James Anderson, Columbia University, E6602

2: Least-squares

• approximate solutions of overdetermined systems

- projection and orthogonality
- QR decomposition
- BLUE property
- least-norm solutions of underdetermined systems

Overdetermined system of linear equations

consider the system of equations

$$y = Ax$$
, where $A \in \mathbb{R}^{m \times n}$ with $m > n$



- more equations (rows of A, y) than unknowns (x)
- typically no solution: no x s.t. y = Ax

instead, find x s.t. $y \approx Ax$:

- define the **residual** r = Ax y
- make r small by minimizing ||r||
- denote x_{ls} as the x that minimizes ||r||: least-squares approximate solution

Norms

the **Euclidean** norm of a vector $z \in \mathbb{R}^n$

$$||z|| = \sqrt{z_1^2 + z_2^2 + \dots + z_n^2} = \sqrt{z^T z}$$

- ||z|| measures the length of a vector (from the origin)
- $\|z x\|$ measures the distance between vectors z and x

more generally, a norm is any function $f: V \to \mathbb{R}$ that satisfies

(1)
$$f(x) \ge 0$$
 for all $x \in V$ with $f(x) = 0$ if and only if $x = 0$

2
$$f(x+y) \leq f(x) + f(y)$$
 for all $x, y \in V$

3
$$f(\lambda x) = |\lambda| f(x)$$
 for all $\lambda \in \mathbb{C}$ and $x \in V$

when satisfied, we replace f with $\|\cdot\|$

Least-squares approximate solution

• for simplicity, assume rank(A) = n

• finding x_{ls} is an unconstrained optimization problem:

$$||r||^{2} = x^{T} A^{T} A x - 2y^{T} A x + y^{T} y$$

• set gradient w.r.t. x to zero:

$$\nabla_x \|r\|^2 = 2A^T A x - 2A^T y = 0$$

• produces the normal equations:

$$A^T A x = A^T y$$

• by assumption the inverse exists, giving

$$x_{\rm ls} = (A^T A)^{-1} A^T y$$

Geometry

- Ax_{ls} is the point in range(A) closest to y
- Ax_{ls} is the **projection** of y onto the subspace range(A)
- the projection is **orthogonal** (see later)



the range of $A \in \mathbb{R}^{m \times n}$ is defined as

$$\operatorname{range}(A) = \{Ax \mid x \in \mathbb{R}^n\} \subseteq \mathbb{R}^m$$

• the set of vectors that can be reached by applying A to vectors in \mathbb{R}^n

in summary, in general there is no x such that y = Ax

• the least-squares approximate solution is

$$x_{\rm ls} = (A^T A)^{-1} A^T y$$

provides the smallest residual error in the Euclidean norm

- $x_{\rm ls}$ is a linear function of y
- if A is square (m = n), then $x_{ls} = A^{-1}y$

• if
$$y \in \operatorname{range}(A)$$
, then $y = Ax_{ls}$

- $A^{\dagger} = (A^T A)^{-1} A^T$ is called the pseudo-inverse of A
- A[†] is a left inverse of (full rank, tall) A, *i.e.*, it satisfies

$$A^{\dagger}A = I$$

• many other left inverses exist: $BA = I \iff B$ is a left inverse

Orthogonality

the **unsigned angle** between two vectors x and y in \mathbb{R}^n is defined as

$$\theta = \cos^{-1}\left(\frac{x^T y}{\|x\| \|y\|}\right)$$

thus $x^Ty = \|x\|\|y\|\cos\theta$

the term x^Ty defines an inner product on $\mathbb{R}^n,$ usually denoted by $\langle x,y\rangle$

special cases:

- x and y are aligned: $\theta = 0$, then $x^T y = ||x|| ||y||$
- x and y are antialigned: $\theta = \pi$, then $x^T y = -\|x\|\|y\|$
- x and y are orthogonal: $\theta = \pm \pi/2$, then $x^T y = 0$

orthogonal vectors are often denoted by $x \perp y$

Projection onto a subspace

 Ax_{ls} is the point in range(A) closest to y, it is the projection of y onto range(A):

$$Ax_{ls} = \mathcal{P}_{range(A)}(y)$$

• $\mathcal{P}_{\operatorname{range}(A)}$ is a linear function:

$$\mathcal{P}_{\mathrm{range}(A)}(y) = Ax_{\mathrm{ls}} = A(A^T A)^{-1} A^T y$$

• the matrix $A(A^TA)^{-1}A^T$ is called the **projection matrix** (associated with range(A)) or **projector**

properties of projection matrices

- If P is a projecor (onto \mathcal{V}), then I P is a projector matrix
- I P projects on \mathcal{V}^{\perp}

•
$$P^2 = P$$
 and so $(I - P)P = 0 = P(I - P)$

Orthogonal projections

the optimal residual

$$r = Ax_{ls} - y = (A(A^T A)^{-1}A^T - I)y$$

is orthogonal to $\operatorname{range}(A)$: to see this, observe that for all $z \in \mathbb{R}^n$,

$$\langle r, Az \rangle = y^T (A(A^T A)^{-1} A^T - I)^T Az = 0, \quad \text{i.e., } r \perp \text{range}(A)$$



properties of orthogonal projectors

• a projector P provides an orthogonal projection if and only if $P = P^T$

Optimality of $x_{\rm ls}$

for any $x \in \mathbb{R}^n$, we have

$$||Ax - y||^{2} = ||(Ax_{ls} - y) + A(x - x_{ls})||^{2}$$
$$= ||Ax_{ls} - y||^{2} + ||A(x - x_{ls})||^{2}$$

therefore, for any $x\neq x_{\mathrm{ls}}, \, \|Ax-y\|>\|Ax_{\mathrm{ls}}-y\|$

show using properties of norms and the fact that $r \perp \operatorname{range}(A)$

QR decomposition

given $H \in \mathbb{R}^{n \times k}$, $n \ge k$, $\operatorname{rank}(H) = k$, then H = QR



- dimensions: $Q \in \mathbb{R}^{n imes k}$, $R \in \mathbb{R}^{k imes k}$
- $Q^T Q = I_k$, R is upper-triangular
- *R_{ii}*'s are non-zero
- most definitions require $R_{ii} > 0$, in this case, Q and R are unique
- express Q in terms of it's columns: $Q = \begin{bmatrix} q_1 & q_2 & \dots & q_k \end{bmatrix}$, the column vectors form an **orthonormal** set:

$$\|q_i\| = 1, \quad q_i^T q_j = 0 \text{ if } i \neq j$$

• columns of Q form an **orthonormal basis** for range(H)

Full QR decomposition

the full QR decomposition of a full-rank matrix $H \in \mathbb{R}^{n \times k}$ is defined as

$$H = \begin{bmatrix} Q_1 & Q_2 \end{bmatrix} \begin{bmatrix} R_1 \\ 0 \end{bmatrix}, \text{ where } [Q_1 \ Q_2] \in \mathbb{R}^{n \times n} \text{ is orthogonal},$$

and R_1 is upper-triangular and invertible



- $H = Q_1 R_1$ is the QR factorization defined previously
- the additional n k columns that form Q_2 are orthogonal to range(H)
- >>[Q,R]=qr(H) %gives full qr factorization
- >>[Q,R]=qr(H,0) %gives qr factorization

Notes

- the above QR decompositions assume H is full-rank and tall
- rank deficiency produces non-invertible $R(R_1)$
- every tall $m \times n$ matrix has a full QR decomposition, and hence a QR decomposition
- every full-rank, tall matrix has a unique QR decomposition with an invertible R
- the projector $A(A^TA)^{-1}A^T$ written in terms of the QR factorization is QQ^T

Solving least-squares via QR decomposition

assuming A is full rank and tall, rewrite x_{ls} in terms of QR factorization A = QR

$$\begin{aligned} x_{\rm ls} &= (A^T A)^{-1} A^T y \\ &= ((QR)^T (QR))^{-1} (QR)^T y \\ &= R^{-1} Q^T y \end{aligned}$$

Algorithm

- **1** compute QR factorization A = QR
- **2** compute $w = Q^T y$
- \odot solve Rx = w // triangular system, use back substitution

there are many ways to compute x_{ls} :

>>[Q,R]= qr(A,O) >>xls = R\(Q'*y)

Least-squares via full QR decomposition

full QR factorization of A:

$$A = \begin{bmatrix} Q_1 & Q_2 \end{bmatrix} \begin{bmatrix} R_1 \\ 0 \end{bmatrix}, \quad Q = \begin{bmatrix} Q_1 & Q_2 \end{bmatrix} \in \mathbb{R}^{m \times m} \text{ is orthogonal},$$

• multiplication by an orthogonal matrix doesn't affect the norm:

$$||Uz||^2 = z^T U^T Uz = z^T z = ||z||^2$$

• substitute QR factorization to get

$$\begin{split} \|Ax - y\|^2 &= \left\| [Q_1 \ Q_2] \left[\begin{array}{c} R_1 \\ 0 \end{array} \right] x - y \right\|^2 \\ &= \left\| [Q_1 \ Q_2]^T [Q_1 \ Q_2] \left[\begin{array}{c} R_1 \\ 0 \end{array} \right] x - [Q_1 \ Q_2]^T y \right\|^2 \\ &= \left\| \left[\begin{array}{c} R_1 x - Q_1^T y \\ -Q_2^T y \end{array} \right] \right\|^2 \\ &= \|R_1 x - Q_1^T y\|^2 + \|Q_2^T y\|^2 \end{split}$$

recall, least-squares objective is to minimize $||r||^2$

• shown that
$$||r||^2 = ||Ax - y||^2 = ||R_1x - Q_1^Ty||^2 + ||Q_2^Ty||^2$$

no optimization performed yet - this involves selecting x:

• clearly selecting $x = x_{ls} = R_1^{-1}Q_1^Ty$ achieves optimality

the residual at $x_{\rm ls}$ is

$$Ax_{\rm ls} - y = -Q_2 Q_2^T y$$

Example: linear regression

objective: estimate the weight of a fish based on its width

- sample size N = 10
- data consists of points (x_i, y_i) where x_i is the diagonal width and y_i the weight



suggests an approximately linear relationship: $y \approx mx + b$, find m, b that minimize

$$\underset{m,b}{\mathsf{minimize}} \quad \sum_{i=1}^{N} (mx_i + b - y_i)^2$$

solving the least squares problem provides estimates $(\bar{m}, \bar{b}) = (36.8, -662.3)$



line of best fit, y = 36.8x - 662.3 minimizes the sum of the square of the residuals

Estimation

many signal reconstruction and estimation problems take the form

$$y = Ax + v, \quad A \in \mathbb{R}^{m \times n}, \quad m \ge n$$

- $x \in \mathbb{R}^n$ is the vector we want to estimate/recover/reconstruct
- $y \in \mathbb{R}^m$ is the sensor measurement/observable
- $v \in \mathbb{R}^m$ is the unknown measurement error/noise

If v was known exactly, or v = 0, then any left inverse of A recovers x

$$y = Ax + v = Ax \quad \Rightarrow \quad By = BAx = x$$

one such choice is $B = A^{\dagger}$

Least-squares estimation

choose an estimate \hat{x} (of x) that minimizes

 $\|A\hat{x} - y\|$

captures the magnitude of the difference between

- the observation
- what would be observed if we could run the model with no noise/error

a linear estimator is any mapping from y to \hat{x} that can be written as $\hat{x} = By$

• least squares estimate $\hat{x} = (A^T A)^{-1} A^T y$ is one example

Best Linear Unbiased Estimate (BLUE) property

• linear: map from observation to estimate is represented by a linear function

$$f(\alpha x + \beta y) = \alpha f(x) + \beta f(y), \quad \alpha, \beta \in \mathbb{R}, \quad x, y \in \mathbb{R}^{d}$$

• unbiased: no estimation error when there is no noise:

 $\hat{x} = x$ when v = 0

estimation error of an estimator is $x - \hat{x}$, for **unbiased linear estimators**:

$$x - \hat{x} = x - By = -Bv$$

clearly, we should make B "small" subject to the constraint BA = I

• **best**: $A^{\dagger} = (A^T A)^{-1} A^T$ is the smallest left inverse of A in the sense that: for any B such that BA = I, it can be shown that

$$||B||_F^2 \ge ||A^{\dagger}||_F^2$$

Underdetermined systems of linear equations

consider the system of equations

$$y = Ax$$
, where $A \in \mathbb{R}^{m \times n}$ with $m < n$



- more unknowns (x) than equations (rows of A, b)
- x is underspecified; many (infinite) choices of x satisfy the equation
- A has a nontrivial nullspace

assume that $\operatorname{rank}(A) = m$, so for every $y \in \mathbb{R}^m$, there is a solution

Nullspace

the **nullspace** of the matrix $A \in \mathbb{R}^{m \times n}$ is defined as

 $\operatorname{null}(A) = \{ x \in \mathbb{R}^n \mid Ax = 0 \} \subseteq \mathbb{R}^n$

- the set of vectors mapped to 0 by the map Az
- $\operatorname{null}(A)$ characterizes ambiguity in x when y = Ax:

$$y = Ax \text{ and } z \in \text{null}A \quad \Rightarrow \quad y = A(x+z)$$

• conversely, if x and \bar{x} are solutions, then $\bar{x} = x + z$ with $z \in null(A)$

all dimensions are related:

$$n = \operatorname{rank}(A) + \dim(\operatorname{null}(A))$$

Nullspace parameterization of solutions

assuming A is wide and full-rank, set of all solutions has the form

$$\{x \mid y = Ax\} = \{x_{p} + z \in \text{null}(A) \text{ and } y = Ax_{p}\}$$

where $x_{\rm p}$ is any "particular" solution

- the vextor z characterizes the choices available
- the dimension of the nullspace gives the number of degrees of freedom in the solution
- choose z to optimize over solutions

Least-norm solution

a particular solution of interest is

$$x_{\ln} = A^T (AA^T)^{-1} y$$

which can be verified by direct substitution: $Ax_{\ln} = y$

the solution x_{\ln} solves the constrained optimization problem

minimize ||x||subject to Ax = y

 $(x_{\ln}$ is the unique solution to this optimization problem)

Optimal least-norm solution

we've directly verified that x_{\ln} is feasible, now show it's optimal

suppose y = Ax, then we must have $A(x - x_{\ln}) = 0$

consider the inner product

$$(x - x_{\ln})^T x_{\ln} = (x - x_{\ln})^T A^T (AA^T)^{-1} y$$

= $(A(x - x_{\ln}))^T (AA^T)^{-1} y$
= 0

the result above shows $(x-x_{
m ln})\perp x_{
m ln}$, so

$$||x||^2 = ||x_{\ln} + x - x_{\ln}||^2 = ||x_{\ln}||^2 + ||x - x_{\ln}||^2 \ge ||x_{\ln}||^2$$

which means x_{ln} has the smallest norm amongst all solutions

Geometry

- x_{\ln} is the orthogonal projection of 0 onto the set $\{x \mid Ax = y\}$
- $x_{\ln} \perp \operatorname{null}(A)$



• $A^{\dagger} - A^T (AA^T)^{-1}$ is a right inverse of (full rank, wide) A: $AA^{\dagger} = I$

underdetermined systems

Solution via QR factorization

apply a QR decomposition to $A^T,$ so that $A^T=QR$

assumed that m < n and rank(A) = m, then

- $R \in \mathbb{R}^{m \times m}$ is invertible
- $Q \in \mathbb{R}^{n \times m}$ with $Q^T Q = I_m$

substituting the factorization into $x_{\mathrm{ln}} = A^T (AA^T)^{-1} y$ we see that

$$x_{\ln} = QR^{-T}y$$
 and $\|x_{\ln}\| = \|R^{-T}y\|$

note: Z^{-T} is shorthand for $(Z^T)^{-1}$

Example: optimal mass transfer

unit mass subject to force x_i at time i, initially at rest



- y_t and v_t: position and velocity at time t
- x_t : constant force applied for time interval [t, t+1]

dynamics given by

 $v_{t+1} - v_t = x_t$ (force equals Δv) $y_{t+1} - y_t = v_t$ (velocity equals Δy)

objective: choose force x to move mass 1 unit of distance in 10s, leaving it at rest [Example from Boyd & Lall]

underdetermined systems

- initial conditions: $y_0 = 0$, $v_0 = 0$
- target: $y_{10} = 1$, $v_{10} = 0$
- expand the system dynamics and express target in terms of initial conditions and the x_i s:

$$\begin{bmatrix} v_{10} \\ y_{10} \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 & 1 \\ 9 & 8 & 7 & \dots & 1 & 0 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_9 \end{bmatrix} + \begin{bmatrix} v_0 \\ y_0 \end{bmatrix}$$

solve for least-norm solution

simulated solution

