

# ENGSCI 721

## INVERSE PROBLEMS

*Oliver Maclarens*  
*oliver.maclarens@auckland.ac.nz*

## MODULE OVERVIEW

Inverse Problems (*Oliver Maclarens*) [*~8 lectures/2-3 tutorials*]

### 1. Basic concepts [3 lectures]

Forward vs inverse problems. Well-posed vs ill-posed problems. Algebra of inverse problems (generalised inverses etc). Regularisation and trade-offs.

### 2. More regularisation [3 lectures]

Higher-order Tikhonov regularisation, truncated singular value decompositions, iterative regularisation.

## MODULE OVERVIEW

### 3. Statistical view of inverse problems I [2 lectures]

Bayesians, Frequentists and all that. Basic frequentist analysis. Linearisation and covariance propagation.

## LECTURE 4: REGULARISATION - TIKHONOV AND BEYOND

Topics:

- Recap of Tikhonov
- Higher-order Tikhonov
- $L_1$  norm regularisation
  - Sparsity
  - Total variation
  - Robust regression

## Eng Sci 741 : Lecture 4.

### Regularisation - Tikhonov & beyond

- Recap: basic Tikhonov
- Higher-order Tikhonov
- $L_1$  norm  
Sparsity
- Total variation
- Robust regression

### Recap: Basic Tikhonov

The standard Tikhonov approach to constructing regularised solutions to inverse problems like

$$\boxed{\text{find } x \text{ satisfying } Ax = y}$$

where  $A$  tends to 'smooth' or 'reduce'

$x \mapsto Ax$  , is to consider

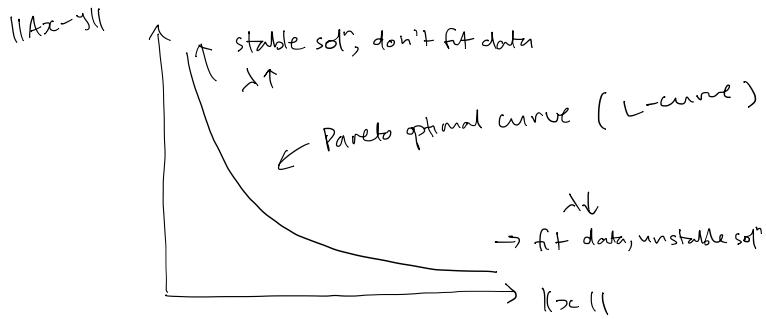
the modified problem:

$$\boxed{\min_{x(\lambda)} \|Ax - y\|^2 + \lambda \|x\|^2}$$

Note that the solution depends on

the regularisation parameter  $\lambda$ , hence we write  $x(\lambda)$

The parameter  $\lambda$  represents a trade-off between the importance of fitting the given data  $y$  & having a 'simple' or 'stable' solution



→ since the data  $y$  is not exact / exactly repeatable, we use  $\lambda$  to satisfy Hadamard's third condition of stability:

'small changes to given info should give small changes to solution'

→ choosing  $\lambda$  is somewhat of an art & part of designing a regularisation procedure

↳ eg { L-curve corner ← we looked at.  
discrepancy principle  
cross-validation  
etc.

When we use the squared  $L_2$  norm for both, it is convenient to write this as the augmented least squares problem:

$$\boxed{\min \|\tilde{r}\|^2 = \min \|\tilde{A}x - \tilde{y}\|^2}$$

where  $\boxed{\begin{aligned} \tilde{A}x &= \tilde{y} \\ \left[ \begin{matrix} A \\ \sqrt{\lambda} I \end{matrix} \right] x &- \left[ \begin{matrix} y \\ 0 \end{matrix} \right] \end{aligned}}$

→ We can then use any standard least-squares software to solve.

↳ Matlab: lsqminnorm

• Python: np.lstsq (does min norm if over-det.).

Can also use pinv since this } slightly gives least squares sol'n. } less eff alg. I think

(Note: Matlab ↴ gives different pseudoinverse!)

## Linear vs nonlinear

The linear case has the explicit solution

$$x = (A^T A + \lambda I)^{-1} A^T y$$

i.e.

$$x = A_\lambda^* y \text{ where } A_\lambda^* = (A^T A + \lambda I)^{-1} A^T$$

In the nonlinear case  $F(x) = y$  we can still solve the 'variational form':

$$\min_{x(\lambda)} \|F(x) - y\|^2 + \lambda \|x\|^2$$

$$\text{or } \min \| \tilde{x} \|^2 = \min \| \tilde{F}(x) - \tilde{y} \|^2$$

where  $\tilde{F}(x) = \tilde{y}$   
is  $\begin{bmatrix} F(x) \\ \sqrt{\lambda} x \end{bmatrix} = \begin{bmatrix} y \\ 0 \end{bmatrix}$

→ no explicit inverse, but...

→ solve via nonlinear least-squares,  
or any minimisation alg (see later)

## Generalisations

The variational form is easiest to generalise

→ nonlinear problems ✓

→ different norms for data/model space ✓

Just gives an objective function to minimise via any minimisation alg.

→ `scipy.optimize` library (Python)

→ `fminsearch` (Matlab; see also opt. toolbox)

- First, however, we look at forms that still fit in least squares framework

- Then forms that fit into other efficient frameworks

↳ convex opt. (Boyd & Vandenberghe)  
<sup>eg</sup>

- Finally (later lectures), generic problems & alternative approaches (iterative reg.)

## Higher-order Tikhonov

The standard model norm that we've seen measures 'size' via

$$\|x\| \text{ or } \|x\|^2$$

→ We can generalise this by considering

$$\|Dx\| \text{ or } \|Dx\|^2$$

for some operator

(or  $\|Dx - x_0\|$  esp. for nonlinear prob)

→ Typically  $D$  represents a first or second derivative operator, (or even a differential eqn!)

↳ can consider diff. operators as 'roughening' operators (cf. integration as smoothing)

↳ we prefer smoother ie less variable solutions, hence penalise roughness.

## General Tikhonov: (can be linear or nonlinear)

$$\boxed{\min \|Ax - y\|_2^2 + \lambda \|Dx\|_2^2}$$

where,

(for discrete, linear diff. operators in 1D  
→ can generalise to higher  $D$ ):

'zeroth order'  
 $D_0 = I = \begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix}$

'first order'

$$D_1 = \begin{bmatrix} -1 & 1 & & & \\ & -1 & 1 & & \\ & & \ddots & \ddots & \\ & & & -1 & 1 \end{bmatrix}$$

Second order

$$D_2 = \begin{bmatrix} 1 & -2 & 1 & & & \\ & 1 & -2 & 1 & & \\ & & \ddots & \ddots & \ddots & \\ & & & 1 & -2 & 1 \end{bmatrix}$$

Note:

$$D_1 x = \begin{bmatrix} -1 & & & \\ & \ddots & & \\ & & \ddots & \\ & & & -1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$= \begin{bmatrix} x_2 - x_1 \\ x_3 - x_2 \\ \vdots \\ x_n - x_{n-1} \end{bmatrix} \quad \text{ie first order forward finite differences}$$

$$D_2 x = \begin{bmatrix} 1 & -2 & 1 & & \\ & \ddots & \ddots & \ddots & \\ & & & & \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$= \begin{bmatrix} x_1 - 2x_2 + x_3 \\ \vdots \\ x_{n-2} - 2x_{n-1} + x_n \end{bmatrix} \quad \text{ie second-order central finite differences}$$

Also:

$$D_1^2 = D_1 \cdot D_1 = \begin{bmatrix} -1 & & & \\ & -1 & & \\ & & -1 & \\ & & & 1 \end{bmatrix} \begin{bmatrix} -1 & & & \\ & -1 & & \\ & & -1 & \\ & & & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & -2 & 1 \\ & 1 & -2 \\ & & 1 \end{bmatrix} \\ = D_2$$

Or 'wrapping' (ie 'circular differences')

$$D_{1c}^2 = D_{1c} \cdot D_{1c} = \begin{bmatrix} -1 & & & \\ & -1 & & \\ & & -1 & \\ & & & 1 \end{bmatrix} \begin{bmatrix} -1 & & & \\ & -1 & & \\ & & -1 & \\ & & & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & -2 & 1 \\ 1 & 1 & -2 \\ -2 & 1 & 1 \end{bmatrix} \\ = D_{2c} \\ \uparrow \\ \text{'circular'}$$

## Least squares form

Just like before, when working with the  $L_2$ -norm we can re-write

$$\min_{x(\lambda)} \|Ax - y\|_2^2 + \lambda \|Dx\|_2^2$$

as

$$\min_{x(\lambda)} \left\| \begin{bmatrix} A \\ \sqrt{\lambda} D \end{bmatrix} x - \begin{bmatrix} y \\ 0 \end{bmatrix} \right\|_2^2$$

$$\text{ie } \min_{x(\lambda)} \|\tilde{A}x - \tilde{y}\|_2^2$$

$$\text{where } \tilde{A} = \begin{bmatrix} A \\ \sqrt{\lambda} D \end{bmatrix}$$

& use standard least-squares software  
(linear or nonlinear)

## Explicit solution (linear only)

In the linear case we have the explicit solution

$$x = (A^T A + \lambda D^T D)^{-1} A^T y$$

i.e.

$$x = A_{D,\lambda}^* y \quad (\text{or just } x = A^* y \text{ if lazy})$$

$$\text{where } A_{D,\lambda}^* = (A^T A + \lambda D^T D)^{-1} A^T$$

Note that we technically have to assume that  $A^T A + \lambda D^T D$  is invertible

→ obvious for  $D = I$  (see before)

→ true for  $D = D_1$  or  $D_2$  too

condition:  $N(A) \cap N(D) = \{0\}$

i.e. null spaces have only trivial sol<sup>n</sup> in common

## Examples

### Simple smoothing

$$x_{\text{noisy}} = x_{\text{smooth}} + \epsilon = Ix_{\text{smooth}} + \epsilon$$

$$= Ax_{\text{smooth}} + \epsilon, \epsilon \text{ unknown}$$

$$\Rightarrow \boxed{A = I} \quad (\text{'deterministic' part; } \epsilon \text{ dealt with by using least squares sol})$$

Note: we don't explicitly model  $\epsilon$ , but relates to choice of 'fit'

$$\text{Here: } \|Ax_{\text{smooth}} - x_{\text{noisy}}\|_2^2$$

→ least squares fit (don't expect exact fit when noisy).

→ but without reg., will fit exactly

Regularised:

$$\boxed{x_{\text{smooth}} = (I + \lambda D^T D)^{-1} I x_{\text{noisy}}}$$

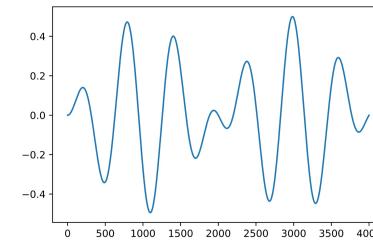
$$\text{Solves: } \underbrace{\|x_{\text{noisy}} - x_{\text{smooth}}\|_2^2}_{\text{fit data}} + \lambda \underbrace{\|Dx_{\text{smooth}}\|_2^2}_{\text{smooth data.}}$$

## Example

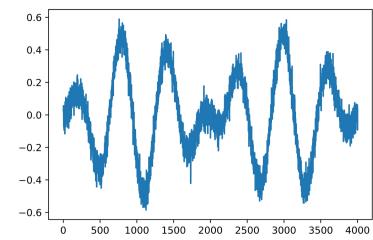
[see Boyd & Vandenberghe 6.3.3 (attached) for original]

My attempt

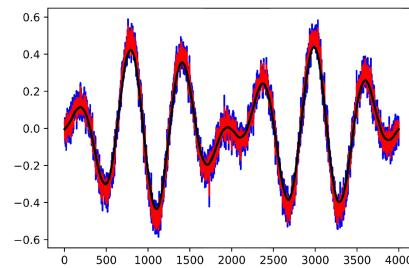
True signal



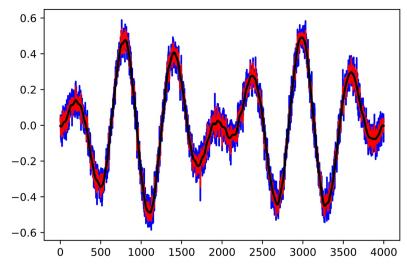
Noisy signal



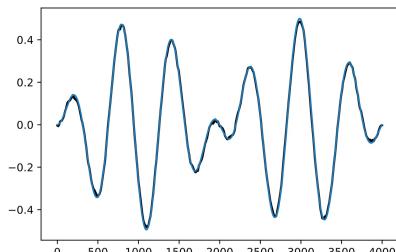
Recovered ( $D_1$  reg.) diff.



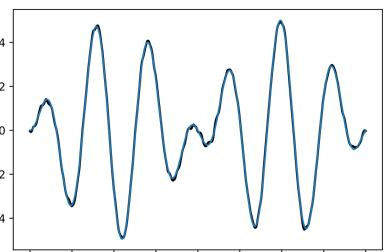
Recovered ( $D_2$  reg.) diff.



'Best  $D_1$ ' & True



'Best  $D_2$ ' & True



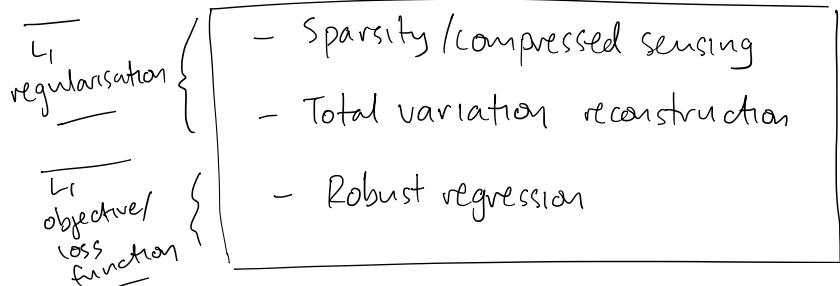
## Other norms etc

'Tikhonov' regularisation generally refers to the  $L_2$  norm for both 'data space' & 'model space'

- Useful computationally (eg re-frame as standard least squares, differentiable...)
- But other norms can be useful to emphasise different data &/or model features

mix &  
match  
data &  
model  
norms

In particular, the  $L_1$  norm is closely related to e.g



(The  $L_0$  &  $L_\infty$  norms also come up fairly often)

## $L_1$ norm & $L_p$ norms

The  $L_1$  norm has the form

$$\|x\|_1 = |x_1| + |x_2| + \dots + |x_n|$$

Both the  $L_1$  &  $L_2$  norms are special cases of  $L_p$  norms:

$$\|x\|_p = \left( |x_1|^p + |x_2|^p + \dots + |x_n|^p \right)^{1/p}$$

for  $p \geq 1$ .

The  $L_\infty$  norm is the limit as  $p \rightarrow \infty$  and given by

$$\|x\|_\infty = \max \{ |x_1|, |x_2|, \dots, |x_n| \}$$

The  $L_0$  'norm' is not technically a norm but means 'number of non-zero' entries in  $x$

$$\text{ie } \|x\|_0 = \text{card}(x)$$

where 'card' means 'cardinality'.

$L_1$  &  $L_0$  'norms': Sparsity.

The  $L_0$  'norm' measures 'sparsity'

→ small  $L_0$  'norm' means } 'sparse'  
few non-zero entries

→ can get 'simple' solutions  
in sense of 'only a few  
nonzero components'

→ Unfortunately is basically just  
enumeration & very difficult  
to work with  $L_0$  directly

A famous & somewhat surprising  
result is that the  $L_1$  norm, when  
used as regularisation also tends to  
produce sparse solutions, i.e. solutions  
with exactly zero elements

Furthermore: (linear)  $L_1$  problems lead to

convex optimisation problems

that can still be solved relatively efficiently  
(eg Linear / Quadratic Programming etc)

Equivalence of norms?

In contrast to  $L_0/L_1$ ,  $L_2$  produces solutions  
with very small but non-zero elements

But aren't all norms in  $\mathbb{R}^n$   
'equivalent'? i.e.

$$\exists \alpha, \beta \text{ s.t. } \alpha \|x\|_a \leq \|x\|_b \leq \beta \|x\|_a$$

for any two norms  $\|\cdot\|_a$  &  $\|\cdot\|_b$

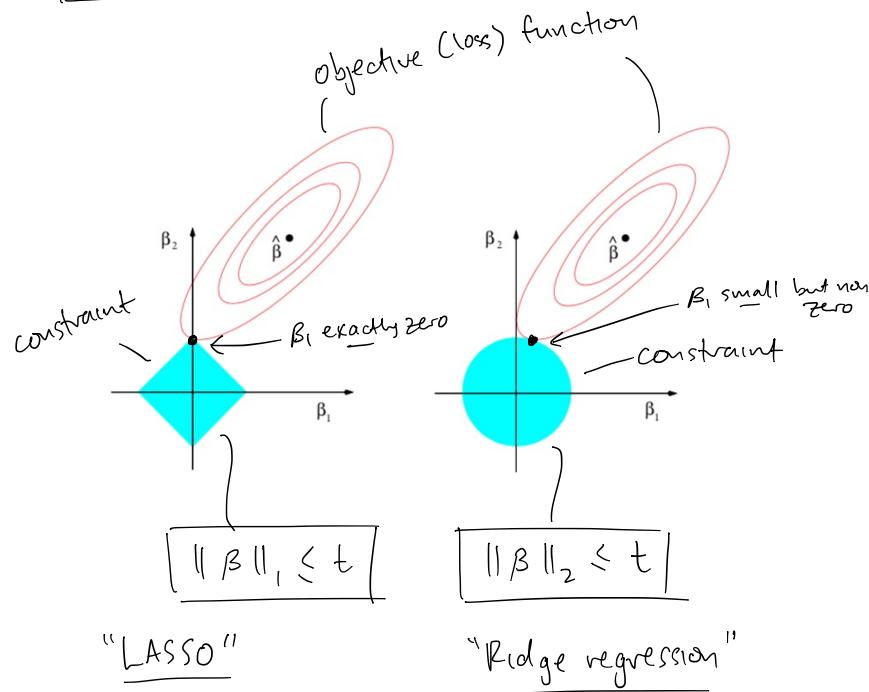
- value of norm can be approximated,  
but the geometry & solutions favoured  
by different norms are different  
↳ also,  $\alpha, \beta$  matter! (as  $n \rightarrow \infty$  approx gets worse)
- quantitatively similar evaluation of given  
solution but problems lead to  
qualitatively solutions as 'best'

$L_2$ : - lots of small but non-zero values  
- related to 'average'

$L_0/L_1$ : - lots of exactly zero values  
- related to 'median'  
- more robust/less sensitive to large  
individuals than  $L_2$ .

## $L_1$ vs $L_2$ regularisation

$$\begin{array}{ll} \min & \|A\beta - y\|_2^2 \\ \text{st.} & \|\beta\|_1 \leq t \end{array} \quad \left[ \begin{array}{ll} \min & \|A\beta - y\|_2^2 \\ \text{st.} & \|\beta\|_2 \leq t \end{array} \right]$$



From 'Statistical learning with sparsity'  
by Hastie et al.

## Applications of $L_1$ : Sparse solutions

('LASSO' regression)

$$\begin{array}{l} \min \|x\|_1 \\ \text{st. } \|Ax - y\|_2 \leq \delta \end{array}$$

same as  
before, but  
emph. model  
norm.

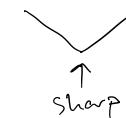
or variational form:

$$\min \|Ax - y\|_2^2 + \lambda \|x\|_1$$

key  
form.

→ Tikhonov-like, but  $L_1$  on  
model/parameters  $x$  instead of  $L_2$

→ Non-differentiable at  $x = 0$



→ But: convex optimisation  
when  $A$  is linear

↳ use e.g. linear/quad. programming,  
iterative reweighted LS etc

↳ many existing libraries for LASSO!

→ can also use differentiable approx. to  $L_1$ ,

→ or just derivative-free optimisation  
etc.

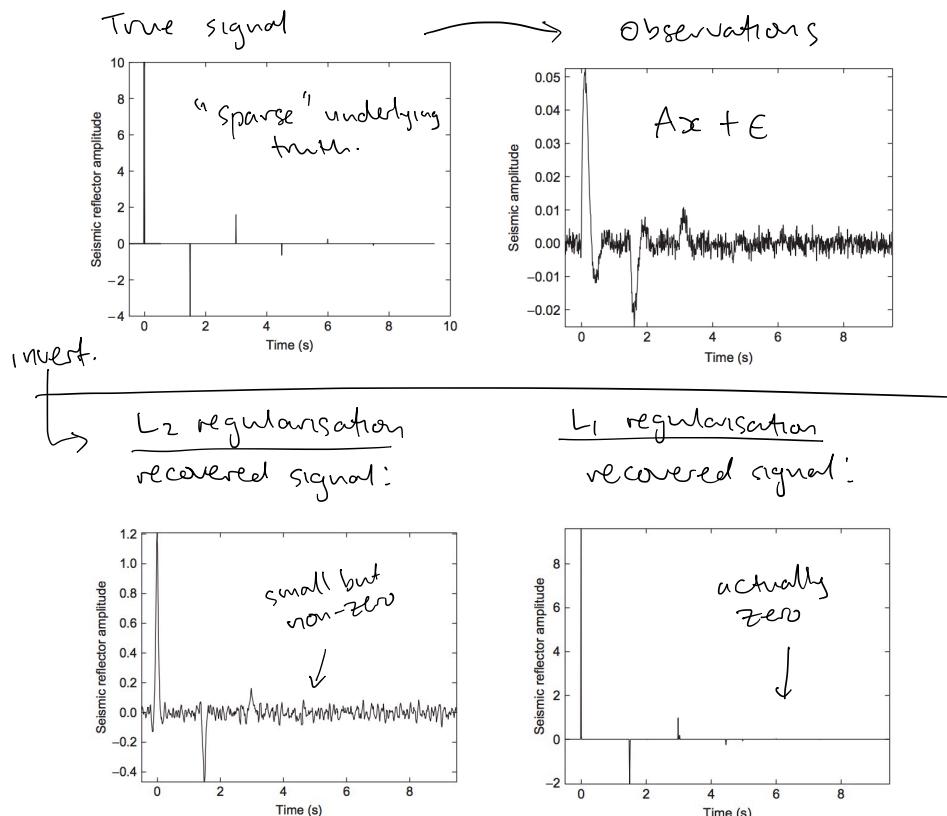
## Examples

- LASSO regression - Google!
  - Hastie et al
  - sklearn (python machine learning)

- Aster et al (Inverse Problems)

Example 7.2 (Full detail attached →)

### Deconvolution seismic sensing



## Applications of $L_1$ : Total variation regularisation

$$\min \|D_1 x\|_1$$

$$\text{s.t. } \|Ax - y\|_2 \leq \delta$$

or variational form:

$$\boxed{\min \|Ax - y\|_2^2 + \lambda \|D_1 x\|_1}$$

→ Tikhonov-like, but  $L_1$  on model 'roughness'  $D_1 x$  instead of  $L_2$

→ piecewise constant (for  $D_1$ ) but allows for small number of discontinuous jumps

↳ good for recovering objects with 'sharp' features like edges in images

→ Solving: same basic idea as before  
(convex / LP / QP / IRLS etc)

Example: Boyd & V. 6.3.3

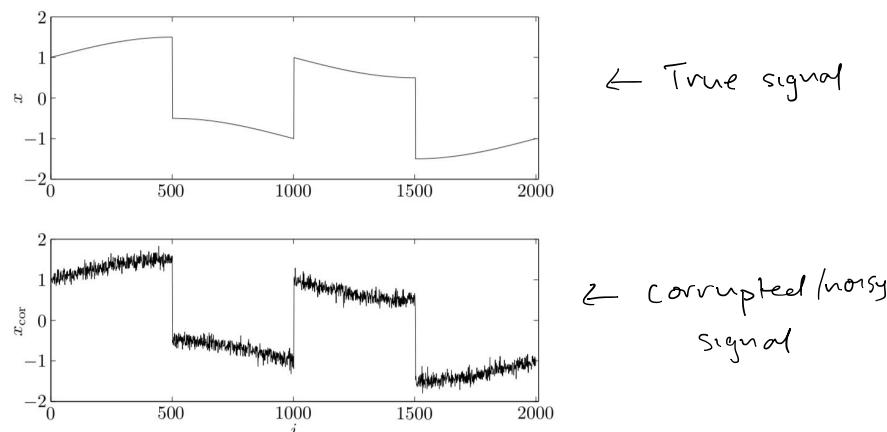


Figure 6.11 A signal  $x \in \mathbb{R}^{2000}$ , and the corrupted signal  $x_{\text{cor}} \in \mathbb{R}^{2000}$ . The noise is rapidly varying, and the signal is mostly smooth, with a few rapid variations.

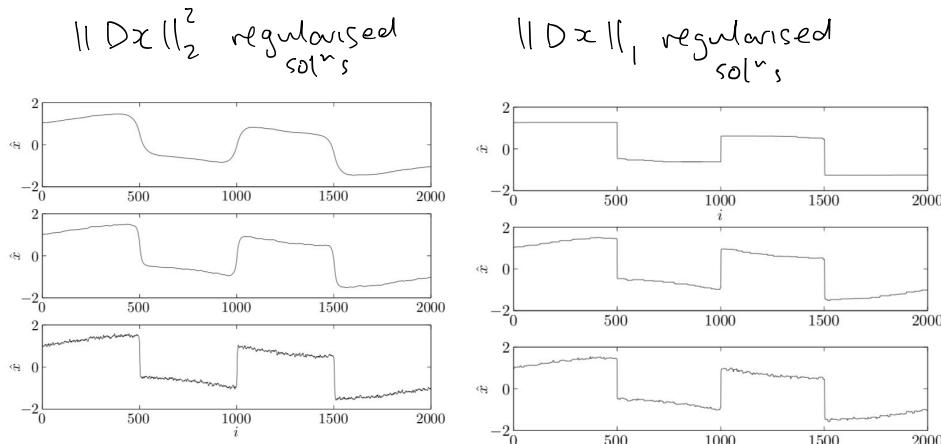


Figure 6.12 Three quadratically smoothed signals  $\hat{x}$ . The top one corresponds to  $\|\hat{x} - x_{\text{cor}}\|_2 = 10$ , the middle one to  $\|\hat{x} - x_{\text{cor}}\|_2 = 7$ , and the bottom one to  $\|\hat{x} - x_{\text{cor}}\|_2 = 4$ . The top one greatly reduces the noise, but also excessively smooths out the rapid variations in the signal. The bottom smoothed signal does not give enough noise reduction, and still smooths out the rapid variations in the original signal. The middle smoothed signal gives the best compromise, but still smooths out the rapid variations.

Figure 6.14 Three reconstructed signals  $\hat{x}$ , using total variation reconstruction. The top one corresponds to  $\|D\hat{x}\|_1 = 5$ , the middle one to  $\|D\hat{x}\|_1 = 8$ , and the bottom one to  $\|D\hat{x}\|_1 = 10$ . The bottom one does not give quite enough noise reduction, while the top one eliminates some of the slowly varying parts of the signal. Note that in total variation reconstruction, unlike quadratic smoothing, the sharp changes in the signal are preserved.

Applications of  $L_1$ : Robust regression

- Here we consider replacing the  $L_2$  norm by the  $L_1$  norm in the data fit term eg

$$\min \|Ax - y\|_1 + \lambda \|x\|_2^2$$

[Variational/Tikhonov form]

LAD:  
Least Absolute Deviation regression

- It turns out that, while minimising wrt the  $L_2$  norm leads to the average as the best data fit, minimising wrt the  $L_1$  norm leads to the median as the best data fit
  - This is much more robust to outliers (extreme observations)
  - Not as easy to work with (non-diff) as before, though still convex for linear
    - ↳ Can eg formulate as a linear programming problem (see appendix)
    - ↳ Again, many packages exist.

First: simple example of 'robustness'

$$y = (2, 3, 5, 7, 8)$$

$$\text{ave}(y) = 5$$

$$\text{med}(y) = 5$$

$$y' = (2, 3, 5, 7, 80)$$

outlier (data here:  
error?)

$$\text{ave}(y') = 19.4 \gg \text{ave}(y)$$

$$\text{med}(y') = 5 = \text{med}(y)$$

Huber et al 'Robust Statistics':

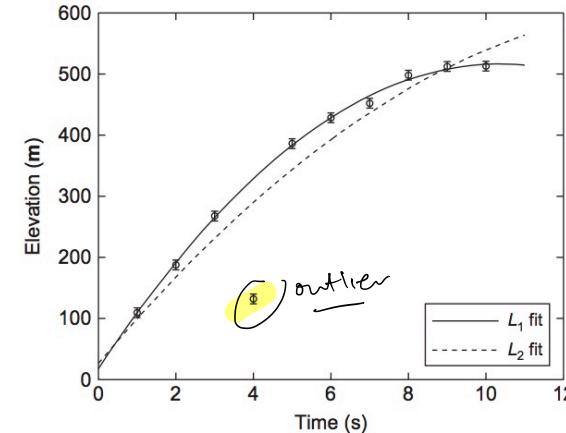
Perhaps the most important purpose of robustness is to safeguard against occasional gross errors. Correspondingly, most approaches to robustness are based on the following intuitive requirement:

*A discordant small minority should never be able to override the evidence of the majority of the observations.*

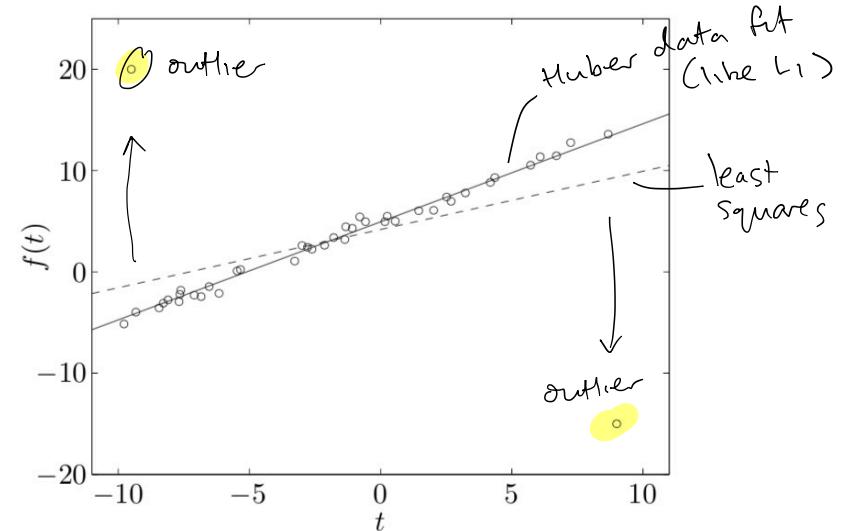
→ Trade-offs! sensitivity vs stability

Examples

(Aster et. al Example 2.4)



(Boyd & V. Example 6.2):



Appendix: reformulation of robust regression  
as linear programming (see eg Boyd & Vandenberghe)

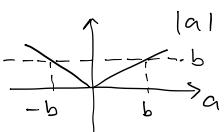
$$\text{V1: } \boxed{\min \|Ax - y\|_1} = |r_1| + |r_2| + \dots + |r_m|$$

$$\text{where } |r_i| = \begin{cases} r_i & \text{if } r_i > 0 \\ -r_i & \text{if } r_i < 0 \end{cases}$$

→ almost linear (piecewise)

$$\text{Note: } |a| \leq b$$

$$\text{equiv to } -b \leq a \leq b$$



$$\& \text{ minimising } |a|$$

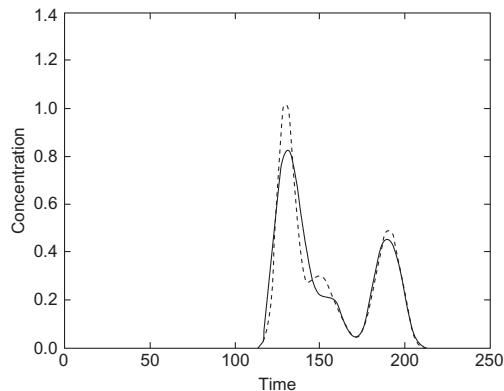
$$\left. \begin{array}{l} \text{equiv to } \min_{a,b} b \\ \text{note: } \rightarrow \text{b var.} \\ \text{st. } |a| \leq b \end{array} \right\} \begin{array}{l} \text{proof?} \\ \text{exercise:} \\ (\text{see also pic } \uparrow) \end{array}$$

i.e.  $-b \leq a \leq b$

So we use

$$\text{V2: } \boxed{\begin{array}{l} \min_{t,x} \mathbf{1}^T t + \mathbf{0}^T x \\ \text{st. } -t \leq Ax - b \leq t \end{array}}$$

Linear  
programming  
(linear objective,  
linear constraints)



**Figure 7.7** Second-order Tikhonov regularization source history solution determined from the L-curve of Figure 7.6, with true history (Figure 7.1) shown as a dashed curve.

bound of 0.36 and an upper bound of 0.73 for the average concentration during this time period. The true concentration average over this interval (Figure 7.1) is 0.57.

## 7.2. SPARSITY REGULARIZATION

In some situations there are reasons to expect that many of the unknown model parameters will be zero. Rather than using Tikhonov regularization to minimize  $\|\mathbf{m}\|_2$ , we may choose to minimize the number of nonzero entries in  $\mathbf{m}$  to obtain a sparse model.

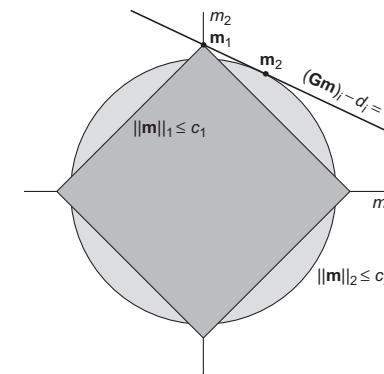
The notation  $\|\mathbf{m}\|_0$  is commonly used to denote the number of nonzero entries in  $\mathbf{m}$  (Note that this 0-norm definition is nonstandard because it is inconsistent with the definition of the  $p$ -norm in (A.85) and does not satisfy the requirements for a vector norm in Section A.7.) We can formulate a corresponding regularized inverse problem as

$$\min \|\mathbf{m}\|_0 \quad (7.4)$$

$$\|\mathbf{Gm} - \mathbf{d}\|_2 \leq \delta.$$

Unfortunately, these kinds of optimization problems can be extremely difficult to solve.

A surprisingly effective alternative to (7.4) is to instead find the least squares solution that minimizes  $\|\mathbf{m}\|_1$ . To see that this is a reasonable approach, consider the set of models with  $\|\mathbf{m}\|_2 = 1$ . Among the models with  $\|\mathbf{m}\|_2 = 1$ , it turns out that the models



**Figure 7.8** Two-dimensional demonstration of the use of model 1-norm minimization (7.5) to obtain sparsity regularization. The square shaded area shows the region  $\|\mathbf{m}\|_1 \leq c_1$ , while the circle shows the region with  $\|\mathbf{m}\|_2 \leq c_2$ . An arbitrary constraint equation in this 2-dimensional model space,  $(\mathbf{Gm})_i - d_i = 0$ , defines a line. The minimum 2-norm residual model satisfying the constraint,  $\mathbf{m}_2$ , will not generally be sparse. However, the minimum 1-norm model satisfying the constraint,  $\mathbf{m}_1 = [0 \ c_1]^T$ , will tend to be sparse due to the presence of corners in the 1-norm contour.

with precisely one nonzero entry of +1 or -1 have the smallest 1-norms (Figure 7.8). Thus, regularizing a least squares problem to minimize its model 1-norm will tend to produce sparse solutions. This tendency for 1-norm regularized models to be sparse becomes even more prominent in higher dimensions. The heuristic approach of minimizing  $\|\mathbf{m}\|_1$  instead of  $\|\mathbf{m}\|_0$  works very well in practice, and recent work [21] has demonstrated reasonable conditions under which the solution to the 1-norm regularized problem is identical to or at least close to the solution of the 0-norm regularized problem (7.4).

The  $L_1$  regularized least squares problem can be written as

$$\begin{aligned} \min \quad & \|\mathbf{m}\|_1 \\ \text{subject to} \quad & \|\mathbf{Gm} - \mathbf{d}\|_2 \leq \delta. \end{aligned} \quad (7.5)$$

Using the standard approach of moving the constraint into the objective function, we can select a positive regularization parameter,  $\alpha$ , so that this is equivalent to

$$\min \|\mathbf{Gm} - \mathbf{d}\|_2^2 + \alpha \|\mathbf{m}\|_1. \quad (7.6)$$

This is a convex optimization problem that can be solved efficiently by many different algorithms. We next present an iterative least squares solution method.

# Boyd & Vandenberghe (Convex Optimization)

312

6 Approximation and fitting

## Quadratic smoothing

The simplest reconstruction method uses the quadratic smoothing function

$$\phi_{\text{quad}}(x) = \sum_{i=1}^{n-1} (x_{i+1} - x_i)^2 = \|Dx\|_2^2,$$

where  $D \in \mathbf{R}^{(n-1) \times n}$  is the bidiagonal matrix

$$D = \begin{bmatrix} -1 & 1 & 0 & \cdots & 0 & 0 & 0 \\ 0 & -1 & 1 & \cdots & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & -1 & 1 & 0 \\ 0 & 0 & 0 & \cdots & 0 & -1 & 1 \end{bmatrix}.$$

We can obtain the optimal trade-off between  $\|\hat{x} - x_{\text{cor}}\|_2$  and  $\|D\hat{x}\|_2$  by minimizing

$$\|\hat{x} - x_{\text{cor}}\|_2^2 + \delta \|D\hat{x}\|_2^2,$$

where  $\delta > 0$  parametrizes the optimal trade-off curve. The solution of this quadratic problem,

$$\hat{x} = (I + \delta D^T D)^{-1} x_{\text{cor}},$$

can be computed very efficiently since  $I + \delta D^T D$  is tridiagonal; see appendix C.

## Quadratic smoothing example

Figure 6.8 shows a signal  $x \in \mathbf{R}^{4000}$  (top) and the corrupted signal  $x_{\text{cor}}$  (bottom). The optimal trade-off curve between the objectives  $\|\hat{x} - x_{\text{cor}}\|_2$  and  $\|D\hat{x}\|_2$  is shown in figure 6.9. The extreme point on the left of the trade-off curve corresponds to  $\hat{x} = x_{\text{cor}}$ , and has objective value  $\|Dx_{\text{cor}}\|_2 = 4.4$ . The extreme point on the right corresponds to  $\hat{x} = 0$ , for which  $\|\hat{x} - x_{\text{cor}}\|_2 = \|x_{\text{cor}}\|_2 = 16.2$ . Note the clear knee in the trade-off curve near  $\|\hat{x} - x_{\text{cor}}\|_2 \approx 3$ .

Figure 6.10 shows three smoothed signals on the optimal trade-off curve, corresponding to  $\|\hat{x} - x_{\text{cor}}\|_2 = 8$  (top), 3 (middle), and 1 (bottom). Comparing the reconstructed signals with the original signal  $x$ , we see that the best reconstruction is obtained for  $\|\hat{x} - x_{\text{cor}}\|_2 = 3$ , which corresponds to the knee of the trade-off curve. For higher values of  $\|\hat{x} - x_{\text{cor}}\|_2$ , there is too much smoothing; for smaller values there is too little smoothing.

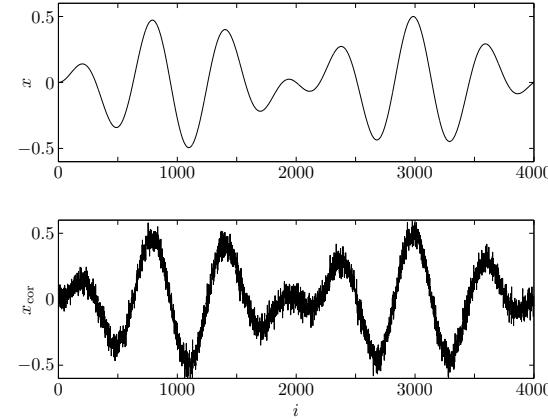
## Total variation reconstruction

Simple quadratic smoothing works well as a reconstruction method when the original signal is very smooth, and the noise is rapidly varying. But any rapid variations in the original signal will, obviously, be attenuated or removed by quadratic smoothing. In this section we describe a reconstruction method that can remove much of the noise, while still preserving occasional rapid variations in the original signal. The method is based on the smoothing function

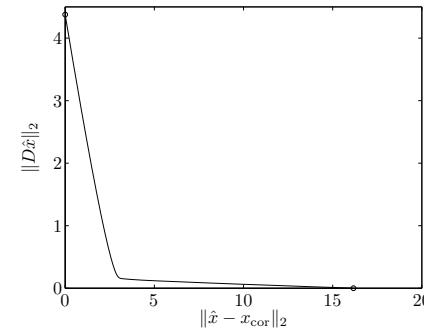
$$\phi_{\text{tv}}(\hat{x}) = \sum_{i=1}^{n-1} |\hat{x}_{i+1} - \hat{x}_i| = \|D\hat{x}\|_1,$$

6.3 Regularized approximation

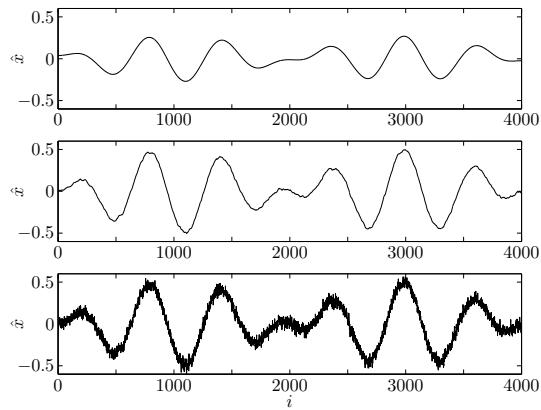
313



**Figure 6.8** Top: the original signal  $x \in \mathbf{R}^{4000}$ . Bottom: the corrupted signal  $x_{\text{cor}}$ .



**Figure 6.9** Optimal trade-off curve between  $\|D\hat{x}\|_2$  and  $\|\hat{x} - x_{\text{cor}}\|_2$ . The curve has a clear knee near  $\|\hat{x} - x_{\text{cor}}\|_2 \approx 3$ .



**Figure 6.10** Three smoothed or reconstructed signals  $\hat{x}$ . The top one corresponds to  $\|\hat{x} - x_{\text{cor}}\|_2 = 8$ , the middle one to  $\|\hat{x} - x_{\text{cor}}\|_2 = 3$ , and the bottom one to  $\|\hat{x} - x_{\text{cor}}\|_2 = 1$ .

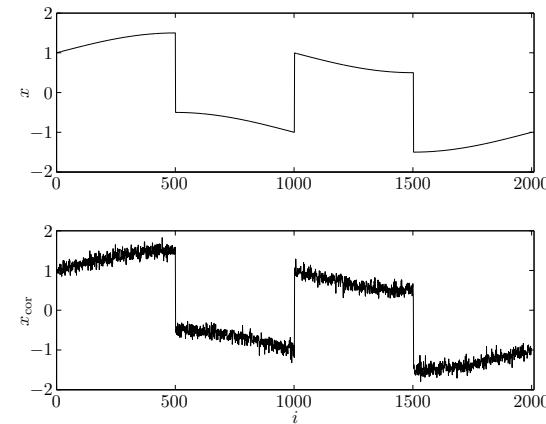
which is called the *total variation* of  $x \in \mathbf{R}^n$ . Like the quadratic smoothness measure  $\phi_{\text{quad}}$ , the total variation function assigns large values to rapidly varying  $\hat{x}$ . The total variation measure, however, assigns relatively less penalty to large values of  $|x_{i+1} - x_i|$ .

#### Total variation reconstruction example

Figure 6.11 shows a signal  $x \in \mathbf{R}^{2000}$  (in the top plot), and the signal corrupted with noise  $x_{\text{cor}}$ . The signal is mostly smooth, but has several rapid variations or jumps in value; the noise is rapidly varying.

We first use quadratic smoothing. Figure 6.12 shows three smoothed signals on the optimal trade-off curve between  $\|D\hat{x}\|_2$  and  $\|\hat{x} - x_{\text{cor}}\|_2$ . In the first two signals, the rapid variations in the original signal are also smoothed. In the third signal the steep edges in the signal are better preserved, but there is still a significant amount of noise left.

Now we demonstrate total variation reconstruction. Figure 6.13 shows the optimal trade-off curve between  $\|D\hat{x}\|_1$  and  $\|\hat{x} - x_{\text{cor}}\|_2$ . Figure 6.14 shows the reconstructed signals on the optimal trade-off curve, for  $\|D\hat{x}\|_1 = 5$  (top),  $\|D\hat{x}\|_1 = 8$  (middle), and  $\|D\hat{x}\|_1 = 10$  (bottom). We observe that, unlike quadratic smoothing, total variation reconstruction preserves the sharp transitions in the signal.

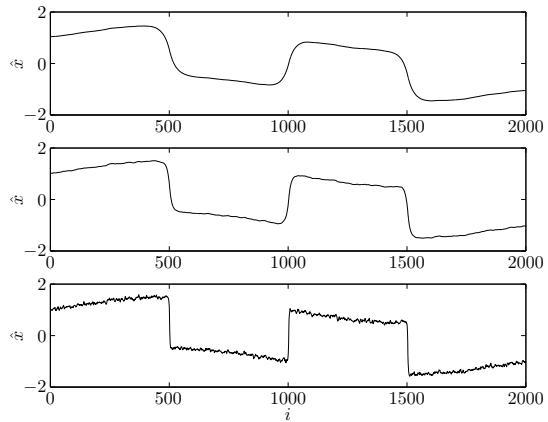


**Figure 6.11** A signal  $x \in \mathbf{R}^{2000}$ , and the corrupted signal  $x_{\text{cor}} \in \mathbf{R}^{2000}$ . The noise is rapidly varying, and the signal is mostly smooth, with a few rapid variations.

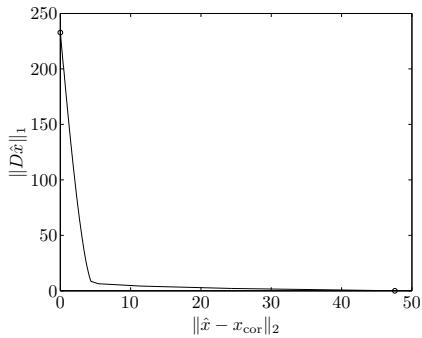
# Boyd & Vandenberghe (Convex Optimization)

316

6 Approximation and fitting



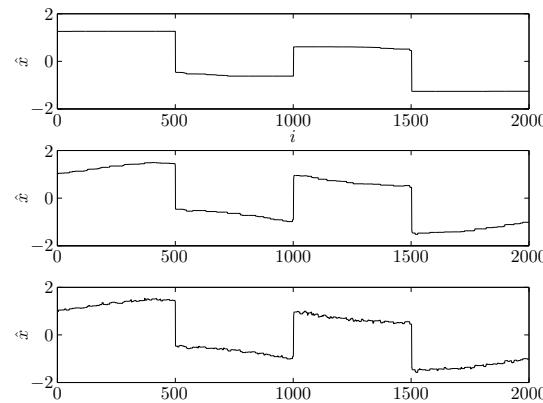
**Figure 6.12** Three quadratically smoothed signals  $\hat{x}$ . The top one corresponds to  $\|\hat{x} - x_{\text{cor}}\|_2 = 10$ , the middle one to  $\|\hat{x} - x_{\text{cor}}\|_2 = 7$ , and the bottom one to  $\|\hat{x} - x_{\text{cor}}\|_2 = 4$ . The top one greatly reduces the noise, but also excessively smooths out the rapid variations in the signal. The bottom smoothed signal does not give enough noise reduction, and still smooths out the rapid variations in the original signal. The middle smoothed signal gives the best compromise, but still smooths out the rapid variations.



**Figure 6.13** Optimal trade-off curve between  $\|D\hat{x}\|_1$  and  $\|\hat{x} - x_{\text{cor}}\|_2$ .

6.3 Regularized approximation

317



**Figure 6.14** Three reconstructed signals  $\hat{x}$ , using total variation reconstruction. The top one corresponds to  $\|\hat{x}\|_1 = 5$ , the middle one to  $\|\hat{x}\|_1 = 8$ , and the bottom one to  $\|\hat{x}\|_1 = 10$ . The bottom one does not give quite enough noise reduction, while the top one eliminates some of the slowly varying parts of the signal. Note that in total variation reconstruction, unlike quadratic smoothing, the sharp changes in the signal are preserved.