# Chapter 3

# Applications

- Least-squares

- Least-norm

- Total least-squares

- Low-rank approximation

The first three sections are discuss the linear system of equations $Ax = y$. The matrix $A \in \mathbb{R}^{m \times n}$ and the vector $y \in \mathbb{R}^m$ are given data. The vector $x \in \mathbb{R}^n$ is an unknown. Assuming that $A$ is full rank, the system $Ax = y$ is called

- *overdetermined* if $m > n$ (in this case it has more equations than unknowns) and

- *underdetermined* if $m < n$ (in this case it has more unknowns than equations).

For most vectors $y \in \mathbb{R}^m$, an overdetermined system has no solution $x$, and for any $y \in \mathbb{R}^m$ an underdetermined system has infinitely many solutions $x$. In the case of an overdetermined system, it is of interested to find an approximate solution. An important example is the least squares approximate solution, which minimizes the 2-norm of the equation error.

In the case of an underdetermined system, it is of interested to find a particular solution. The least-norm solution is an example of a particular solution, It minimizes the 2-norm of the solution. Note that the least-squares approximate solution is (most of the time) not a solution, while the least-norm solution is (aways) one of infinitely many solutions.

## 3.1   Least-squares

The least-squares method for solving approximately an overdetermined system $Ax = y$ of equations is defined as follows. Choose $x$ such that the 2-norm of the residual (equation error)
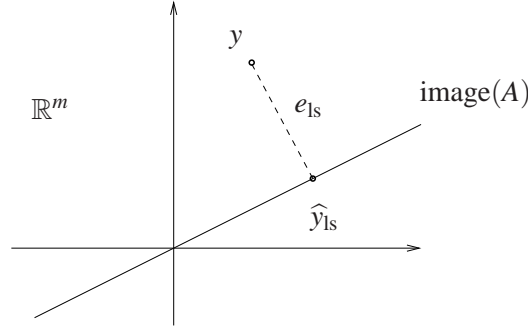
$$e(x) := y - Ax$$

is minimized. A minimizer

$$\widehat{x}_{\mathrm{ls}} := \arg\min_x \| \underbrace{y - Ax}_{e(x)} \|_2 \tag{3.1}$$

is called a *least-squares approximate solution* of the system $Ax = y$.

A geometric interpretation of the least-squares approximation problem (3.1) projection of $y$ onto the image of $A$.

Here $\widehat{y}_{\mathrm{ls}} := A\widehat{x}_{\mathrm{ls}}$ is the projection, which is the least-squares approximation of $y$ and $e_{\mathrm{ls}} := \widehat{y}_{\mathrm{ls}} - A\widehat{x}_{\mathrm{ls}}$ is the approximation error.

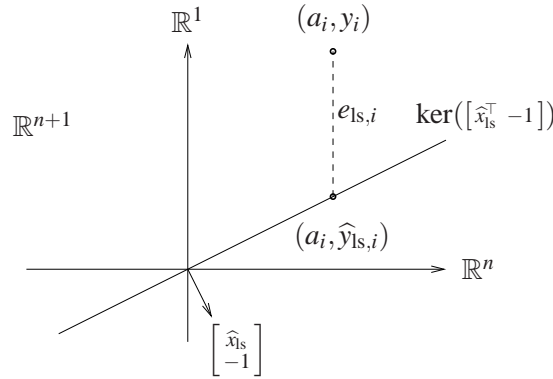Let $a_i$ be the $i$th row of $A$. We refer to the vector $\mathrm{col}(a_i, y_i)$ as a "data point". We have,

$$A\widehat{x}_{\mathrm{ls}} = \widehat{y}_{\mathrm{ls}} \quad \Longleftrightarrow \quad \begin{bmatrix} A & \widehat{y}_{\mathrm{ls}} \end{bmatrix} \begin{bmatrix} \widehat{x}_{\mathrm{ls}} \\ -1 \end{bmatrix} = 0$$

$$\Longleftrightarrow \quad \begin{bmatrix} a_i & \widehat{y}_{\mathrm{ls},i} \end{bmatrix} \begin{bmatrix} \widehat{x}_{\mathrm{ls}} \\ -1 \end{bmatrix} = 0, \quad \text{for } i = 1,\ldots,m$$

so that for all $i$, $(a_i, \widehat{y}_{\mathrm{ls},i})$ lies on the subspace perpendicular to $(\widehat{x}_{\mathrm{ls}}, -1)$. $(a_i, \widehat{y}_{\mathrm{ls},i})$ is an the least-squares approximation of the $i$ data point $\mathrm{col}(a_i, y_i)$.

$$(a_i, \widehat{y}_{\mathrm{ls},i}) = (a_i, \widehat{y}_{\mathrm{ls},i}) + (0, e_{\mathrm{ls},i}),$$

and $(0, e_{\mathrm{ls},i})$ is the least-squares approximation error. Note that $e_{\mathrm{ls},i}$ is the vertical distance from $(a_i, y_i)$ to the subspace.

The above derivation suggestions another geometric interpretation of the least-squares approximation.



Note that the former geometric interpretation is in the space $\mathbb{R}^m$, while the latter is in the (data space) $\mathbb{R}^{n+1}$.

*Exercise problem* 49. [Derivation of solution $x_{\mathrm{ln}}$ via Lagrange multipliers] Assuming that $m \geq n = \mathrm{rank}(A)$, i.e., $A$ is full column rank, show that

$$\widehat{x}_{\mathrm{ls}} = (A^\top A)^{-1} A^\top y.$$

$\square$

Notes:

- $A_{\mathrm{ls}} := (A^\top A)^{-1} A^\top$ is a left-inverse of $A$

- $\widehat{x}_{\mathrm{ls}}$ is a linear function of $y$ (given by the matrix $A_{\mathrm{ls}}$)

- If $A$ is square, $\widehat{x}_{\mathrm{ls}} = A^{-1} y$ (i.e., $A_{\mathrm{ls}} = A^{-1}$)

- $\widehat{x}_{\mathrm{ls}}$ is an exact solution if $Ax = y$ has an exact solution

- $\widehat{y}_{\mathrm{ls}} := A\widehat{x}_{\mathrm{ls}} = A(A^\top A)^{-1} A^\top y$ is a least-squares approximation of $y$

**Projector onto the image of *A* and orthogonality principle**

The $m \times m$ matrix

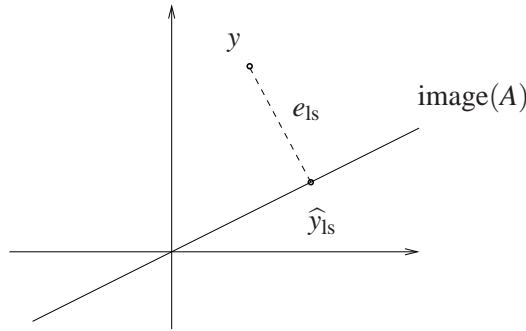$$\Pi_{\text{image}(A)} := A(A^\top A)^{-1} A^\top$$

is the orthogonal projector onto the subspace $\mathscr{L} := \text{image}(A)$. Suppose that the columns of $A$ form an orthonormal basis for $\mathscr{L}$. Then, recall that $\Pi_{\text{image}(Q)} := AA^\top$.

The least-squares residual vector

$$e_{\text{ls}} := y - A\widehat{x}_{\text{ls}} = \underbrace{\left( I_m - A(A^\top A)^{-1} A^\top \right)}_{\Pi_{(\text{image}(A))^\perp}} y$$

is orthogonal to $\text{image}(A)$

$$\langle e_{\text{ls}}, A\widehat{x}_{\text{ls}} \rangle = y^\top \left( I_m - A(A^\top A)^{-1} A^\top \right) A\widehat{x}_{\text{ls}} = 0. \tag{3.2}$$



*Exercise problem* 50. Show that the orthogonality condition (3.2) is a necessary and sufficient condition for $\widehat{x}_{\text{ls}}$ being a least squares approximate solution to $Ax = b$.

□

**Least-squares via QR factorization**

Let $A = QR$ be the QR factorization of $A$. We have,

$$(A^\top A)^{-1} A^\top = (R^\top Q^\top QR)^{-1} R^\top Q^\top$$
$$= (R^\top Q^\top QR)^{-1} R^\top Q^\top = R^{-1} Q^\top,$$

so that

$$\widehat{x}_{\text{ls}} = R^{-1} Q^\top y \quad \text{and} \quad \widehat{y}_{\text{ls}} := Ax_{\text{ls}} = QQ^\top y.$$

*Exercise problem* 51 (Least-squares with an increasing number of columns in *A*). Let $A =: \begin{bmatrix} a_1 & \cdots & a_n \end{bmatrix}$ and consider the sequence of least squares problems

$$A^i x^i = y, \qquad \text{where } A^i := \begin{bmatrix} a_1 & \cdots & a_i \end{bmatrix}, \quad \text{for } i = 1, \ldots, n$$

Define $R_i$ as the leading $i \times i$ submatrix of $R$ and let $Q_i := \begin{bmatrix} q_1 & \cdots & q_i \end{bmatrix}$. Show that

$$\widehat{x}_{\text{ls}}^i = R_i^{-1} Q_i^\top y.$$

□

**Weighted least-squares**

Given a positive definite matrix $W \in \mathbb{R}^{m \times m}$, define the wighted 2-norm

$$\|e\|_W^2 := e^\top W e.$$

and the weighted least-squares approximate solution

$$\widehat{x}_{W,\mathrm{ls}} := \arg\min_x \|y - Ax\|_W^2.$$

*Exercise problem* 52. Show that
$$\widehat{x}_{W,\mathrm{ls}} = (A^\top W A)^{-1} A^\top W y,$$

and that the least-squares orthogonality principle holds for the weighted least-squares problem as well by replacing the inner product $\langle e, y \rangle$ by the weighted inner product

$$\langle e, y \rangle_W := e^\top W y.$$

□

**Recursive least-squares**

The least-squares criterion is

$$\|y - Ax\|_2^2 = \sum_{i=1}^m (y_i - a_i^\top x)^2$$

where $a_i^\top$ is the $i$th row of $A$. We consider the sequence of least-squares problems

$$\text{minimize} \quad \sum_{i=1}^k (y_i - a_i^\top x)^2$$

the solutions of which are

$$\widehat{x}_{\mathrm{ls}}(k) := \left( \sum_{i=1}^k a_i a_i^\top \right)^{-1} \sum_{i=1}^m a_i y_i.$$

The meaning is that the measurements $(a_i, y_i)$ come sequentially (in time) and we aim to compute a solution each time a new data point arrives. Instead of recomputing the solution from scratch, we can recursively update $\widehat{x}_{\mathrm{ls}}(k-1)$ in order to obtain $\widehat{x}_{\mathrm{ls}}(k)$.

Recursive algorithm

- Initialization: $P(0) = 0 \in \mathbb{R}^{n \times n}$, $q(0) = 0 \in \mathbb{R}^n$

- For $m = 0, 1, \ldots, m$

- $P(k+1) := P(k) + a_{k+1} a_{k+1}^\top$, $q(k+1) := q(k) + a_{k+1} y_{k+1}$

- If $P(k)$ is invertible, $\widehat{x}_{\mathrm{ls}}(k) = P^{-1}(k) q(k)$.

On each step, the algorithm requires inversion of an $n \times n$ matrix, which requires $O(n^3)$ operations. At certain $k$, $P(k)$ being invertible implies that $P(k')$ is invertible, for all $k' > k$.

The computational complexity of the algorithm can be decreased to $O(n^2)$ operations per step by using the following result about the inverse of matrix with rank-1 update

$$(P + aa^\top)^{-1} = P^{-1} - \frac{1}{1 + a^\top P^{-1} a} (P^{-1} a)(P^{-1} a)^\top.$$

## Multiobjective least-squares

Least-squares minimizes the cost function

$$J_1(x) := \|Ax - y\|_2^2.$$

Consider a second cost function

$$J_2(x) := \|Bx - z\|_2^2,$$

which we want to minimize together with $J_1$. Usually the criteria $\min_x J_1(x)$ and $\min_x J_2(x)$ are competing. A common example is $J_2(x) := \|x\|_2^2$ — minimize $J_1$ with small $x$.
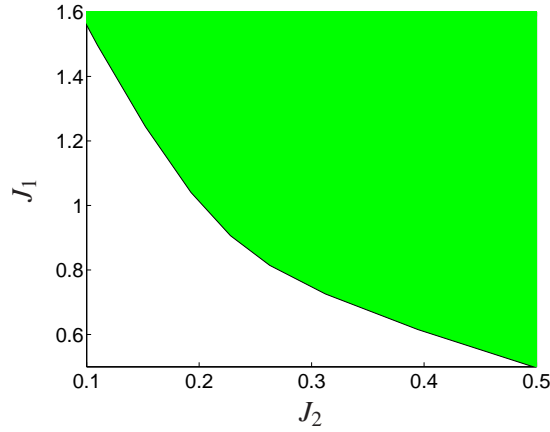
The set of achievable objectives is

$$\{(\alpha, \beta) \in \mathbb{R}^2 \mid \exists x \in \mathbb{R}^n \text{ subject to } J_1(x) = \alpha, \; J_2(x) = \beta\}$$

Its boundary is the optimal trade-off curve and the corresponding $x$'s are called *Pareto optimal*.

A common method for "solving" multiobjective optimization problems is secularization. For any $\mu \geq 0$, the problem

$$\widehat{x}(\mu) = \arg\min_x J_1(x) + \mu J_2(x)$$

produces a Pareto optimal point. For a convex problem (such as the the multiobjective least-squares), by varying $\mu \in [0, \infty)$, $\widehat{x}(\mu)$ sweeps all Pareto optimal solutions.



## Regularized least-squares

*Exercise problem* 53. Show that the solution of the *Tychonov regularization* problem

$$\widehat{x}_{\text{reg}} = \arg\min_x \|Ax - b\|_2^2 + \mu \|x\|_2^2$$

is

$$\widehat{x}_{\text{reg}} = (A^\top A + \mu I_n)^{-1} A^\top y.$$

□

Note that $\widehat{x}_{\text{reg}}$ exists for any $\mu > 0$, independent on size and rank of $A$. The parameter $\mu$ controls the trade-off between

- fitting accuracy $\|Ax - b\|_2$, and

- solution size $\|x\|_2$.

For small $\mu$, the solution is larger but gives better fit. For large $\mu$, the solution is smaller but the fit is worse. In the extreme case $\mu = 0$, assuming that the system $Ax = b$ is overdetermined, the regularized least-squares problem is equivalent to the standard least-squares problem, which does not constrain the size of $x$. In the other extreme $\mu \to 0$, assuming that $Ax = b$ is underdetermined, the regularized least-squares problem tends to the least-norm problem.

## 3.2  Least-norm

Consider an underdetermined system $Ax = y$, with full rank $A \in \mathbb{R}^{m \times n}$. The set of solutions is

$$\mathscr{A} := \{ x \in \mathbb{R}^n \mid Ax = y \} = \{ x_p + z \mid z \in \ker(A) \} = x_p + \ker(A).$$

where $x_p$ is a particular solution, i.e., $Ax_p = y$. The least-norm solution is defined by the optimization problem
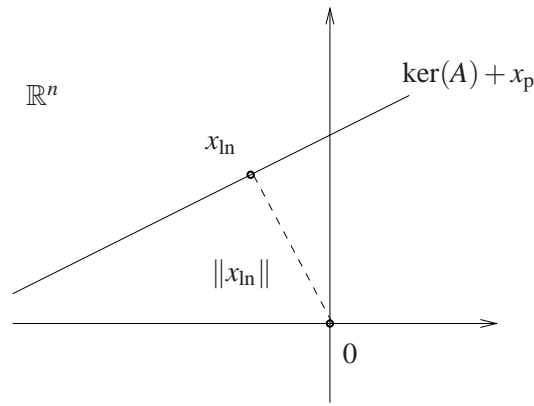
$$x_{\text{ln}}^2 := \arg\min_x \|x\|_2 \quad \text{subject to} \quad Ax = y. \tag{3.3}$$

*Exercise problem* 54 (Derivation of solution $x_{\text{ln}}$ via Lagrange multipliers).  Assuming that $n \geq m = \text{rank}(A)$, i.e., $A$ is full row rank, show that

$$x_{\text{ln}} = A^\top (AA^\top)^{-1} y.$$

$\square$

A geometric interpretation of (1.3) is the projection of 0 onto the solution set $\mathscr{A}$.



*Exercise problem* 55.  The orthogonality principle for least-norm is $x_{\text{ln}} \perp \ker(A)$.  Show that it is a necessary and sufficient condition for optimality of $x_{\text{ln}}$

$\square$

Let $A^\top = QR$ be the QR factorization of $A^\top$.  The right inverse of $A$ is

$$A^\top (AA^\top)^{-1} = QR(R^\top Q^\top QR)^{-1} = Q(R^\top)^{-1},$$

so that

$$x_{\text{ln}} = Q(R^\top)^{-1} y.$$

## 3.3  Total least-squares

The least-squares method minimizes the 2-norm of the equation error $e(x) := y - Ax$

$$\min_{x,e} \|e\|_2 \quad \text{subject to} \quad Ax = y - e$$

Alternatively, the equation error $e$ can be viewed as a correction on $y$.  The total least-squares method is motivated by the asymmetry of the least-squares method: both $A$ and $b$ are given data, but only $b$ is corrected.  The total least squares problem is defined by the optimization problem

$$\text{minimize}_{x,\widetilde{A},\widetilde{y}} \quad \left\| \begin{bmatrix} \widetilde{A} & \widetilde{y} \end{bmatrix} \right\|_F \quad \text{subject to} \quad (A + \widetilde{A})x = y + \widetilde{y}$$

Here $\widetilde{A}$ is the correction on $A$ and $\widetilde{y}$ is the correction on $y$.  The Frobenius norm $\|C\|_F$ of $C \in \mathbb{R}^{m \times n}$ is defined as

$$\|C\|_F := \sqrt{\sum_{i=1}^{m} \sum_{j=1}^{n} c_{ij}^2}.$$

## Geometric interpretation of the total least squares criterion

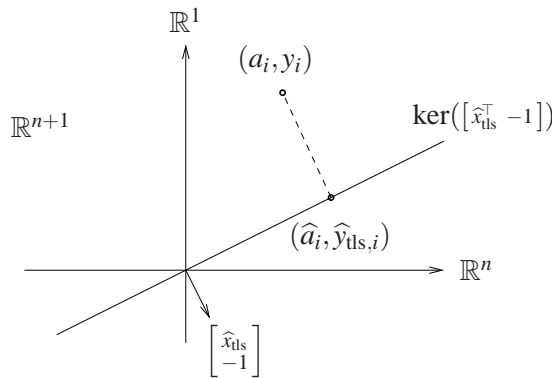In the case $n = 1$, the problem of solving approximately $Ax = y$ is

$$\begin{bmatrix} a_1 \\ \vdots \\ a_m \end{bmatrix} x = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}, \qquad \text{where} \quad x \in \mathbb{R}. \tag{3.4}$$

A geometric interpretation of the total least squares problem (3.4) is: fit a line

$$\mathscr{L}(x) := \{ (a,b) \mid ax = b \}$$

passing through the origin to the points $(a_1, y_1), \ldots, (a_m, y_m)$.

- least squares minimizes sum of squared *vertical* distances from $(a_i, y_i)$ to $\mathscr{L}(x)$,

- total least squares minimizes sum of squared *orthogonal* distances from $(a_i, y_i)$ to $\mathscr{L}(x)$.



## Solution of the total least squares problem

**Theorem 56.** *Let* $\begin{bmatrix} A & y \end{bmatrix} = U\Sigma V^\top$ *be the SVD of the data matrix* $\begin{bmatrix} A & y \end{bmatrix}$ *and*

$$\Sigma := \mathrm{diag}(\sigma_1, \ldots, \sigma_{n+1}), \quad U := \begin{bmatrix} u_1 & \cdots & u_{n+1} \end{bmatrix}, \quad V := \begin{bmatrix} v_1 & \cdots & v_{n+1} \end{bmatrix}.$$

*A total least squares solution exists if and only if* $v_{n+1,n+1} \neq 0$ *(last element of* $v_{n+1}$*) and is unique if and only if* $\sigma_n \neq \sigma_{n+1}$.

*In the case when a total least squares solution exists and is unique, it is given by*

$$\widehat{x}_{\mathrm{tls}} = -\frac{1}{v_{n+1,n+1}} \begin{bmatrix} v_{1,n+1} \\ \vdots \\ v_{n,n+1} \end{bmatrix}$$

*and the corresponding total least squares corrections are*

$$\begin{bmatrix} \widetilde{A}_{\mathrm{tls}} & \widetilde{y}_{\mathrm{tls}} \end{bmatrix} = -\sigma_{n+1} u_{n+1} v_{n+1}^\top.$$

## 3.4 Low-rank approximation

The low-rank approximation problem is defined as: Given a matrix $A \in \mathbb{R}^{m \times n}$, $m \geq n$, and an integer $r$, $0 < r < n$, find

$$\widehat{A}^* := \arg\min_{\widehat{A}} \| A - \widehat{A} \| \quad \text{subject to} \quad \mathrm{rank}(\widehat{A}) \leq r. \tag{3.5}$$

$\widehat{A}^*$ is an optimal rank-$r$ approximation of $A$ with respect to the norm $\| \cdot \|$, e.g.,

$$\| A \|_{\mathrm{F}}^2 := \sum_{i=1}^{m} \sum_{j=1}^{n} a_{ij}^2 \qquad \text{or} \qquad \| A \|_2 := \max_x \frac{\| Ax \|_2}{\| x \|_2}$$

**Theorem 57** (Solution via SVD). *Let $A = U\Sigma V^\top$ be the SVD of A and define*

$$U =: \begin{bmatrix} \overset{r}{U_1} & \overset{r-n}{U_2} \end{bmatrix} \; n \;, \quad \Sigma =: \begin{bmatrix} \overset{r}{\Sigma_1} & \overset{r-n}{0} \\ 0 & \Sigma_2 \end{bmatrix} \begin{matrix} r \\ r-n \end{matrix} \quad and \quad V =: \begin{bmatrix} \overset{r}{V_1} & \overset{r-n}{V_2} \end{bmatrix} \; n \;.$$

*An solution to (3.5) is*

$$\widehat{A}^* = U_1 \Sigma_1 V_1^\top.$$

*It is unique if and only if $\sigma_r \neq \sigma_{r+1}$.*

## 3.5  Notes and references

Least-squares and least-norm are standard topics in both numerical linear algebra and engineering. Numerical aspects of the problem are considered in [Bjö96]. For an overview of total least squares problem, see [MV07]

## Bibliography

[Bjö96]  Å. Björck. *Numerical Methods for Least Squares Problems*. SIAM, 1996.

[MV07]  I. Markovsky and S. Van Huffel. Overview of total least squares methods. *Signal Processing*, 87:2283–2302, 2007.