

10. Approximation and data fitting

learning model given data

- norm approximation
- least-norm problems
- • regularized approximation
- • robust approximation

Norm approximation

$$\underset{x}{\text{minimize}} \quad \|Ax - \underline{b}\|$$

($A \in \mathbf{R}^{m \times n}$ with $m \geq n$, $\|\cdot\|$ is a norm on \mathbf{R}^m)

interpretations of solution $x^* = \operatorname{argmin}_x \|Ax - \underline{b}\|$:

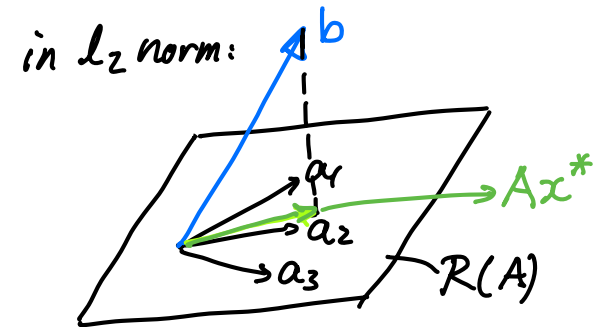
- **geometric:** Ax^* is point in $\mathcal{R}(A)$ closest to \underline{b}
- **estimation:** linear measurement model

$$y = Ax + v$$

y are measurements, x is unknown, v is measurement error

given $y = \underline{b}$, best guess of x is x^*

- **optimal design:** x are design variables (input), Ax is result (output)
 x^* is design that best approximates desired result \underline{b}



examples

- least-squares approximation ($\|\cdot\|_2$): solution satisfies

$$A^T A x = A^T b$$

$$(x^* = (A^T A)^{-1} A^T b \text{ if } \mathbf{rank} A = n)$$

- Chebyshev approximation ($\|\cdot\|_\infty$): can be solved as an LP

$$\|Ax - b\|_\infty \leq t$$

$$|a_i^T x - b_i| \leq t$$

$$\begin{cases} \underset{t, x}{\text{minimize}} & t \\ \text{subject to} & -t\mathbf{1} \preceq Ax - b \preceq t\mathbf{1} \end{cases}$$

- sum of absolute residuals approximation ($\|\cdot\|_1$): can be solved as an LP

$$\begin{cases} \underset{y, x}{\text{minimize}} & \mathbf{1}^T y \\ \text{subject to} & -y \preceq Ax - b \preceq y \end{cases}$$

$$\|Ax - b\|_1 = \sum_{i=1}^n \underbrace{|a_i^T x - b_i|}$$

$$\begin{cases} \min & \sum y_i \\ \text{s.t.} & |a_i^T x - b_i| \leq y_i, \forall i \end{cases}$$

Penalty function approximation

$$\begin{aligned} & \text{minimize} \quad \phi(\underline{r}_1) + \cdots + \phi(\underline{r}_m) \\ & \text{subject to} \quad \underline{r} = Ax - b \\ & x \in \mathbb{R}^n, r \in \mathbb{R}^m \\ & (A \in \mathbb{R}^{m \times n}, \phi : \mathbb{R} \rightarrow \mathbb{R} \text{ is a } \underline{\text{convex}} \text{ penalty function}) \end{aligned}$$

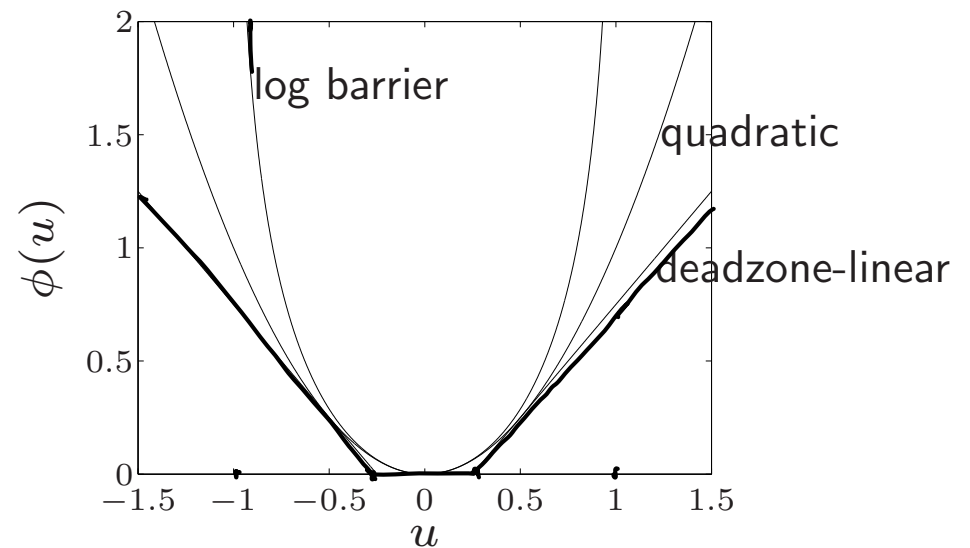
examples

- quadratic: $\phi(u) = u^2$
- deadzone-linear with width a :

$$\phi(u) = \max\{0, |u| - a\}$$

- log-barrier with limit a :

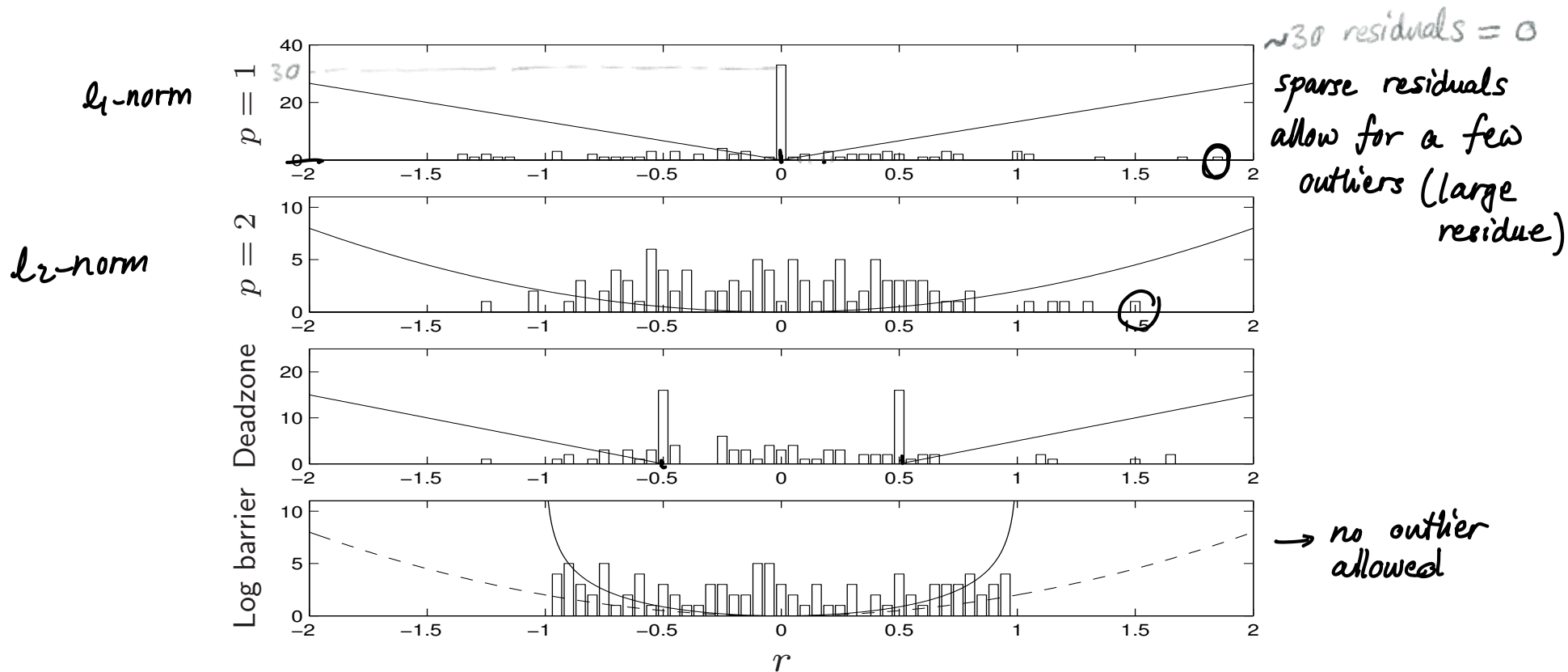
$$\phi(u) = \begin{cases} -a^2 \log(1 - (u/a)^2) & |u| < a \\ \infty & \text{otherwise} \end{cases} \quad a=1$$



$$r \in \mathbb{R}^{100}, \quad x \in \mathbb{R}^{30} \quad A \in \mathbb{R}^{100 \times 30}$$

example ($m = 100, n = 30$): histogram of residuals for penalties

$$\phi(u) = |u|, \quad \phi(u) = u^2, \quad \phi(u) = \max\{0, |u| - a\}, \quad \phi(u) = -\log(1 - u^2)$$



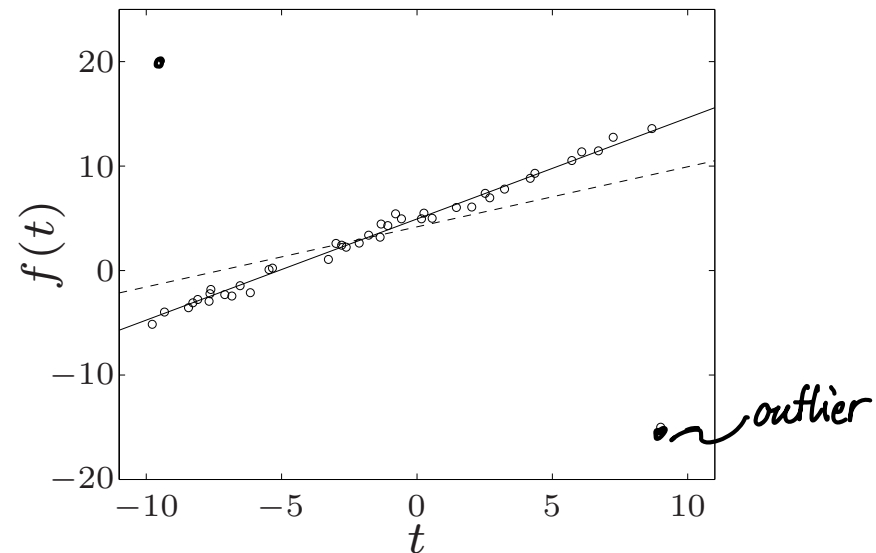
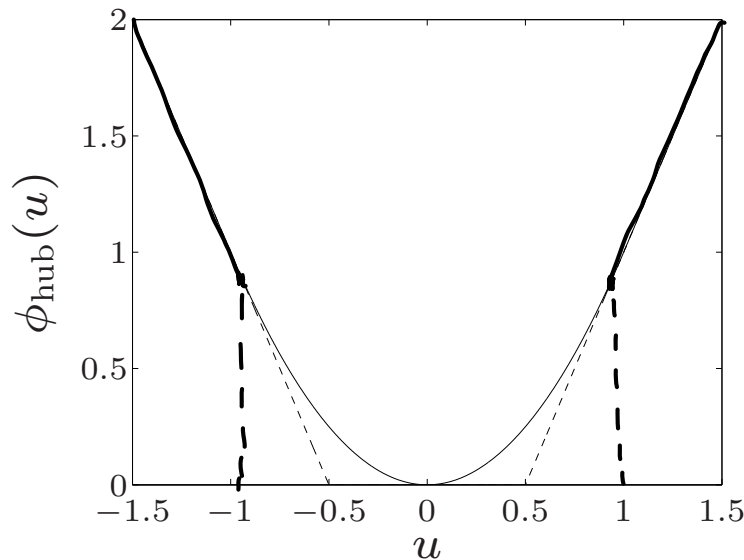
shape of penalty function has large effect on distribution of residuals

Huber penalty function (with parameter M)

*combines l_1 & l_2 norms:
behaves as l_1 norm for larger
residuals (away from zero)*

$$\phi_{\text{hub}}(u) = \begin{cases} u^2 & |u| \leq M \\ M(2|u| - M) & |u| > M \end{cases}$$

linear growth for large u makes approximation less sensitive to outliers



- left: Huber penalty for $M = 1$
- right: affine function $f(t) = \alpha + \beta t$ fitted to 42 points t_i, y_i (circles) using quadratic (dashed) and Huber (solid) penalty

$$\phi(u) = u^2$$

- Robust statistics (robust to "outliers")

Announcements

- HW8 (last hw) due Wed by midnight, solutions will be posted 24 hrs after (or sooner if all HWs are submitted)
- **Final exam**: 10-hours take-home, start @ 9am
either 3/9 or 3/10, see Canvas announcement.
- Final exams + solutions from previous years posted (Files/exams/)
BUT this year's final is a different style (and includes some CVX/PY problems as well)
- Course evaluations are open!
Please fill: depts use these in instructor/TA evaluations & allocating resources.
I always read all feedback & much appreciate it — thanks!
- Office hours: as usual today (**on Zoom**), tomorrow, and Fri instead of TA session.

Least-norm problems

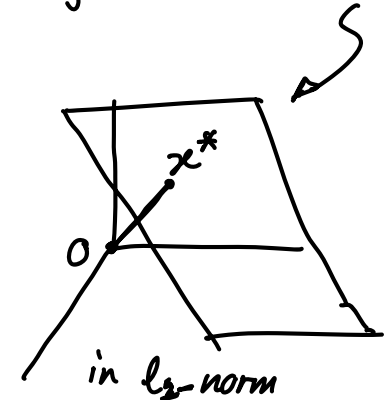
wide matrix

$$\left[\begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \end{array} \right]_{m \times n}$$

($A \in \mathbf{R}^{m \times n}$ with $m \leq n$, $\|\cdot\|$ is a norm on \mathbf{R}^n)

$$\begin{cases} \text{minimize} & \|x\| \\ \text{subject to} & Ax = \underline{b} \end{cases}$$

$\{x \mid Ax = b\}$ is an affine set



interpretations of solution $x^* = \operatorname{argmin}_{Ax=b} \|x\|$:

- **geometric:** $\underline{x^*}$ is point in affine set $\{x \mid Ax = b\}$ with minimum distance to 0
- **estimation:** $b = Ax$ are (perfect) measurements of x ; $\underline{x^*}$ is smallest ('most plausible') estimate consistent with measurements
- **design:** \underline{x} are design variables (inputs); $\overset{\text{desired}}{\underline{b}}$ are required results (outputs)
 x^* is smallest ('most efficient') design that satisfies requirements

$$x(t) \sum_t \|x(t)\|^2 \sim \text{signal energy}$$

$$A \in \mathbb{R}^{m \times n} \quad m < n$$

examples

- least-squares solution of linear equations ($\|\cdot\|_2$):
can be solved via optimality conditions

$$\begin{cases} \min_x & \|x\|_2^2 \\ \text{s.t.} & Ax = b \end{cases} \leftarrow \nu$$

$$L(x, \nu) = \|x\|_2^2 + \nu^T (Ax - b)$$

$$\underline{2x + A^T \nu = 0}, \quad \underline{Ax = b} \quad \nabla_x L(x, \nu) \big|_{x=x^*} = 0$$

- minimum sum of absolute values ($\|\cdot\|_1$): can be solved as an LP

$$\begin{cases} \min_{x, y} & \mathbf{1}^T y \\ \text{subject to} & -y \preceq x \preceq y, \quad Ax = b \end{cases}$$

$$x^* = A^+ b$$

$$\begin{cases} \min_x & \sum |x_i| \\ \text{s.t.} & Ax = b \\ & |x_i| \leq y_i \end{cases}$$

tends to produce sparse solution x^*

extension: least-penalty problem

$$\begin{cases} \min & \phi(x_1) + \dots + \phi(x_n) \\ \text{subject to} & \underline{Ax = b} \end{cases}$$

$\phi : \mathbf{R} \rightarrow \mathbf{R}$ is convex penalty function

training a NNNet:

overparameterized \Rightarrow many perfect fits to data

but additional properties of learned x are desired

\rightarrow design appropriate loss functions

(generalization)

Regularized approximation

$$\left[\underset{x}{\text{minimize (w.r.t. } \mathbf{R}_+^2)} \quad \underbrace{(\|Ax - b\|)}_{f_1(x)}, \underbrace{\|x\|}_{f_2(x)} \right]$$

fitting error
"size"

$A \in \mathbf{R}^{m \times n}$, norms on \mathbf{R}^m and \mathbf{R}^n can be different

interpretation: find good approximation $Ax \approx b$ with small x

- **estimation:** linear measurement model $y = Ax + \underline{v}$, with prior knowledge that $\|x\|$ is small $y = Ax$
- **optimal design:** small x is cheaper or more efficient, or the linear model $y = Ax$ is only valid for small x
- **robust approximation:** good approximation $Ax \approx b$ with small x is less sensitive to errors in A than good approximation with large x

Ax
errors in A get multiplied by x_i 's

Scalarized problem

$$\left[\underset{x}{\text{minimize}} \quad \|Ax - b\| + \underline{\gamma} \|x\| \right]$$

- solution for $\gamma > 0$ traces out optimal trade-off curve
- other common method: minimize $\|Ax - b\|^2 + \delta \|x\|^2$ with $\delta > 0$

Tikhonov regularization : ℓ_2 -norm (ridge-regression)

$$\left[\underset{x}{\text{minimize}} \quad \|Ax - b\|_2^2 + \delta \|x\|_2^2 \right]$$

can be solved as a least-squares problem

$$\left[\underset{x}{\text{minimize}} \quad \left\| \underbrace{\begin{bmatrix} A \\ \sqrt{\delta} I \end{bmatrix}}_{\tilde{A}} x - \underbrace{\begin{bmatrix} b \\ 0 \end{bmatrix}}_{\tilde{b}} \right\|_2^2 \right]$$

$$x^* = \tilde{A}^\dagger \tilde{b}$$

solution $x^* = (A^T A + \delta I)^{-1} A^T b$

Signal reconstruction

minimize (w.r.t. \mathbf{R}_+^2) $(\|\hat{x} - x_{\text{cor}}\|_2, \phi(\hat{x}))$



- $x \in \mathbf{R}^n$ is unknown signal
- $x_{\text{cor}} = x + \underline{v}$ is (known) corrupted version of x , with additive noise v
- variable \hat{x} (reconstructed signal) is estimate of x
- $\phi : \mathbf{R}^n \rightarrow \mathbf{R}$ is regularization function or smoothing objective

examples: quadratic smoothing, total variation smoothing:

$$\hat{x} = \begin{bmatrix} \hat{x}_1 \\ \vdots \\ \hat{x}_n \end{bmatrix}$$

$$\underline{\phi_{\text{quad}}(\hat{x})} = \sum_{i=1}^{n-1} \underline{(\hat{x}_{i+1} - \hat{x}_i)^2},$$

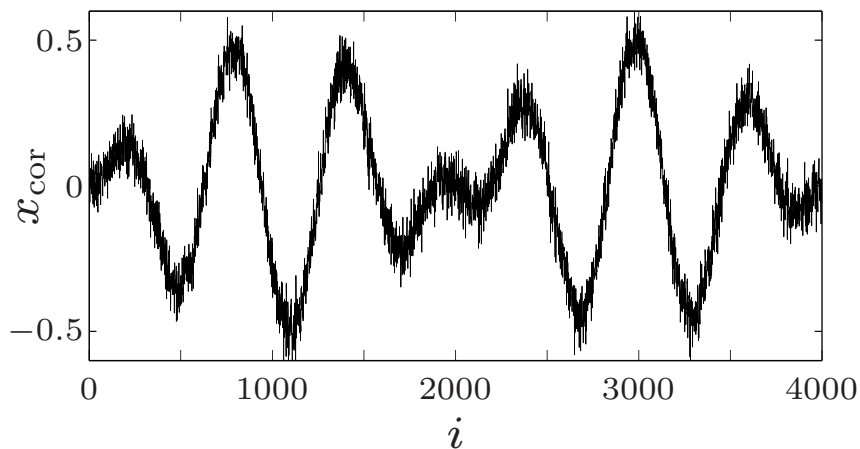
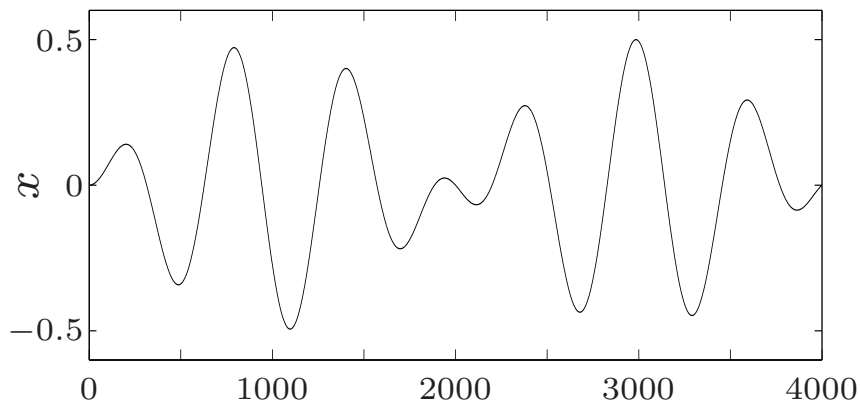
$$\underline{\phi_{\text{tv}}(\hat{x})} = \sum_{i=1}^{n-1} |\hat{x}_{i+1} - \hat{x}_i|$$

$$= \|d\|_1$$

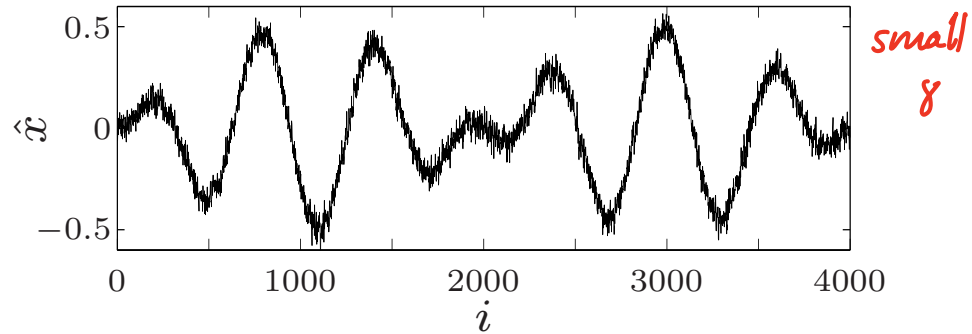
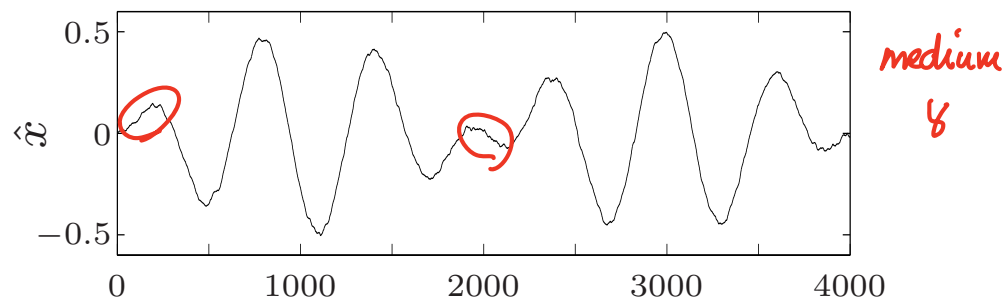
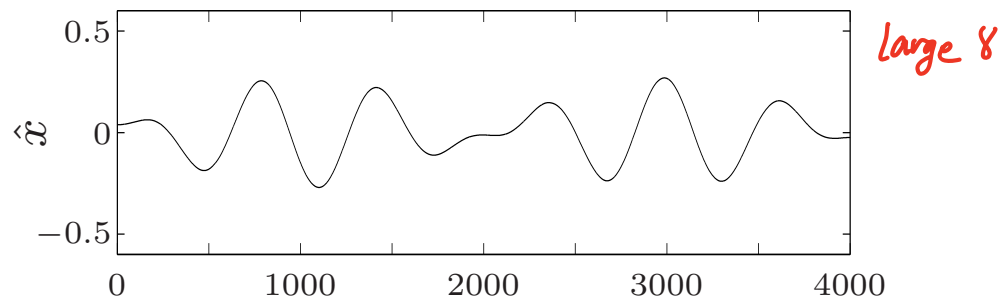
$$d = \begin{bmatrix} \hat{x}_2 - \hat{x}_1 \\ \hat{x}_3 - \hat{x}_2 \\ \vdots \\ \hat{x}_n - \hat{x}_{n-1} \end{bmatrix} \in \mathbf{R}^{n-1}$$

$$\phi_{\text{quad}}(\hat{x}) = \|d\|_2^2$$

quadratic smoothing example $\hat{x} = \underset{x}{\operatorname{argmin}} \underbrace{\|x - x_{\text{cor}}\|_2^2}_{\text{data fidelity}} + \underbrace{\gamma \phi_{\text{quad}}(x)}_{\text{smoothness}}$



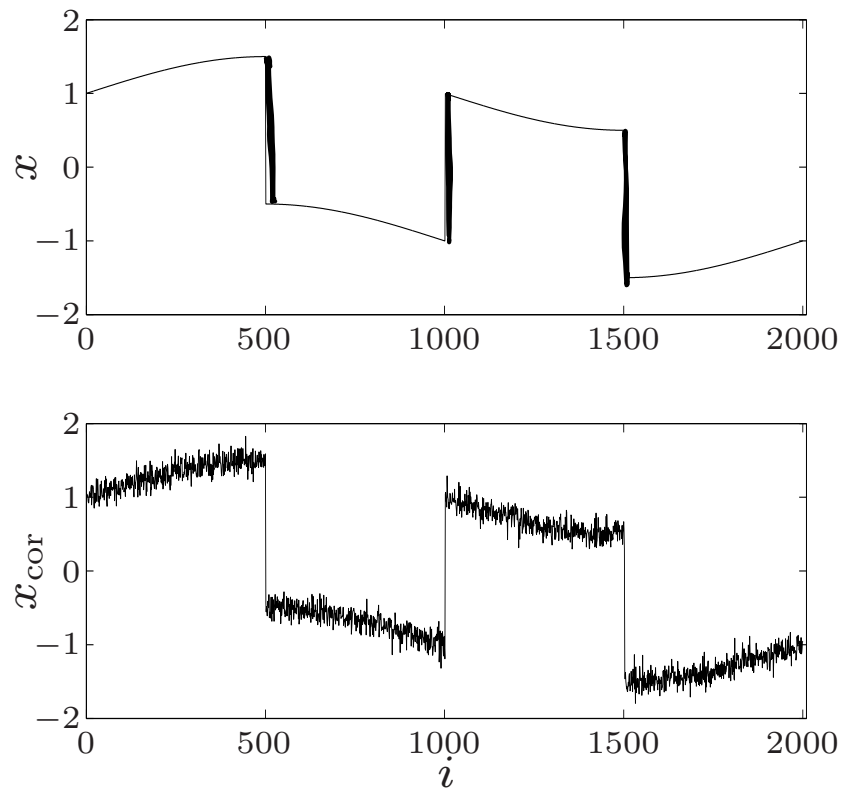
original signal x and noisy
signal x_{cor}



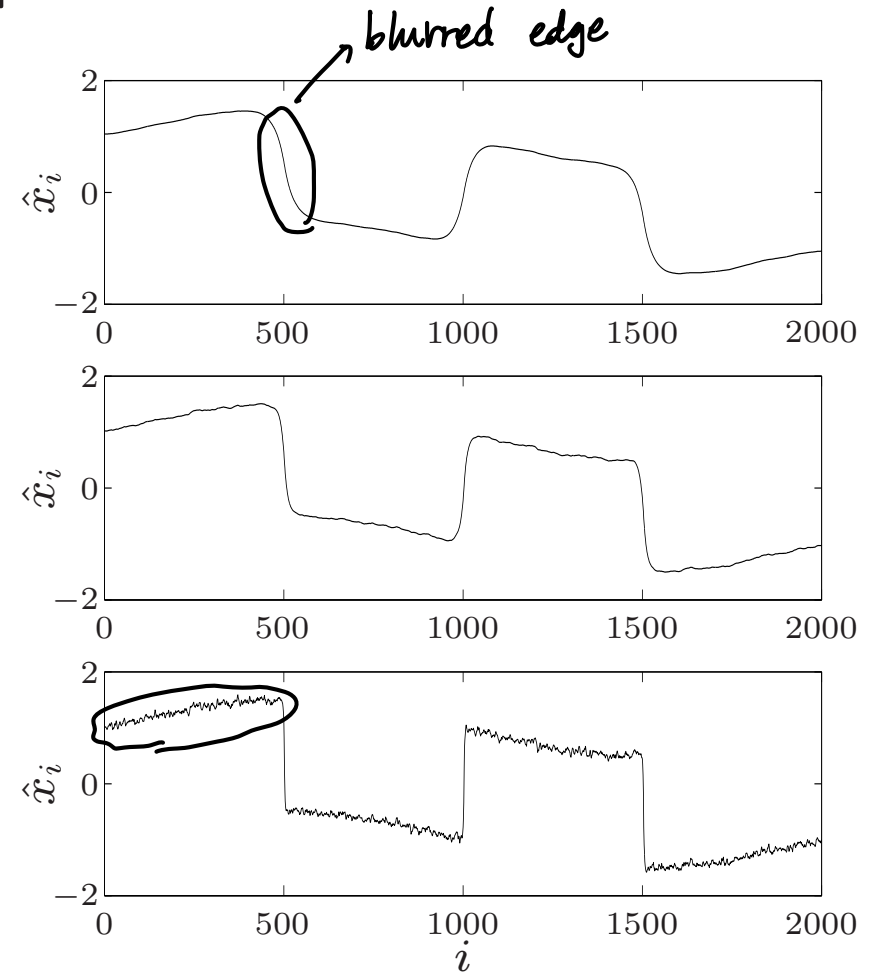
three solutions on trade-off curve
 $\|\hat{x} - x_{\text{cor}}\|_2$ versus $\phi_{\text{quad}}(\hat{x})$

total variation reconstruction example

$$\|x - x_{\text{cor}}\|_2^2 + \gamma \phi_{\text{tv}}(x)$$

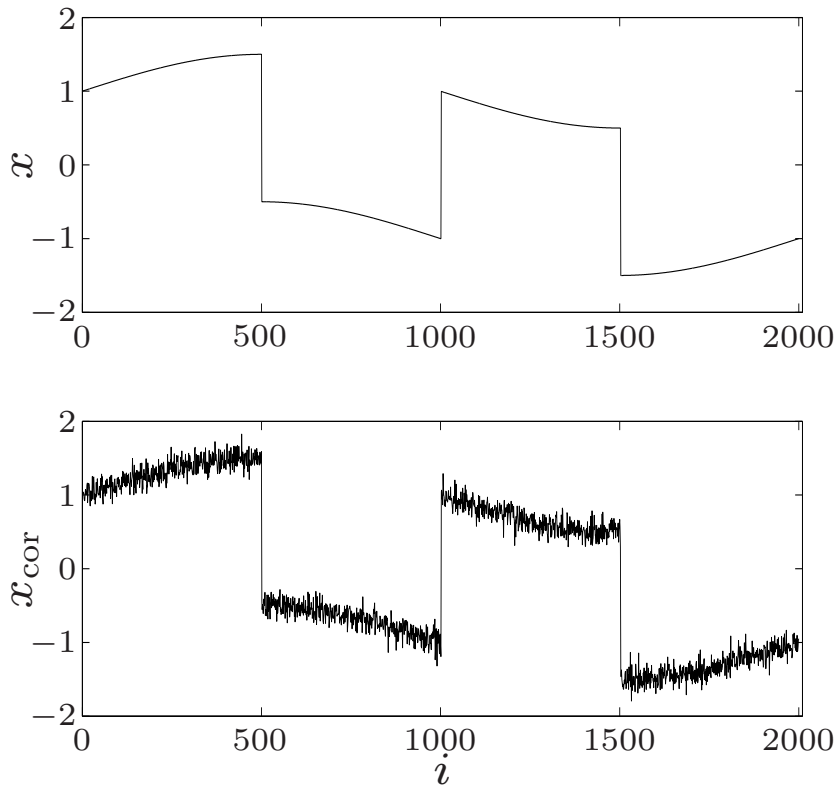


original signal x and noisy
signal x_{cor}

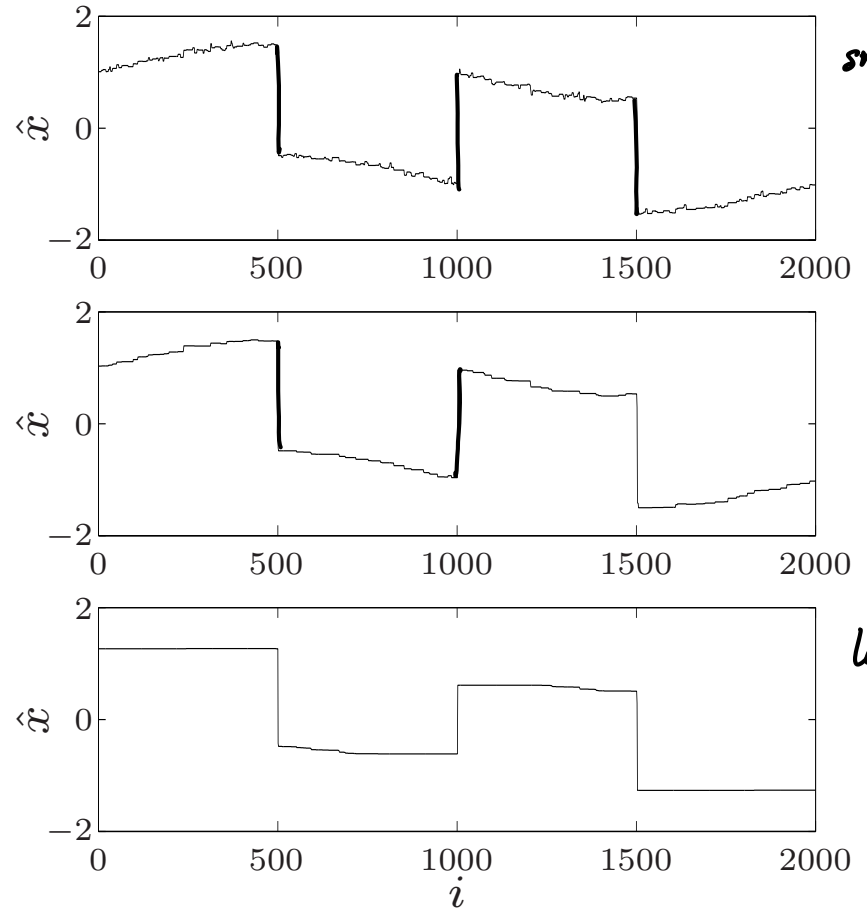


three solutions on trade-off curve
 $\|\hat{x} - x_{\text{cor}}\|_2$ versus $\phi_{\text{quad}}(\hat{x})$

quadratic smoothing smooths out noise **and** sharp transitions in signal



original signal x and noisy
signal x_{cor}



three solutions on trade-off curve

$$\phi_{\text{tv}}(\hat{x}) = \sum_{i=1}^{n-1} |\hat{x}_{i+1} - \hat{x}_i|$$

total variation smoothing preserves sharp transitions in signal

Robust approximation

minimize $\|Ax - b\|$ with uncertain A

two approaches:

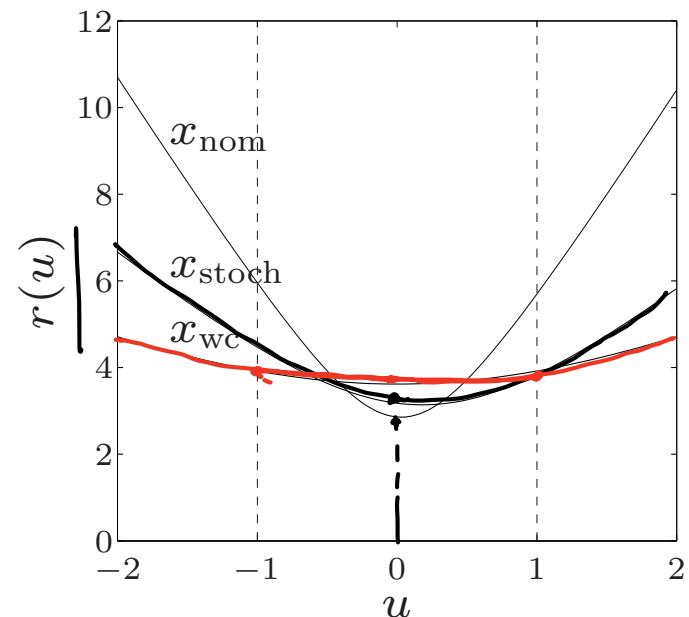
- stochastic: assume A is random, minimize $\mathbf{E}_A \|Ax - b\|$
- worst-case: set \mathcal{A} of possible values of A , minimize $\sup_{A \in \mathcal{A}} \|Ax - b\|$

tractable only in special cases (certain norms $\|\cdot\|$, distributions, sets \mathcal{A})

example: $A(u) = A_0 + uA_1$

- • x_{nom} minimizes $\|A_0x - b\|_2^2$
- • x_{stoch} minimizes $\mathbf{E} \|A(u)x - b\|_2^2$
with u uniform on $[-1, 1]$
- • x_{wc} minimizes $\sup_{-1 \leq u \leq 1} \|A(u)x - b\|_2^2$

figure shows $r(u) = \|A(u)x - b\|_2$

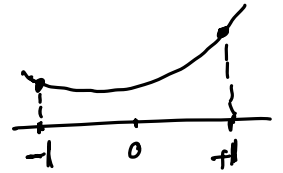


- $$\begin{aligned}
 \mathbb{E} \|A(u)x - b\|_2^2 &= \mathbb{E}_u \| (A_0 + uA_1)x - b \|_2^2 \\
 &= \mathbb{E}_u \| A_0x - b + uA_1x \|_2^2 \\
 &= \mathbb{E}_u (\|A_0x - b\|_2^2 + \|uA_1x\|_2^2) \\
 &= \|A_0x - b\|_2^2 + \|A_1x\|_2^2 \underbrace{\mathbb{E}_u u^2}_{\text{variance of } U[-1,1]}
 \end{aligned}$$

$\rightarrow \left[\min_x \|A_0x - b\|_2^2 + (\mathbb{E}_u u^2) \|A_1x\|_2^2 \right] \rightsquigarrow \text{least squares}$

- convex quadratic fit of $u \leftarrow \| (A_0 + uA_1)x - b \|_2^2$

$$\begin{aligned}
 \min_x \sup_{-1 \leq u \leq 1} \| (A_0 + uA_1)x - b \|_2^2 \\
 = \max \left\{ \underbrace{\| (A_0 + A_1)x - b \|_2^2}_{\text{convex in } x}, \underbrace{\| (A_0 - A_1)x - b \|_2^2}_{\text{convex in } x} \right\}
 \end{aligned}$$



stochastic robust LS with $A = \bar{A} + \underline{U}$, U random, $\mathbf{E} U = 0$, $\mathbf{E} U^T U = \underline{P}$

$$\text{minimize } \mathbf{E} \|(\bar{A} + U)x - b\|_2^2$$

- explicit expression for objective:

$$\begin{aligned} \mathbf{E} \|Ax - b\|_2^2 &= \mathbf{E} \|\bar{A}x - b + Ux\|_2^2 \\ &= \|\bar{A}x - b\|_2^2 + \mathbf{E}_{\mathbf{U}} x^T \underline{U^T U} x \\ &= \|\bar{A}x - b\|_2^2 + \underline{x^T P x} \end{aligned}$$

- hence, robust LS problem is equivalent to LS problem

$$\underset{x}{\text{minimize}} \quad \|\bar{A}x - b\|_2^2 + \|P^{1/2}x\|_2^2$$

- for $P = \delta I$, get Tikhonov regularized problem

$$\left[\text{minimize} \quad \|\bar{A}x - b\|_2^2 + \delta \|x\|_2^2 \right.$$

skipped

worst-case robust LS with $\mathcal{A} = \{\bar{A} + u_1 A_1 + \cdots + u_p A_p \mid \|u\|_2 \leq 1\}$

$$\text{minimize} \quad \sup_{A \in \mathcal{A}} \|Ax - b\|_2^2 = \sup_{\|u\|_2 \leq 1} \|P(x)u + q(x)\|_2^2$$

where $P(x) = \begin{bmatrix} A_1 x & A_2 x & \cdots & A_p x \end{bmatrix}$, $q(x) = \bar{A}x - b$

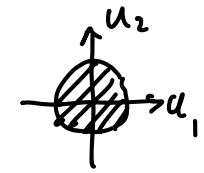
- from page 7-32, strong duality holds between the following problems

$$\begin{array}{ll} \text{maximize} & \|Pu + q\|_2^2 \\ \text{subject to} & \|u\|_2^2 \leq 1 \end{array} \qquad \begin{array}{ll} \text{minimize} & t + \lambda \\ \text{subject to} & \begin{bmatrix} I & P & q \\ P^T & \lambda I & 0 \\ q^T & 0 & t \end{bmatrix} \succeq 0 \end{array}$$

- hence, robust LS problem is equivalent to SDP

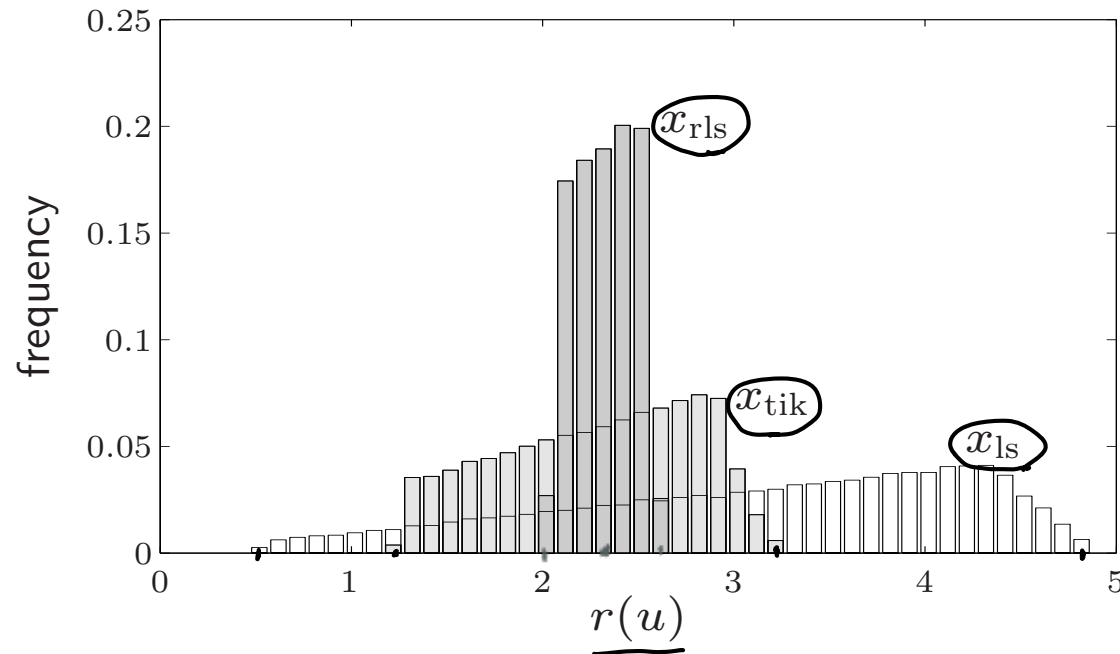
$$\begin{array}{ll} \text{minimize} & t + \lambda \\ \text{subject to} & \begin{bmatrix} I & P(x) & q(x) \\ P(x)^T & \lambda I & 0 \\ q(x)^T & 0 & t \end{bmatrix} \succeq 0 \end{array}$$

example: histogram of residuals



$$r(u) = \|(\underline{A_0} + \underline{u_1}A_1 + \underline{u_2}A_2)x - b\|_2$$

with u uniformly distributed on unit disk, for three values of x



- • x_{ls} minimizes $\|A_0x - b\|_2$
- • x_{tik} minimizes $\|A_0x - b\|_2^2 + \|x\|_2^2$ (Tikhonov solution)
- • x_{wc} minimizes $\sup_{\|u\|_2 \leq 1} \|A_0x - b\|_2^2 + \|x\|_2^2$
 $\quad \quad \quad = x_{rls}$