

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 2: GENERALISED INVERSES

Topics:

- The algebra (and some calculus) of inverse problems
- Why generalised inverses aren't enough!

Eng Sci 721 : Lecture 2.

Algebra (& a little calculus)

of inverse problems

↳ Resolving lack of existence &/or uniqueness using generalised inverses

↳ Formulating & solving as optimisation problems

Two step vs simultaneous optimisation

→ or, why the generalised inverse also needs regularisation

Algebra (& some calc) of Inverse Problems

Generalised Inverses

Our basic problem can be defined

as :

'solve', ie 'invert',
equations like $F(x) = y$
for x , given y

where :

- x & y could be vectors, functions, images etc
- solutions might not exist, might not be unique &/or might not be stable

Note : mappings, measurements & vectors

we might have an 'exact' model of the form

$$\boxed{y = a + bx}, \text{ for } \boxed{\text{scalars}} x, y$$

If we take a series of 'noisy' measurements, we get eg:

$$y_{\text{obs},i} = (a + bx)_i + e_i$$

$$\Rightarrow y_{\text{obs},i} = a + x_i + e_i$$

for $i = 1, \dots, n$.

These lead to vector eqns in terms of the noisy observations/realisations:

$$\begin{bmatrix} y_{\text{obs},1} \\ \vdots \\ y_{\text{obs},n} \end{bmatrix} = a \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} + b \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix}$$

eg $\boxed{\bar{y}_{\text{obs}} \approx a\bar{1} + b\bar{x}}$ } $\bar{y}, \bar{x}, \bar{1}$ vectors

(I usually drop the explicit overbars on vectors)

Linear or nonlinear? Finite or infinite?

- We will discuss some basic algebra of the problem in the linear & finite dimensional setting
→ ie using Linear Algebra in \mathbb{R}^n
- This can be considered as the 'algebra of linear mappings' in \mathbb{R}^n

→ An important question is:
do these results carry over to the nonlinear &/or infinite-dimensional setting(s)?

Short answer: yes!

(key concepts &
as long as careful)

Linear or nonlinear? Finite or infinite?

- Algebra of arbitrary mappings?

↳ Category theory

(see eg Nashed / MacLaren & Nicholson
for generalised inverses in
category theory setting
→ existence = axiom of
choice!)

- Discontinuous infinite dimensional linear mappings ('ill-posed')

↳ ill-conditioned finite dimensional matrices

→ Theory is a bit beyond scope,
but most of the key
tools are applicable to
nonlinear setting
(we'll solve some of these!)

Linear setting

$$\mathbb{R}^n \xrightarrow{A} \mathbb{R}^m$$

Consider the system of equations

$Ax = y$ } - A is $m \times n$ matrix
- $x \in \mathbb{R}^n$ vector
- $y \in \mathbb{R}^m$ vector

eg

$$\begin{matrix} n \\ m \end{matrix} \xrightarrow{\quad A \quad} \begin{matrix} n \\ m \end{matrix} = \begin{matrix} m \end{matrix}$$

rows: eqns

cols: unknowns

How do we solve when $m \neq n$?

→ $m > n$, more rows than cols.

existence? $\begin{matrix} n \\ m \end{matrix}$ } → eqns > unknowns
→ possibly inconsistent/overdetermined
→ 'more data than parameters'

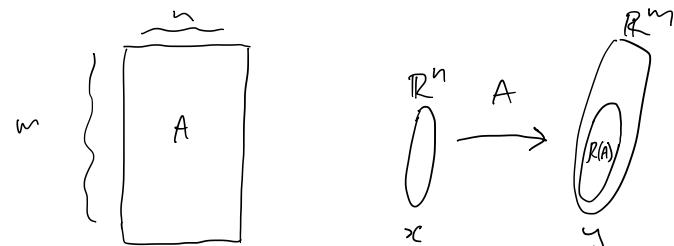
→ $m < n$, more cols than rows

uniqueness? $\begin{matrix} n \\ m \end{matrix}$ } → unknowns > eqns
→ possibly many sol's
→ 'more param. than data'

Tools developed to 'solve' each case

Case 1: possibly inconsistent: more equations than unknowns.

Consider $Ax = y$, A is $m \times n$ & $m > n$



→ Assume w.l.o.g. that all n columns are linearly independent

- $y \in R(A) \Leftrightarrow$ there is a unique solution
- $y \notin R(A) \Rightarrow$ no (exact) solution



1. Inconsistent equations: approximate solution

→ Define $r = y - Ax$ { residual 'error'

→ Measure size of error with a norm $\| \cdot \|$ (see handout for diff. types)

→ Typically assume $\| \cdot \|_2$

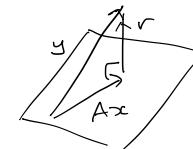
New problem:

minimise $\| y - Ax \|$, A & y given
 $x \in R^n$

- "best approximation"

- "closest approx"

etc.



Minimiser of $\| \cdot \|$ & minimiser of $\| \cdot \|_2^2$

are same ($x \mapsto x^2$ is monotonic for $x \geq 0$)

⇒ least squares approximation!

minimise $\| y - Ax \|_2^2$ (equiv. problem)

Norms, products, summation etc

$$\circ \quad \|x\|_2^2 = \langle x, x \rangle \\ = x^T x$$

$$= \sum_i x_i^2$$

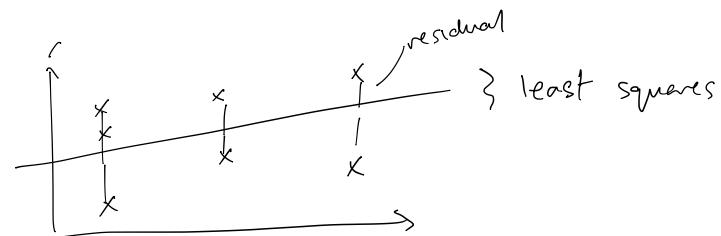
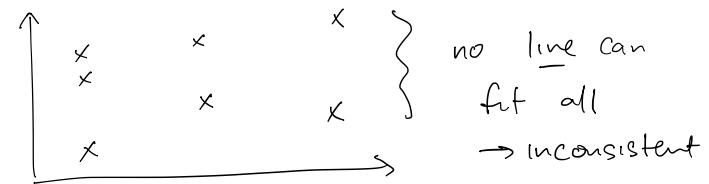
$$\circ \quad Ax = \sum A_{ij} x_j \\ = A_{ij} x_j \quad (\text{Einstein summation convention} \\ \rightarrow \text{sum over repeated indices})$$

so:

$$\circ \quad y^T Ax = \sum_i A_{ij} x_j$$

See handout (?) (we'll do some practice!)

Illustration



Derivation via calc (can do geometrically too!)

$$\min_x \|y - Ax\|_2^2 = \min_x \langle y - Ax, y - Ax \rangle \\ = \min_x f(x)$$

$$\text{where } f(x) = \langle y - Ax, y - Ax \rangle \\ [= (y - Ax)^T (y - Ax)]$$

$$\left. \begin{aligned} &= \langle y, y \rangle - 2 \langle y, Ax \rangle + \langle Ax, Ax \rangle \\ &= y^T y - 2 y^T Ax + x^T A^T Ax \\ &= y_i y_i - 2 y_i A_{ij} x_j + (A_{ij} x_j)(A_{ik} x_k) \end{aligned} \right\} \text{same, different notation.} \\ \text{etc}$$

Differentiating vectors, matrices, tensors wrt vectors, matrices, tensors --

Requires:

(Matrix calculus), (Tensor calculus) etc
→ multiple conventions/notation
(see e.g. wiki page on Matrix Calculus)

I will sketch here, but then give you three key rules! You can use instead of remembering details!

Conventions

→ I'll use $\boxed{d_x f}$ for derivative of f
wrt x , regardless of whether f, x
are scalar, vector etc.

→ use 'Jacobian layout', e.g. vector f, x :

$$\boxed{[d_x f]_{ij} = \frac{\partial f_i}{\partial x_j} = \begin{matrix} f_1 \\ \vdots \\ f_m \end{matrix} \left[\begin{matrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{matrix} \right]}$$

Note, if f is scalar-valued, this implies
 $d_x f$ is a row vector:

$$\boxed{d_x f = \left[\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right]}$$

Hence this is sometimes written

$$\boxed{d_x f = \frac{\partial f}{\partial x^T}}$$

to indicate the 'layout' direction is 'row like'
in x .

We can define the gradient of f as

$$\boxed{\nabla_x f = (d_x f)^T = \frac{\partial f^T}{\partial x}} \quad \text{(often drop transpose on } f \text{)}$$

which again indicates how to 'layout' results.

3 Key rules : (A, a independent of x)

Constant

$$d_x(a) = \frac{\partial}{\partial x^T}(a) = 0^T$$

Linear

$$d_x(Ax) = \frac{\partial}{\partial x^T}(Ax) = A$$

Quadratic

$$\begin{aligned} d_x(x^T Ax) &= \frac{\partial}{\partial x^T}(x^T Ax) \\ &= x^T(A + A^T) \end{aligned}$$

I know these !

Extra:

Additional principles for 'deriving' results
(see handout):

o [Derivative of scalar], possibly dep. on x :

$$d_x(a(x)) = d_x(a) = d_x(a^T)$$

o [Multivariable vector) chain rule]

$$d_x(f(h(x), g(x))) = \frac{\partial f}{\partial x^T}(h, g)$$

$$= (d_g f) \cdot (d_x g) + (d_h f) \cdot (d_x h)$$

$$= \frac{\partial f}{\partial g^T} \frac{\partial g}{\partial x^T} + \frac{\partial f}{\partial h^T} \frac{\partial h}{\partial x^T}$$

Imply eg 'product rule' of form:

$$d_x(h^T g) = h^T d_x g + g^T d_x h$$

$$= h^T \frac{\partial g}{\partial x^T} + g^T \frac{\partial h}{\partial x^T}$$

Back to least squares!

$$\min_x f(x) = y^T y - 2y^T A x + x^T A^T A x$$

$$\Rightarrow \text{set } d_x f = 0^T \quad (1) \quad (\text{row vector})$$

(or $\nabla_x f = 0 \dots$)

3 Rules:

$$d_x (y^T y) = 0^T$$

$$d_x (-2y^T A x) = -2y^T A$$

$$\begin{aligned} d_x (x^T A^T A x) &= x^T (A^T A + (A^T A)^T) \\ &= 2x^T A^T A \end{aligned}$$

$$(1) \Rightarrow -2y^T A + 2x^T A^T A = 0^T$$

$$\Rightarrow -A^T y + A^T A x = 0$$

$$\Rightarrow \boxed{A^T A x = A^T y}$$

Least squares approximation

We have seen this leads to....

The Normal equations:

normal?
 $A^T(Ax-y)=0$
geometric

$$\boxed{A^T A x = A^T y}$$

→ Since we assume the n cols of A
are linearly independent then
 $A^T A$ is invertible (see handout) &
so get unique approximate solⁿ

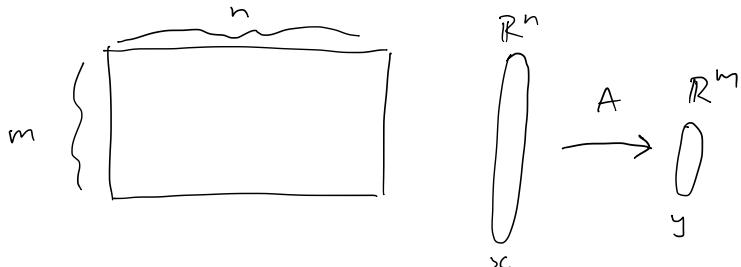
$$\boxed{x^* = (A^T A)^{-1} A^T y}$$

(we will look at what happens if
not LI soon!).

Tools developed to 'solve' each case

Case 2: possibly non-unique: more unknowns than equations

Consider $Ax = y$, A is $m \times n$ & $m < n$



→ Assume w.l.o.g. that all m rows are linearly independent

→ Assume $Ax = y$ has at least one solution, i.e. $y \in R(A)$ (range / image of A)

General solution of $Ax = y$:

$$x = x^* + x_0$$

where $x_0 \in N(A)$ & x^* is any solⁿ of $Ax = y$

(Recall: Nullspace: $N(A) = \{x \mid Ax = 0\}$)

2. Underdetermined equations: least norm solⁿ

- If we want to 'pick out' a single solⁿ (we don't always!), a 'natural' choice is to choose the 'smallest' (think: 'simplest' or 'most efficient' solⁿ)
- Again, this is relative to a particular norm, e.g. $\| \cdot \|_2 = \text{l}_2 \text{ norm.}$

⇒ Problem to solve:

$$\begin{array}{ll} \min & \|x\| \\ \text{s.t.} & Ax = y \end{array}$$

} note: here
assuming
exactly solvable

equivalent to:

$$\begin{array}{ll} \min & \|x\|^2 \\ \text{s.t.} & Ax = y \end{array}$$

} 'least squares
solⁿ'
(cf least squares
approximation)

Least norm problem ('Model reduction')

Using duality or Lagrange multipliers (see eg \rightarrow)

you can show that the least-squares
solution to the minimum norm problem
requires solving the dual problem

$$A A^T v = -2y \text{ for } v^*$$

$$\text{& then obtaining } x^* = -\frac{1}{2} A^T v^*$$

\rightarrow Since we assume the rows of A
are linearly independent then
 $A A^T$ is invertible (see handout) &
so get unique minimum norm

$$x^* = A^T (A A^T)^{-1} y$$

(we will look at what happens if
not L.I. soon!).

[solⁿ sketch : Lagrange multipliers (beyond scope...?)]

$$d_x (x^T x + \lambda^T (Ax - y))$$

$$= 2x^T + \lambda^T A = 0 \quad (1)$$

$$d_\lambda (x^T x + \lambda^T (Ax - y)) = d_\lambda (\lambda^T (Ax - y))$$

$$= d_\lambda ((Ax - y)^T \lambda)$$

$$= (Ax - y)^T = 0^T \Leftrightarrow Ax - y = 0 \quad (2)$$

$$\lambda^T A = -2x^T \quad (1)$$

$$x = -\frac{1}{2} A^T \lambda$$

$$Ax = -\frac{1}{2} A A^T \lambda = y \quad (2)$$

$$A A^T \lambda = -2y$$

$$\lambda = -2(A A^T)^{-1} y$$

$$x = -\frac{1}{2} A^T (-2(A A^T)^{-1} y)$$

$$= A^T (A A^T)^{-1} y$$

]

Summary so far:

- Over-determined: can find least squares approx.
- Under-determined: can find least squares/norm soln

→ no proper 'full' inverse exists in each case, but each is either a left inverse or a right inverse

$$LA = I \quad (\text{left})$$

$$AR = I \quad (\text{right})$$

solve rectangular systems from 'one direction/side'.

Left inverses

$$\mathbb{R}^n \xrightarrow[A]{L} \mathbb{R}^m$$

L is a left inverse (a retraction in category theory) for A if it satisfies

$$LA = I \quad \left\{ \begin{array}{l} A: m \times n \\ L: n \times m \\ I: n \times n \end{array} \right.$$

For least squares data approximation, we have

$$L = (A^T A)^{-1} A^T \quad \& \quad x = Ly$$

A left inverse exists when ~~rows~~ > ~~cols~~ of A & the cols are LI (think \downarrow map into bigger space)

$$\begin{matrix} n \\ m \end{matrix} \xrightarrow{\quad} \begin{matrix} n \\ n \end{matrix} = \begin{matrix} m \\ m \end{matrix}$$

$$\begin{matrix} \mathbb{R}^n \\ \text{models} \end{matrix} \xrightarrow[\text{② retract}]{\text{① map } (A)} \begin{matrix} \mathbb{R}^m \\ \text{data} \end{matrix}$$

'shrink data to make solvable'

① then ② vs. ② then ①
is identity on model space

shrinks/projects \mathbb{R}^m data to subspace of \mathbb{R}^m

Right inverses

$$\mathbb{R}^n \xrightarrow[A]{R} \mathbb{R}^m$$

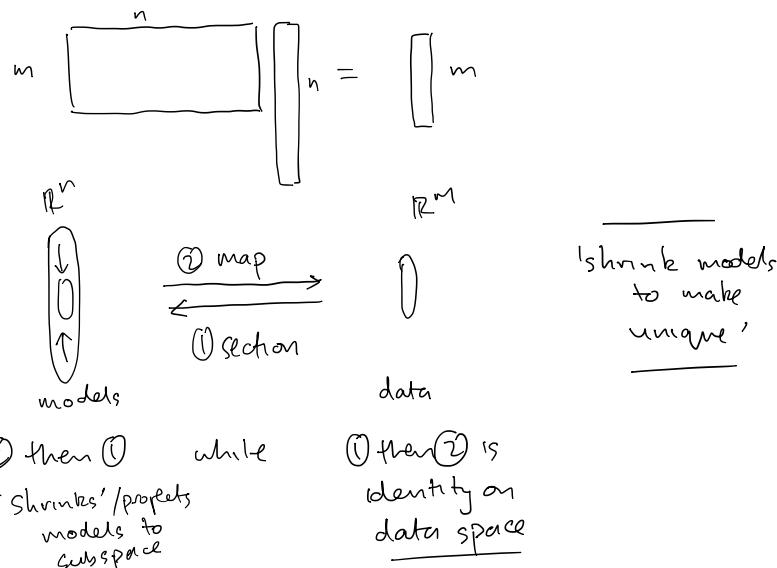
R is a right inverse (a section in category theory) for A if it satisfies

$$AR = I \quad \left\{ \begin{array}{l} A: m \times n \\ R: n \times m \\ I: m \times m \end{array} \right.$$

For least norm problems, we have

$$R = A^T (A A^T)^{-1} \quad \& \quad b_c^* = R y$$

A right inverse exists when cols \Rightarrow rows of A & the rows are LI



Unification: generalised inverses

→ We can unify the solution of these problems with 'generalised inverses'

→ This will solve both existence & uniqueness issues at the same time!

└ spoiler alert: but not stability issues!

| we will see an alternative way

| to combine/unity the solutions when we discuss regularisation |)

The (actually, a - see later) generalised inverse can be characterised as solving the minimum norm approximation

problem,

⇒ it solves

$$\min \|y - Ax\|$$

AND

$$\min \|x\|$$

(in a particular way)

Unification: generalised inverses

- Generalised (or pseudo, for the special cases we look at) inverses solve the two step optimisation prob:

• Stage 1: minimise $\underset{x}{\|y - Ax\|}$ or $\|y - Ax\|^2$

Then

• Stage 2: minimise $\|x\|$ or $\|x\|^2$ among all solutions to stage 1.

The result is a matrix A^+ ← dagger symbol that

- is a left inverse if one exists
- is a right inverse if one exists
- 'close to' left/right inverse if neither exist.

Generalised vs pseudo-inverses

The most general algebraic characterisation of a generalised inverse A^+ is just

$$(1) \boxed{AA^+A = A}$$

Note: • suppose $A^T A$ is invertible

⇒ leastsq left inverse $L = (A^T A)^{-1} A^T$ exists

$$A^T \cdot (1) \rightarrow A^T A A^+ A = A^T A$$

$$\begin{aligned} A^+ A &= (A^T A)^{-1} \underbrace{A^T A}_{M} \\ &= (M)^{-1} M \\ &= I \quad \left. \begin{array}{l} \{ A^+ \text{ is a left inverse} \\ \text{too} \end{array} \right. \end{aligned}$$

• now suppose $A A^T$ is invertible

⇒ leastnorm right inverse $R = A^T (A A^T)^{-1}$ exists

$$(1) \cdot A^T \rightarrow A A^+ A A^T = A A^T$$

$$\begin{aligned} A A^+ &= (A A^T) (A A^T)^{-1} \\ &= I \quad \left. \begin{array}{l} \{ A^+ \text{ is a right inverse} \\ \text{too} \end{array} \right. \end{aligned}$$

Generalised vs pseudo-inverses

We don't have to take

- $A^+ = L = (A^T A)^{-1} A^T = \text{least squares}$
- $A^+ = R = A^T (A A^T)^{-1} = \text{least norm}$

in each case.

Generalised inverses defined via (1)

are not unique

→ add extra conditions } eg different
norms etc
→ depends on goals } diff. inverses

Pseudo?

→ (& will assume unless otherwise)

Eg what we've been using (least sq/norm)
is the Moore-Penrose Pseudo-inverse

This satisfies

$$\boxed{\begin{array}{l} A A^+ A = A \\ A^+ A A^+ = A^+ \\ (A^+ A)^T = A^+ A \\ (A A^+)^T = A A^+ \end{array}}$$

least squares,
least norm
generalised
inverse
 $\stackrel{=}{}$
'pseudo-inverse'

Model resolution, data resolution operators

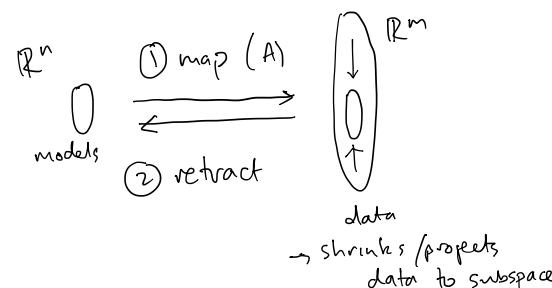
$$\boxed{R_D = A A^+} \quad \left. \right\} \text{ how much data is 'shrunk' or smeared}$$

$$\boxed{R_M = A^+ A} \quad \left. \right\} \text{ how much model is 'shrunk' or smeared}$$

Least squares data approx (A^+ is left inverse)

$$\begin{matrix} n \\ m \end{matrix} \xrightarrow{} \begin{matrix} n \\ n \end{matrix} = \begin{matrix} m \\ m \end{matrix}$$

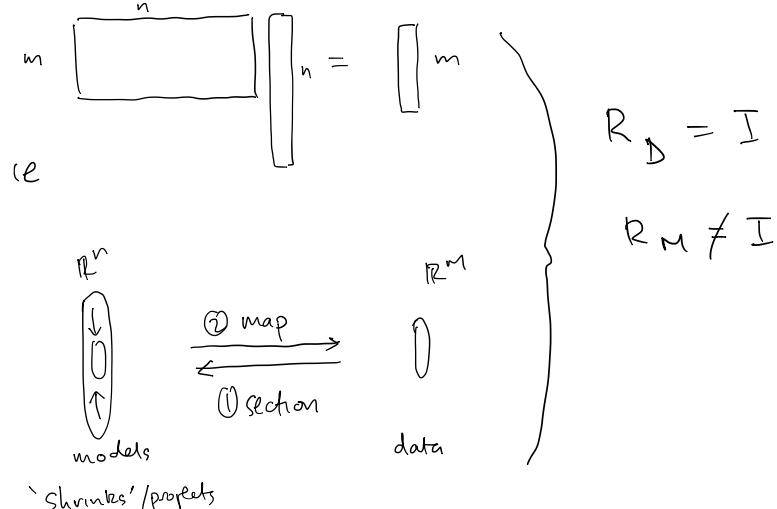
ie



$$R_D \neq I$$

$$R_M = I$$

least squares model reduction (A^+ is right inverse)



'shrink's' / projects
models to
subspace

Subtle point (see later):

we actually want $R_D \neq I$

in above, if we have noise!

Example

projectile motion (Aster et al Ex. 1.1)

$$\boxed{y(t) = a + bt - 0.5ct^2}$$

→ estimate a, b, c given we observe at m times: } def $\theta = \begin{pmatrix} a \\ b \\ c \end{pmatrix}$

$$\theta \in \mathbb{R}^n = \mathbb{R}^3$$

$$y_1 = a + bt_1 - 0.5t_1^2$$

$$y_2 = a + bt_2 - 0.5t_2^2$$

$$\vdots$$

$$y_m = a + bt_m - 0.5t_m^2$$

i.e.

$$\begin{bmatrix} 1 & t_1 & -0.5t_1^2 \\ 1 & t_2 & -0.5t_2^2 \\ \vdots & \vdots & \vdots \\ 1 & t_m & -0.5t_m^2 \end{bmatrix} \begin{bmatrix} a \\ b \\ c \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

i.e.

$$\boxed{F \cdot \theta = y}$$

[Note: linear in parameters!]

Example : least squares (param < data)

→ More than three observations
eg 4 obs., 3 param

```
def fmap(tobs):
    A = np.zeros((len(tobs), 3))
    for i, ti in enumerate(tobs):
        A[i, :] = np.array([1, ti, -0.5*ti**2])
    return A

#true parameters
theta_true = np.array([10, 100, 9.81])

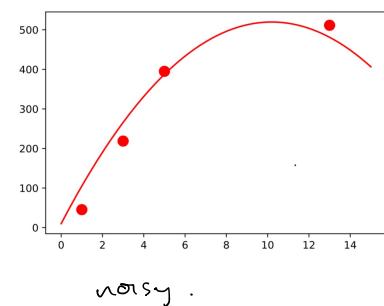
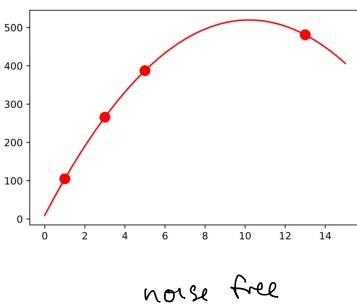
#fine time grid
t = np.linspace(0, 15, 1000)

#observation times
#tobs = np.array([1, 13]) #under-determined
tobs = np.array([1, 3, 5, 13]) #over-determined

#forward map
Aobs = fmap(tobs)

#observed data
yobs = np.dot(Aobs, theta_true) #noise-free
#yobs = np.dot(Aobs, theta_true) + np.random.normal(0, 30, size=len(tobs))

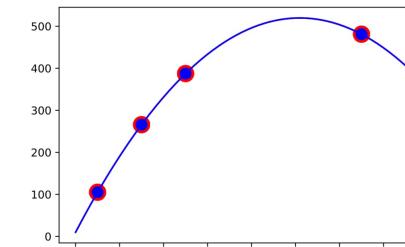
#plots
plt.plot(tobs, yobs, 'ro', markersize=10)
plt.plot(t, x, 'r')
plt.show()
```



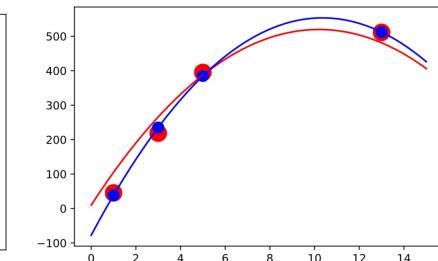
Invert :

Apinv = np.linalg.pinv(Aobs)

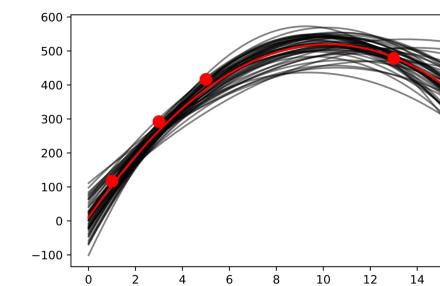
noise free



noisy



repeated sampling



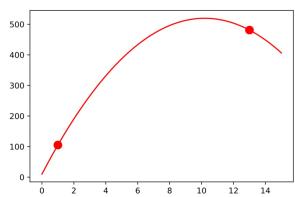
Exercise! explore what the parameters are doing compared to true values (this & next)

Example : least norm (param > data)

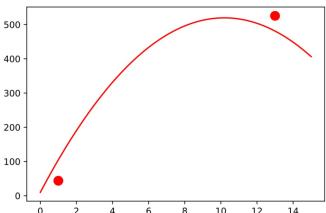
→ less than three observations

Solutions

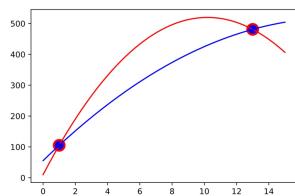
noise free data



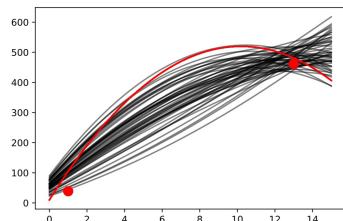
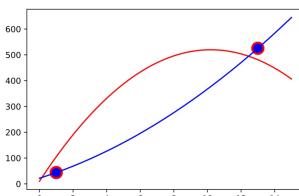
noisy data



recovered (blue)



recovered (blue)



repeated sampling

Note: we've used the naive model norm $\|\theta\|_2$ including the constant term → here curved actually 'smaller' than straight (see later)

So are we done?

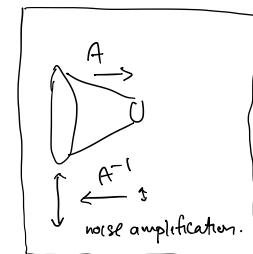
Recall : two stage defⁿ of generalised inverse

- First minimise data fit
- Then minimise model

Inverse problems:

Typically

- underdetermined
- also have noise



So: generalised inverse will exactly fit noise!

→ unstable (just like usual inverse)

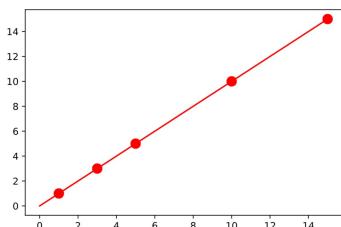
Example : instability despite model 'reduction'

- Eg - high degree polynomial \Rightarrow possibly non unique.
- plus observation noise \Rightarrow possibly inconsistent

$$\boxed{y = a_0 + a_1 x + a_2 x^2 + \dots + a_n x^n} \text{ class}$$

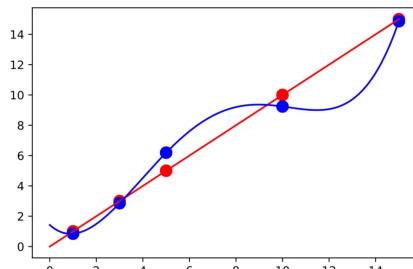
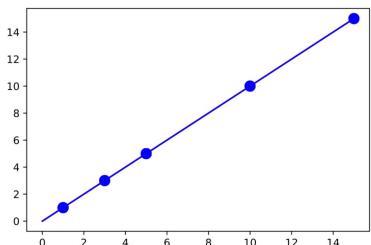
But true model: $a_i = 0$ if $i \neq 1$
 $a_1 = 1$

'True' data: (5 observations)



noise
free
recovery

add noise &
recover



unstable! Not reduced
enough!

Trade-offs & regularisation

→ Just as we saw in L1 that having an inverse in principle is not enough,

$\boxed{\text{having a generalised inverse}} \\ \text{is } \underline{\text{not}} \text{ enough}$

→ stability is still a key issue.

However, the 'extremal' or variational characterisation of generalised inverses provides a as to how to control stability as well!

Trade-offs & regularisation

o Instead of two-step:

o Stage 1: minimise $\underset{x}{\|y - Ax\|}$ or $\|y - Ax\|^2$

Then

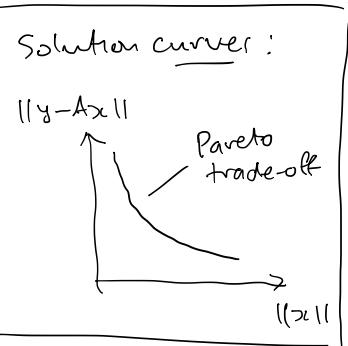
o Stage 2: minimise $\|x\|$ or $\|x\|^2$ among
all solutions to stage 1.

o Try simultaneous minimisation:

vector/multi objective problem

$$\min_x (\|y - Ax\|, \|x\|)$$

simultaneously



→ allows us to filter noise in
underdetermined case, while
still 'shrink' or reducing
models to get 'simple' solutions