}}\mathbf{A})^{-1}.$$

Singular value decomposition (SVD)

Every rectangular matrix $\mathbf{A} \in \mathbb{R}^{m imes n}$ can be expressed as

 $\mathbf{A} = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^\mathsf{T},$

where **U** and **V** are unitary matrices, and Σ is an $m \times n$ diagonal matrix with rank **A** nonnegative diagonal elements σ_i (called singular values) ordered in a nonincreasing way.

Given matrix \mathbf{A} , its rank can be marked explicitly in the SVD:

$$\mathbf{A} = \begin{bmatrix} \mathbf{U}_1 \mathbf{U}_2 \end{bmatrix} \begin{bmatrix} \mathbf{\Sigma}_0 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{V}_1^{\mathsf{T}} \\ \mathbf{V}_2^{\mathsf{T}} \end{bmatrix} = \mathbf{U}_1 \mathbf{\Sigma}_0 \mathbf{V}_1^{\mathsf{T}},$$

where the number of columns of U_1 , V_1 and Σ_0 is rank A. The reduced decomposition $A = U_1 \Sigma_0 V_1^T$ is called the thin SVD of A.

Conditioning

Let $\mathbf{A} = \mathbf{U}_1 \mathbf{\Sigma}_0 \mathbf{V}_1^T$ be the thin SVD of $\mathbf{A} \in \mathbb{R}^{m \times n}$. The least-squares problem

find **x**, which minimizes $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2$,

Regularization

Large Toeplitz systems

has always a unique minimum norm solution^a given by

$$\mathbf{x}_{\mathsf{svd}} = \mathbf{V}_1 \mathbf{\Sigma}_0^{-1} \mathbf{U}_1^\mathsf{T} \mathbf{y}.$$

^aThat is, if solution to the least-squares problem is nonunique (rank A < n), the solution with the minimum norm is chosen.

The matrix

Outline

Least-squares problem

$$\mathbf{A}^+ = \mathbf{V}_1 \mathbf{\Sigma}_0^{-1} \mathbf{U}_1^{\mathsf{T}}$$

is called the pseudo-inverse or the Moore-Penrose inverse of A.

SVD-based least-squares solution

Conditioning

Outline

Least-squares problem

• The SVD expresses the matrix **A** as a sum of rank-one matrices:

$$\mathbf{A} = \mathbf{U}_1 \mathbf{\Sigma}_0 \mathbf{V}_1^\mathsf{T} = \sum_{i=1}^{\mathsf{rank}\,\mathbf{A}} \sigma_i \, \mathbf{u}_i \mathbf{v}_i^\mathsf{T},$$

Regularization

Large Toeplitz systems

where \mathbf{u}_i and \mathbf{v}_i are the *i*th columns of **U** and **V**, respectively.

• The minimum norm minimizer of $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2$ can be thus expressed as

$$\mathbf{x}_{\mathsf{svd}} = \mathbf{V}_1 \mathbf{\Sigma}_0^{-1} \mathbf{U}_1^{\mathsf{T}} \mathbf{y} = \left[\sum_{i=1}^{\mathsf{rank}\,\mathbf{A}} \frac{\mathbf{v}_i \mathbf{u}_i^{\mathsf{T}}}{\sigma_i}\right] \mathbf{y}_i$$

SVD-based approach vs. normal equations

In comparison to the normal equations

- + In the SVD-based approach, the initial matrix **A** is used directly. Therefore, the condition number is not squared.
- If m >> n, the matrix **A** is much larger than $\mathbf{A}^{\mathsf{T}}\mathbf{A}$ and the numerical cost of computing the SVD can be prohibitive.

Iterative methods

The linear least-squares problem, expressed as

```
find x, which minimizes \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2,
```

is in fact the problem of unconstrained minimization of a convex quadratic objective function

$$\phi(\mathbf{x}) = \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 = \mathbf{x}^\mathsf{T}\mathbf{A}^\mathsf{T}\mathbf{A}\mathbf{x} - 2\,\mathbf{y}^\mathsf{T}\mathbf{A}\mathbf{x} + \mathbf{y}^\mathsf{T}\mathbf{y},$$

where $\mathbf{A}^{\mathsf{T}}\mathbf{A}$ is positive semidefinite.

Therefore, any iterative optimization scheme can be used to find an approximation to the exact solution, in particular the conjugate gradient least-squares (CGLS) method described in Lecture B-3.

Outline



- Vector and matrix norms
- Conditioning and estimation of accuracy
- Deconvolution as a common source of ill-conditioning

Norm on a linear space

A norm $\|\cdot\|$ on a linear space \mathcal{X}

Let $\mathcal X$ be a real linear space. A function $\|\cdot\|:\mathcal X\to\mathbb R$ with the properties

$$\begin{split} \|\mathbf{x}\| &\geq 0 & (\text{positivity}) \\ \|\mathbf{x}\| &= 0 \text{ iff } \mathbf{x} = \mathbf{0} & (\text{definitness}) \\ \|k\mathbf{x}\| &= |k| \|\mathbf{x}\| & (\text{homogeneity}) \\ \|\mathbf{x} + \mathbf{y}\| &\leq \|\mathbf{x}\| + \|\mathbf{y}\| & (\text{triangle inequality}) \end{split}$$

for all $\mathbf{x}, \mathbf{y} \in \mathcal{X}$ and all $k \in \mathbb{R}$ is called a norm on \mathcal{X} .

Vector norms

Vector *p*-norms

The most often used vector norms are the so-called *p*-norms $\|\cdot\|_p$:

$$\|\mathbf{x}\|_{\boldsymbol{p}} = \left[\sum_{i} |x_{i}^{\boldsymbol{p}}|\right]^{\frac{1}{p}}.$$

In particular,

$$\|\mathbf{x}\|_1 = \sum_i |x_i|,$$
$$\|\mathbf{x}\| = \|\mathbf{x}\|_2 = \sqrt{\sum_i x_i^2},$$
$$\|\mathbf{x}\|_{\infty} = \max_i |x_i|.$$

Matrix norms

Operator norms

A vector norm $\|\cdot\|$ induces the corresponding matrix operator norm

$$\|\mathbf{A}\| = \max_{\mathbf{x}\neq\mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|}.$$

In particular, the vector *p*-norms induce the matrix norms:

$$\begin{aligned} \|\mathbf{A}\|_{1} &= \max_{j} \sum_{i} |a_{ij}|, \\ \|\mathbf{A}\|_{2} &= \sqrt{\max_{i} \lambda_{i} (\mathbf{A}^{\mathsf{T}} \mathbf{A})} = \max_{i} \sigma_{i} (\mathbf{A}), \\ \|\mathbf{A}\|_{\infty} &= \max_{i} \sum_{j} |a_{ij}|. \end{aligned}$$

The Frobenius matrix norm is defined by

$$\|\mathbf{A}\|_{\mathsf{F}} = \sqrt{\sum_{i,j} a_{ij}^2} = \sqrt{\sum_i \lambda_i^2}.$$

Consistent matrix and vector norms

It is important to use the matrix norms that are *consistent* with the used vector norms, that is

 $\|\mathbf{A}\mathbf{x}\| \leq \|\mathbf{A}\|\|\mathbf{x}\|.$

The following norms are consistent:

$$\begin{split} \|\mathbf{x}\|_{1} \dots \|\mathbf{A}\|_{1} \\ \|\mathbf{x}\|_{2} \dots \|\mathbf{A}\|_{2}, \|\mathbf{A}\|_{F} \\ \|\mathbf{x}\|_{\infty} \dots \|\mathbf{A}\|_{\infty} \end{split}$$

Condition number

The notion of the condition number of a problem has been introduced in Lecture B-1:

- The condition number is the amplification factor of the relative error between input and output data. It measures how much the errors in the input data can affect the errors in the output data.
- Conditioning is a property of the problem, which (independent of any algorithm used to solve it) can be well-conditioned or ill-conditioned.
- If a problem is ill-conditioned, an algorithm can give better results only by chance. The errors in input data will inevitably propagate to the output data.

Condition number — diagonal matrices

Conditioning

Let $\mathbf{D} = \text{diag}(d_1, d_2, \dots, d_n)$ be a nonsingular diagonal matrix. The linear system $\mathbf{D}\mathbf{x} = \mathbf{y}$ is a system of simple decoupled equations:

$$d_i x_i = y_i, \qquad i = 1, 2, \ldots, n,$$

Regularization

Reading

Large Toeplitz systems

which has the exact solution

Outline

Least-squares problem

$$x_i = rac{y_i}{d_i}, \qquad i=1,2,\ldots,n.$$

In practice, the right-hand side is known only approximately (due to measurement errors, floating-point representation errors, etc.), so that the system equation takes the form

$$d_i(x_i + \Delta x_i) = y_i + \Delta y_i, \qquad i = 1, 2, \ldots, n,$$

where $\frac{|\Delta y_i|}{\|\mathbf{y}\|} \approx \epsilon > 0$. Notice that $d_i \Delta x_i = \Delta y_i$.

Condition number — diagonal matrices

Conditioning

Outline

Least-squares problem

Since $d_i \Delta x_i = \Delta y_i$, the absolute error of the solution is

$$\Delta x_i = \frac{\Delta y_i}{d_i}, \qquad i = 1, 2, \dots, n.$$

Regularization

Large Toeplitz systems

Reading

Hence, the relative error level of the solution is bounded as follows

$$\frac{|\Delta x_i|}{\|\mathbf{x}\|} = \frac{1}{d_i} \frac{|\Delta y_i|}{\|\mathbf{x}\|} = \frac{1}{d_i} \frac{\|\mathbf{y}\|}{\|\mathbf{x}\|} \frac{|\Delta y_i|}{\|\mathbf{y}\|} \le \frac{\max_j |d_j|}{d_i} \frac{|\Delta y_i|}{\|\mathbf{y}\|} \le \frac{\max_j |d_j|}{\min_j |d_j|} \frac{|\Delta y_i|}{\|\mathbf{y}\|}$$

The condition number of the linear system Dx = y is thus

$$\frac{|\Delta x_i|}{\|\mathbf{x}\|} / \frac{|\Delta y_i|}{\|\mathbf{y}\|} \le \frac{\max_j |d_j|}{\min_j |d_j|} = \kappa(\mathbf{D})$$

and can be very large, if d_i are of very different magnitudes.

Condition number — diagonal matrices — 2D example



23/67

Condition number — diagonal matrices — 2D example



Condition number

The diagonal matrix is an exemplary case, since all matrices can be transformed to diagonal matrices using the SVD. The linear system

$$\mathbf{A}\mathbf{x} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\mathsf{T}\mathbf{x} = \mathbf{y}$$

by the following transformation of variables²

$$\mathbf{\hat{x}} = \mathbf{V}^{\mathsf{T}}\mathbf{x}, \qquad \qquad \mathbf{\hat{y}} = \mathbf{U}^{\mathsf{T}}\mathbf{y}$$

is reduced to the diagonal system

$$\mathbf{\Sigma} \mathbf{\hat{x}} = \mathbf{\hat{y}}$$

The condition number of both systems is thus

$$\kappa(\mathbf{A}) = \kappa(\mathbf{\Sigma}) = \frac{\sigma_{\max}}{\sigma_{\min}}.$$

 $^2 {\rm Transformations}$ by unitary matrices have always the best-possible condition number of 1.



Condition number

Although it is easy to construct a counterexample, highly ill-conditioned matrices of small dimensions are rather rare in practice. However, ill-conditioning of medium and large matrices is much more common.

Logarithmic plot of the singular values of two 1000×1000 matrices of random numbers N(0,1) and U(0,1):



Deconvolution as a common source of ill-conditioning

In proper infinite-dimensional function spaces, the convolution with a continuous function h(t),

$$y(t) = \int_0^T h(t-s)x(s)\,\mathrm{d}s, \qquad t\in[0,\,T],$$

can be stated in the form of a simple linear operator equation $y = \mathcal{H}x$. The operator \mathcal{H} is compact³, so its inverse \mathcal{H}^{-1} is unbounded and noncontinuous. Thus, $\kappa(\mathcal{H}) = \infty$ and the problem of finding x(t), given y(t) and h(t), is extremely ill-conditioned (ill-posed).

³Intuitively, compact \approx smoothing.

Deconvolution as a common source of ill-conditioning

$$(\mathcal{H}x)(t) = \int_0^t x(t) dt \qquad (\mathcal{H}^{-1}y)(t) = y'(t)$$

Both \mathcal{H} and its inverse \mathcal{H}^{-1} are linear. However, the inverse is everywhere non-continuous.



Deconvolution as a common source of ill-conditioning

In practice, all computations are performed using finite-dimensional data, so that discretized h and y are used to compute the discretized x. If the data originate from measurements or simulations, the considered time interval [0, T] is usually sampled into N equal time steps $t_i = i\Delta t$, so that

$$y_i = \sum_{j=0}^{N-1} h_{i-j} x_j, \qquad i = 0, \dots, N-1,$$

where $y_i = y(t_i)$, $x_j = x(t_j)$ and $h_k = \Delta th(t_k)$.

The discretized equations can be stated in the matrix form as a single large finite-dimensional linear system

Deconvolution as a common source of ill-conditioning

 $\mathbf{y} = \mathbf{H}\mathbf{x},$

Regularization

Large Toeplitz systems

Reading

where

Outline

$$\mathbf{y} = [y_1, \dots, y_N]^\mathsf{T}$$
$$\mathbf{x} = [x_1, \dots, x_N]^\mathsf{T}$$

and $\mathbf{H} = [h_{ij}]$ is a diagonal-constant matrix, that is

$$h_{ij}=h_{i-j}.$$

Teoplitz matrix

Least-squares problem

Conditioning

A diagonal-constant matrix is a matrix $\mathbf{H} = [h_{ij}]$, which satisfies $h_{ij} = h_{i-j}$. Such matrices are also called Toeplitz matrices and are examples of structured matrices.

Deconvolution as a common source of ill-conditioning

Discretized version of the deconvolution problem:

$\mathbf{y}=\mathbf{H}\mathbf{x},$

where $\mathbf{H} = [h_{ij}]$ is a Toeplitz matrix, that is $h_{ij} = h_{i-j}$.

- Since the original continuous system was extremely ill-conditioned (unbounded inverse, hence $\kappa(\mathcal{H}) = \infty$ and the inverse problem is even ill-posed), the condition number $\kappa(\mathbf{H})$ grows to ∞ as the discretization time step Δt tends to zero.
- An important conclusion: finer mesh (discretization) need not lead to a better accuracy.
- The singular values of a Toeplitz matrix that "arises from the discretization of first-kind Fredholm integral equations [...] decay gradually, until they level off at a plateau approximately at the machine precision times σ_{max} (in infinite precision they would decay to zero)" (P. Ch. Hansen).

Deconvolution — example

Consider the following deconvolution problem:

$$y(t) = \int_0^t \cos(t-s)x(s)ds,$$

where y(t) can be interpreted as the velocity response of an undamped single degree of freedom system to a force excitation described by x(t).

Assume the observed response

$$y(t)=rac{1}{2}t\sin t,\qquad t\in [0,2\pi].$$

The unique exact solution is $x(t) = \sin t$. To find it numerically, the time interval $[0, 2\pi]$ has been discretized into 100 time steps.

Deconvolution — example



The singular values of **H** plotted in the logarithmic scale. The singular values decay quickly and span across 1.8 orders of magnitude, so that $\kappa(\mathbf{H}) \approx 67$.



Deconvolution — example

Three discretized right-hand sides (measurements y(t)): exact and contaminated with uncorrelated Gaussian measurement errors at 1% and 10% rms levels.



Corresponding computed solutions (the exact solution is $x(t) = \sin t$). Note the increasing noise.



Outline



O Numerical regularization

- Regularization for direct methods
- Regularization for iterative methods
- Regularization parameter

Numerical regularization

Consider a linear system $\mathbf{A}\mathbf{x} = \mathbf{y}$ with a severely ill-conditioned \mathbf{A} .

- Such a system cannot be solved accurately based on y only, no matter what algorithm is used, because even tiny inaccuracies of y are highly amplified and dominate the computed x.
- However, the accuracy can be improved if additional information about some expected or typical characteristics of x is exploited. This is called numerical regularization of the solution.
- The amount of regularization is controlled by the so-called regularization parameter α :
 - When α is too large, the regularization information tends to dominate and distort the computed regularized solution \mathbf{x}_{α} .
 - When α is too small, the measurement noise (amplified via ill-conditioning) tends to dominate and distort the computed regularized solution \mathbf{x}_{α} .

Numerical regularization

The additional information is usually related to smoothness of \mathbf{x} and takes the form of the requirement of a limited magnitude of $\|\mathbf{D}^{p}\mathbf{x}\|$, where \mathbf{D} is the matrix of the *p*th order differences,

$$\begin{split} \mathbf{D}^{0} &= \mathbf{I}, \\ \mathbf{D}^{1} &= \begin{bmatrix} -1 & 1 & 0 & \cdots & 0 \\ 0 & -1 & 1 & \cdots & 0 \\ 0 & 0 & -1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & -1 \end{bmatrix}, \\ \mathbf{D}^{2} &= \begin{bmatrix} 1 & -2 & 1 & \cdots & 0 \\ 0 & 1 & -2 & \cdots & 0 \\ 0 & 1 & -2 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{bmatrix}. \end{split}$$

Regularization for direct methods

Consider $\boldsymbol{A}\boldsymbol{x}=\boldsymbol{y}$ and the formula for its SVD-based solution,

$$\mathbf{x} = \mathbf{V}_1 \mathbf{\Sigma}_0^{-1} \mathbf{U}_1^\mathsf{T} \mathbf{y} = \left[\sum_{i=1}^{\mathsf{rank}\,\mathbf{A}} \frac{\mathbf{v}_i \mathbf{u}_i^\mathsf{T}}{\sigma_i} \right] \mathbf{y}_i$$

- The formula clearly explains the reason of error amplification: divisions by the increasingly tinier singular values.
- In many practical problems, the components of x and y, which correspond to tiny singular values, are of high frequency⁴. As a result, these components are often disproportionally affected by the noise.

⁴ The singular vectors \mathbf{u}_i and \mathbf{v}_i [of a matrix that arises from the discretization of first-kind Fredholm integral equations] have an increasing number of sign changes in their elements as *i* increases, i.e., as the corresponding singular values σ_i decrease. Often, the number of sign changes is precisely i - 1. (P. Ch. Hansen, 1995)

TSVD and standard Tikhonov regularization

However, the noisy components can be filtered out from the solution by explicit damping of the excessively amplified components,

$$\mathbf{\tilde{x}} = \left[\sum_{i=1}^{\mathsf{rank}\,\mathbf{A}} \phi(\sigma_i) \mathbf{v}_i \mathbf{u}_i^\mathsf{T}\right] \mathbf{y},$$

where $\phi(\sigma)$ is a filtering function. The most popular choices are

• Truncated Singular Value Decomposition (TSVD)

$$\phi(\sigma) = \begin{cases} \sigma^{-1} & \text{if } \sigma \ge \alpha \\ 0 & \text{otherwise} \end{cases}$$

• Standard Tikhonov regularization with the parameter α

$$\phi(\sigma) = \frac{\sigma}{\sigma^2 + \alpha^2}.$$

In fact, these methods regularize the solution by limiting $\|\mathbf{x}\|$.

Tikhonov regularization

The Tikhonov regularized solution of Ax = y is defined as

find **x**, which minimizes
$$\|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 + \alpha^2 \|\mathbf{D}\mathbf{x}\|^2$$

where $\alpha \geq 0$ is the regularization parameter. The solution is regularized with respect to the norm $\|\mathbf{D}\mathbf{x}\|$, where the matrix **D** represents the additionally available information. The solution satisfies the regularized counterpart of the normal equation,

$$\left(\mathbf{A}^{\mathsf{T}}\mathbf{A} + \alpha^{2}\mathbf{D}^{\mathsf{T}}\mathbf{D}\right)\mathbf{x} = \mathbf{A}^{\mathsf{T}}\mathbf{y},$$

which corresponds to the following least-squares problem:

minimize
$$\left\| \begin{bmatrix} \mathbf{A} \\ \alpha \mathbf{D} \end{bmatrix} \mathbf{x} - \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} \right\|^2$$

Given the SVD of the matrix **A**, the standard case of $\mathbf{D} = \mathbf{I}$ can be computed directly (see the previous slide).

Regularization for iterative methods

Iterative methods are often used for large-scale problems (direct methods are then impractical).

- The solution is approximated iteratively by computing a sequence of approximated solutions **x**_k.
- In regularizing iterative methods, the components corresponding to large singular values (low frequency) are retrieved before the components corresponding to small singular values (high frequency). Therefore, the number of iterations plays the role of the regularization parameter: the more iterations, the less regularized the solution.

Regularization for iterative methods

The conjugate gradient least squares method (CGLS) (Lectures B-3 and C-4) is probably the most popular regularizing iterative method.

• In the standard formulation, the CGLS method iteratively minimizes

$$\phi(\mathsf{x}) = \|\mathsf{y} - \mathsf{A}\mathsf{x}\|^2$$

and the regularization is with respect to $\|\mathbf{x}\|$.

• A regularization with respect to $\|\mathbf{D}\mathbf{x}\|$ can be performed via a change of variables,

$$\mathbf{w} = \mathbf{D}\mathbf{x}, \qquad \qquad \mathbf{x} = \mathbf{D}^{-1}\mathbf{w},$$

so that the method minimizes iteratively

$$\phi(\mathbf{w}) = \|\mathbf{y} - \mathbf{A}\mathbf{D}^{-1}\mathbf{w}\|^2$$

with respect to $\|\mathbf{w}\| = \|\mathbf{D}\mathbf{x}\|$.

Regularization parameter

The "amount of regularization" is controlled by the regularization parameter α :

TSVD number of truncated singular values Tikhonov weighting parameter α

iterative number of iterations

In all regularization methods, direct or iterative, the proper choice of the regularization parameter α is crucial:

- \bullet Too small α results in a computed solution being too noisy.
- If α is too large, the computed solution is overly distorted by the regularizing condition (overregularized).

Regularization parameter — example

20

-2

0

Tikhonov-regularized solution to the deconvolution example at 10% measurement error (exact solution $x(t) = \sin t$)



40

60

80

100

Regularization parameter — the L-curve method

In the absence of information about the noise level in \mathbf{y} , all methods of choice of the regularization parameter are more or less heuristic. The most popular method is the L-curve method⁵.

- The L-curve is a log-log plot of the norm of the regularization condition ||Dx|| versus the norm of the residuum ||y – Ax|| in dependence on the regularization parameter α.
- The L-curve usually consists of two branches corresponding to
 - excessive (horizontal, lower) and
 - insufficient (vertical, upper)

regularization. The corner is assumed to correspond to the proper value of α ; it is usually defined as the point of maximum curvature.

⁵For an overview and comparison of other methods, see F. Bauer, M.A. Lukas (2011) Comparing parameter choice methods for regularization of ill-posed problems. *Mathematics and Computers in Simulation*, 81(9):1795–1841. doi:10.1016/j.matcom.2011.01.016



Regularization parameter — the L-curve — example

The deconvolution example at 10% measurement error ($\mathbf{D} = \mathbf{D}^1$)



Regularization parameter — the L-curve — example

The deconvolution example at 10% measurement error, the regularized solution compared to the exact solution ($\mathbf{D} = \mathbf{D}^1$, $\alpha = 1.3$ found by the L-curve method)



Outline



- Toeplitz systems
- Common problems
- A not-so-large example $(10\,000 \times 10\,000)$



Toeplitz matrices

A matrix $\mathbf{H} = [h_{ij}]$ is called a Toeplitz matrix, if it satisfies $h_{ij} = h_{i-j}$, $\mathbf{H} = \begin{bmatrix} h_0 & h_{-1} & h_{-2} & \cdots & h_{-n+1} \\ h_1 & h_0 & h_{-1} & \cdots & h_{-n+2} \\ h_2 & h_1 & h_0 & \cdots & h_{-n+3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ h_{n-1} & h_{n-2} & h_{n-3} & \cdots & h_0 \end{bmatrix}$

Systems with Toeplitz matrices often arise from the disretization of linear integral equations in time domain, e.g. in inverse problems of input identification.

Block matrices with Toeplitz blocks

Discretization of matrix linear integral equations, e.g. in problems of input identification in $MIMO^6$ systems, may yield linear equations with a block matrix with Toeplitz blocks.



⁶Multiple inputs, multiple outputs.

Toeplitz systems — common problems

Common problems:

- Memory If many time steps are considered, such a matrix and the space required for its storage can be huge.
 - Time Most of direct operations on such large matrices can be too time-consuming.
- Accuracy Large Toeplitz matrices are usually extremely ill-conditioned.

Toeplitz matrices correspond to discrete convolutions: use frequency domain, if possible! But in transient analysis, other problems can then appear:

- spectral leakage
- windowing functions
- regularization

Memory

$$\mathbf{H} = \begin{bmatrix} h_0 & h_{-1} & h_{-2} & \cdots & h_{-n+1} \\ h_1 & h_0 & h_{-1} & \cdots & h_{-n+2} \\ h_2 & h_1 & h_0 & \cdots & h_{-n+3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ h_{n-1} & h_{n-2} & h_{n-3} & \cdots & h_0 \end{bmatrix}$$

- A general $n \times n$ matrix requires n^2 storage for all elements h_{ij} .
- An n × n Toeplitz matrix is defined by 2n − 1 elements: h_{-n+1}, ..., h_{n-1}. Such a matrix can be thus stored in a reduced form. For large matrices, the difference is tremendous: e.g. 7.5 GB vs. 176 kB for a 90000 × 90000 double matrix.

• This is typical for structured systems (and for sparse systems). In case of discretization of time-convolutions, lower-triangular Toeplitz systems may arise, which are defined by *n* elements only.



Most of direct operations on large Toeplitz matrices, including all decompositions and factorizations, would be too time-consuming. Iterative methods will yield a reasonably accurate solution in a much shorter time.



Time Quick matrix-vector multiplication

Most of iterative methods (like CGLS) make extensive use of matrix-vector multiplications. Standard multiplications requires $O(n^2)$ operations. However, Toeplitz matrices are discrete versions of convolution operators.

• In continuous time, the convolution can be performed in frequency domain by simple multiplication of spectra:

$$\mathcal{F}(f\star g)=(\mathcal{F}f)(\mathcal{F}g),$$

so that

$$(f \star g)(t) = \mathcal{F}^{-1}\left[(\mathcal{F}f)(\mathcal{F}g)\right](t). \tag{(*)}$$

 In discrete time, the fast Fourier transform (FFT) can be used, which requires only O(n log n) time. However, the discrete counterpart of (*) is valid only for *circulant matrices*, which form a subset of Toeplitz matrices.

Outline 0	Least-squares problem	Conditioning 00000000000000000000000000000000000	Regularization 00000000000000	Large Toeplitz systems	Reading 00	HW8 00
Time						

Circulant matrices

Teoplitz matrix

A matrix $\mathbf{H} = [h_{ij}]$, which satisfies $h_{ij} = h_{i-j}$, is called a Toeplitz matrix.

Circulant matrix

A Toeplitz matrix $\mathbf{H} = [h_{ij}]$, where $h_{ij} = h_{i-j}$, is called a circulant matrix, if $h_{-k} = h_{n-k}$.



Time Quick matrix-vector multiplication (circulant matrices)

An $n \times n$ circulant matrix $\hat{\mathbf{H}}$ can be multiplied by any vector $\hat{\mathbf{x}}$ in $O(n \log n)$ time instead of $O(n^2)$ time via the FFT⁷:

$$\boldsymbol{\hat{H}}\boldsymbol{\hat{x}} = \mathsf{FFT}^{-1}\left[\mathsf{FFT}(\boldsymbol{\hat{h}})\,\mathsf{FFT}(\boldsymbol{\hat{x}})\right],$$

where $\hat{\boldsymbol{h}}$ is the first column of $\hat{\boldsymbol{H}}.$



Toeplitz matrices that arise from the discretization of time convolutions are usually lower-triangular and so they are not circulant.

⁷Depending on the exact definition of the FFT, sometimes the right-hand side has to be multiplied by \sqrt{n} .

Time

Quick matrix-vector multiplication (non-circulant Toeplitz matrices)

However, non-circulant Toeplitz matrices can be augmented to become circulant:



Time

Quick matrix-vector multiplication (non-circulant Toeplitz matrices)

A product of a non-circulant Toeplitz matrix with a vector, $\mathbf{H}\mathbf{x}$, can be computed via the augmented version of the matrix by dropping the trailing zeros from the product $\hat{\mathbf{H}}\hat{\mathbf{x}}$, where $\hat{\mathbf{x}} = \begin{bmatrix} \mathbf{x}^{\mathsf{T}}\mathbf{0}^{\mathsf{T}} \end{bmatrix}^{\mathsf{T}}$, that is, where $\hat{\mathbf{x}}$ is the original vector \mathbf{x} padded with zeros.

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 \\ 3 & 2 & 1 & 1 \\ 4 & 3 & 2 & 1 \end{bmatrix} \begin{bmatrix} w \\ x \\ y \\ z \end{bmatrix} \sim \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 4 & 3 & 2 \\ 2 & 1 & 0 & 0 & 0 & 0 & 4 & 3 \\ 3 & 2 & 1 & 0 & 0 & 0 & 0 & 4 \\ 4 & 3 & 2 & 1 & 0 & 0 & 0 & 0 \\ 0 & 4 & 3 & 2 & 1 & 0 & 0 & 0 \\ 0 & 0 & 4 & 3 & 2 & 1 & 0 & 0 \\ 0 & 0 & 0 & 4 & 3 & 2 & 1 & 0 \\ 0 & 0 & 0 & 0 & 4 & 3 & 2 & 1 & 0 \\ 0 & 0 & 0 & 0 & 4 & 3 & 2 & 1 & 0 \\ 0 & 0 & 0 & 0 & 4 & 3 & 2 & 1 & 0 \\ \end{bmatrix} \begin{bmatrix} w \\ x \\ y \\ z \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Accuracy

Large Toeplitz systems are usually extremely ill-conditioned

- use a regularizing iterative method, like CGLS, which will retrieve well-conditioned components before the ill-conditioned components
- The number of iterations will play the role of the regularization parameter

the more iterations, the less regularized the solution

A not-so-large example $(10\,000 \times 10\,000)$

Consider the deconvolution problem:

$$y(t) = \int_0^t \cos(t-s)x(s)ds,$$

where y(t) is the velocity response of an undamped single DOF system to a force excitation x(t).

Assume the observed response

$$y(t) = rac{1}{2}t\sin t, \qquad t\in [0,10 imes 2\pi].$$

The unique exact solution is $x(t) = \sin t$. To find it numerically, the time interval $[0, 10 \times 2\pi]$ has been discretized into 10000 time steps.

A not-so-large example $(10\,000 \times 10\,000)$

Two discretized right-hand sides (measurements y(t)): exact and contaminated with uncorrelated Gaussian measurement errors at 50% rms level.



The exact solution is $x(t) = \sin t$. The solution computed in the noisy case (without regularization, here the 10 000 iteration) would be useless.



A not-so-large example $(10\,000 \times 10\,000)$

Residuum norms of the iterates

A common stop condition for CGLS is based on the norm of the residuum $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|$.





Orange dots mark the iterates no. $25, 50, 75, \ldots, 300$.

A not-so-large example $(10\,000 \times 10\,000)$

CGLS iterates



Outline



Further reading

- Linear least-squares problems
 - G. Dahlquist, Å. Björck, Linear Least Squares Problems, [in:] *Numerical Methods in Scientific Computing*, vol. 2.
- Conditioning, regularization
 - P.Ch. Hansen, Discrete inverse problems: Insight and Algorithms, SIAM 2010.
 - P.Ch. Hansen, Deconvolution and regularization with Toeplitz matrices, Numerical Algorithms 29:323–378, 2002.
 - F. Bauer, M.A. Lukas, Comparing parameter choice methods for regularization of ill-posed problems. Mathematics and Computers in Simulation 81(9):1795–1841, 2011.
- Large Toeplitz systems:
 - I. Gohberg, V. Olshevsky, Fast algorithms with preprocessing for matrix-vector multiplication problems, J. Complexity, 10(4):411–427, 1994.
 - I. Gohberg, V. Olshevsky, Complexity of multiplication with vectors for structured matrices, Linear Algebra Appl., 202:163–192, 1994.

Outline



Homework 8 (25 points)

Regularization and iterative linear solvers

Available soon at http://info.ippt.pan.pl/~ljank.

E-mail the answer and the source code to ljank@ippt.pan.pl.