

ENGSCI 721

INVERSE PROBLEMS

Oliver Maclaren
oliver.maclaren@auckland.ac.nz

MODULE OVERVIEW

Inverse Problems (Oliver Maclaren) [7-8 lectures/2 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 [1-2 lectures]

Intro to the statistical viewpoint on inverse problems. Bayesians, Frequentists and all that. Basic frequentist ideas.

LECTURE 7/8: INTRODUCTION TO STATISTICAL APPROACHES TO INVERSE PROBLEMS

Topics:

- From approximate models to probabilistic models
- Expectation as a generalised inverse
- Likelihood and maximum likelihood estimation
- Bayesians, Frequentists and all that
- Basic ideas of frequentist evaluation of methods

EngSci 721 Lecture 7.18

Intro. to statistical approach to
inverse problems.

Warning: fraught topic!

- Statistics lies at the intersection of
 - applied science methodology
 - experimental design
 - philosophy of science
 - mathematics
 - etc!
- statisticians (& non-statisticians) constantly argue about foundations!

↳ rule of thumb: assume everyone (incl. me.., but esp. the dogmatic advocates ...) are wrong about some parts.

↳ different approaches have diff. pros & cons & domains of application

How to navigate?

Wittgenstein's ladder:

- use some 'ladder' concepts to help 'climb' towards understanding
- be aware that you likely need to throw them away to reach a 'higher level' of understanding!

(True of all learning, but especially for statistics)

Starting point: constant model

" $y \approx \theta$ ": $\theta \rightarrow \boxed{F} \rightarrow y$

Above is theoretical 'template' for model runs.
(realisations)

Consider 'running' model m times, get:

$$\{(\theta, y)_1, (\theta, y)_2, \dots (\theta, y)_m \}$$

If θ is constant (fixed each time)

but only have $y \approx \theta$, get

$\{(\theta, y_1), (\theta, y_2), \dots (\theta, y_m) \}$, where $y_i \neq y_j$ for $i \neq j$
(in gen.)

ie

$$\bar{y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_m \end{bmatrix} \approx \begin{bmatrix} 1 \\ 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \theta = \bar{1}_m \theta$$

↑
using explicit overbar
here to distinguish from scalar y , but will drop

vector of m ones

From the perspective we've taken so far, to det. θ from y , we want to 'solve' the tall system:

" $y = \bar{1}_m \theta$ " (vector eqn of obs., have dropped overbars)

ie

$$\begin{bmatrix} y \\ 1 \\ \vdots \\ 1 \end{bmatrix} = \begin{bmatrix} \theta \end{bmatrix}$$

for θ given y (vector of obs.).

Note: Overdetermined system (m eqns, one unknown)

$A = \begin{bmatrix} 1 \end{bmatrix}$ is a tall matrix with $L I$ cols (single col.)
⇒ no proper inverse, but has generalised inverse

Also: $\bar{1}_m^T \bar{1}_m = \underbrace{\begin{bmatrix} 1 & 1 & \dots \end{bmatrix}}_m \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \}_{m \times 1} = m$

& $\bar{1}_m^T y = \{1 \ 1 \ \dots 1\} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} = \sum y_i$

→

so:

$$A^+ = (A^T A)^{-1} A^T = (I_m^T I_m)^{-1} I_m^T$$

$$= (m)^{-1} I_m^T$$

$$= \frac{1}{m} I_m^T$$

& our (pseudo-inverse) solⁿ is

$$\theta^+ = A^+ y = \frac{1}{m} I_m^T y$$

$$= \frac{1}{m} \sum_i y_i$$

$$= \text{mean}(y)$$

$$= \overline{E_m(y)}$$

sample based expected
value.

(see Ru's
material
for more
on E)

→ ie when estimating a constant
from noisy data, the generalised
inverse solves the over-determined
tall system by using:

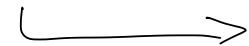
$$\theta = \text{sample mean}(y)$$

$$= \overline{E_m(y)}.$$

→ The data resol. operator $R_D = A A^+ = I_m \overline{E_m}$
projects y to the constant vector with same
mean as y : can only 'resolve' average.

This approach is quite general: don't need explicit
assumptions on the 'distribution of errors'
→ pros & cons

→ many (not all!) statisticians
prefer to explicitly model the
errors probabilistically



Probabilistic 'error':

- Capture all sources of 'noise'
- measurement error
- missing factors etc.

For example, additive error model:

$$y = \theta + e, \quad E[e] = 0$$

↑ ↑ ↑ ↙
 random still random mean of
 variable const. variable error is zero
 ↙ ↙ ↙
 obs = signal + noise

Realisations:

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \theta + \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_m \end{bmatrix}$$

↑ ↑ ↑
 observe want don't observe, but
 know something
 about, $E[e] = 0$

- What does this mean?
- How to solve in the statistical sense?

Meaning?

- From approx. models to probability models:

Old story

$$\theta \rightarrow \boxed{F} \rightarrow y$$

'template'
 for approximate
 model ($y_{obs} \approx y$)

New story

$$\theta \rightarrow \boxed{F} \rightarrow y = \theta + e$$

'template' for
 realisations from
 random vars additive error
 probabilistic model.

More generally

$$\theta \rightarrow \boxed{F} \rightarrow y \sim P_Y(y)$$

'template' for
 realisations
 from arbitr.
 probabilistic
 model for y .

- Forward mapping maps from parameters, theories etc to random variables/prob. dist. for these random variables.
- We observe a finite number of realisations of these random vars/from these distr.

How to solve in statistical approach?

Many approaches! Depends on 'philosophy'.

Approach 1.

1A. Define 'ideal' (infinite data) target } population
1B. 'Plug-in' sample analog } vs. sample

Consider:

$$Y = \theta + e, \mathbb{E}[e] = 0$$

} 'ideal' / 'template' /
'recipe' for generating
as many realisations
as you want

$$\begin{aligned} (1A). \mathbb{E}[Y] &= \mathbb{E}[\theta + e] \\ &= \mathbb{E}[\theta] + \mathbb{E}[e] \text{ (linearity)} \\ &= \theta + 0 \\ &\quad \text{---} \\ &\quad \text{constant by assumption} \end{aligned}$$

$$\Rightarrow \theta = \mathbb{E}[Y] \quad \left. \begin{array}{l} \text{note expected} \\ \text{value under} \\ \text{'infinite' observations} \end{array} \right\}$$

(1B) Sample version of expected value:

$$\theta = \mathbb{E}[Y] \approx \mathbb{E}_m[Y] = \frac{1}{m} \sum_i y_i$$

$$\Rightarrow \theta_{\text{est}} = \theta^+ = \text{sample mean}$$

→ same as before!

Morals: → There is a close correspondence
between generalised inverses, projections
etc & expectations of random vars.

→ If can't solve in usual sense
(e.g. over-determined), we
can solve in 'average sense'
↳ match expected value.

→ We can 'plug in' sample versions of
our 'ideal' quantities to get estimates

→ Same sort of ideas carry over
to non-trivial models,
conditional expectations etc.

→ The 'functional analysis' approach
to probability and statistics

- 'Probability via expectation' by Whittle
- Hilbert space methods in probability
& statistics by Small & McLeish

[Uses prob. ideas but only 'unusually']

Estimate?

We used a sample-based estimate of θ :

$$\hat{\theta}_{\text{est}} = \frac{1}{m} \sum_i y_i \approx \mathbb{E}[y] = \theta$$

Is this a good estimate? What other ways are there of estimating things?

First: More detailed probability modelling

Instead of just

$$Y = \theta + e, \quad \mathbb{E}[e] = 0$$

introduce a more detailed model of e :

$$Y = \theta + e, \quad e \sim N(0, \sigma^2)$$

$$\text{where } e \sim N(0, \sigma^2)$$

means 'normally distributed'

with zero mean & variance σ^2

$$\mathbb{E}[e] = 0 \quad \text{Var}(e) = \sigma^2$$

→ see Ru's part for (much) more!

→

Normal/Gaussian distribution

→ Ru's favourite topic

→ just mention briefly for now...

The probability density for a vector y representing an m -dimensional IID sample of size m from a normal (Gaussian) distribution with constant mean θ & constant variance σ^2 is:

$$P_Y(y; \theta) = \frac{1}{(2\pi\sigma^2)^{m/2}} \cdot \exp\left\{-\frac{1}{2} \left\|\frac{y - \theta}{\sigma}\right\|^2\right\}$$

('given θ ' ; when θ isn't vec. a random var. some, esp. Bayesians just use $P(y|\theta)$ instead of $P(y;\theta)$).

(assuming σ known here)

(θ scalar, 1θ vector)

Likelihood

Given an observation/realisation y of Y ,

the likelihood function is defined as

$$L(\theta; y) \propto P(Y=y; \theta) dy$$

($P(Y=y) dy$
 $\approx P(Y=y)$)

i.e

- is a function of the parameter
- for each parameter value is proportional to the probability of the observed data under the corresponding model.

Intuitively, the likelihood measures how well the model corresponding to a given parameter value 'fits' the data.

→ if $P(Y=y; \theta)$ is high for a given θ then the model gives high probability to the data & hence provides a good 'fit'

→ note: often write $P(Y; \theta) = P(Y|\theta)$
but still:

$$P(Y=y; \theta) \neq P(\theta | Y)$$

→ likelihood is not a probability of the parameter (is prob. of data)

→ it measures how 'compatible' a model & data are

→ doesn't require θ to be random

→ isn't additive
etc.

'Maximum likelihood' estimation determines

a point estimate by solving:

$$\hat{\theta}_{ML} = \arg \max_{\theta} L(\theta; y) \quad \left. \right\} \begin{matrix} \text{"parameter that} \\ \text{makes data} \\ \text{most probable"} \end{matrix}$$

Note: $\arg \max_{\theta} F(\theta) = \arg \max_{\theta} \ln F(\theta)$ under natural log.
under general conditions. (ln monotonic)

$$\arg \max_{\theta} A \cdot F(\theta) + B = \arg \max_{\theta} F \quad \left. \right\} \begin{matrix} \text{for } A, B \text{ indep. of } \theta, A > 0 \end{matrix}$$

& $\arg \max_{\theta} F(\theta) = \arg \min_{\theta} -F(\theta)$

Now:

Define $\lambda(\theta; y) = \ln L(\theta; y)$
= "log-likelihood function"

$$\Rightarrow \left. \begin{matrix} \hat{\theta}_{ML} = \arg \max L(\theta; y) \\ = \arg \max \lambda(\theta; y) \\ = \arg \min -\lambda(\theta; y) \end{matrix} \right|$$

For normal distribution as before:

$$L(\theta; y) = C \cdot \exp \left\{ -\frac{1}{2} \left\| \frac{y - \theta}{\sigma} \right\|^2 \right\}$$

for some constant C .

$$\Rightarrow \ln L(\theta; y) = \ln C - \frac{1}{2} \left\| \frac{y - \theta}{\sigma} \right\|^2$$

$$\begin{aligned} \Rightarrow \arg \max_{\theta} \ln L(\theta; y) &= \arg \min \{-\ln L(\theta; y)\} \\ &= \arg \min_{\theta} \frac{1}{2\sigma^2} \|y - \theta\|^2 \\ &= \arg \min_{\theta} \|y - \theta\|^2 \end{aligned}$$

Take away: max likelihood estimate
under additive normal errors
= least squares estimate

(= estimate from generalised
inverse
= sample plug-in)
[≠ for fixed σ .]

Regularisation?

See 'lecture on maximum likelihood estimation' on canvas (from a short course I presented at)

for more on likelihood-based estimation

→ Make up more observations, stack & do max likelihood/least squares
('mixed estimation' in freq. econometrics,
'prior likelihood' in 'pure likelihood' approach)

→ prior distribution & 'MAP' estimation
(Bayesian approach → see Ru's part)

(see L3 comments)

So ... many approaches give same point estimate

→ Is this a 'good' estimate?

→ what about 'interval' estimates & uncertainty?

But will cover Bayes... So...

I will sketch some frequentist ideas.

Note: Beware of Bayesian descriptions of frequentist inference & vice-versa
--- often dogmatic, often wrong!
Also applies to Bayesian descriptions of Bayes & freq. description of freq.!

With that in mind ...

What is a 'good estimate' of θ ?

Bayesians

Frequentists

- want a 'posterior probability' $P(\theta|y)$ that represents their personal belief / 'state of information' etc about the true value of θ given data y & prior beliefs $P(\theta)$. Point estimate might be max, mean, median of $P(\theta|y)$
- to make it more 'objective', need to agree on the prior
- can think of in terms of (personal) decision theory & maximising subjective expected utility.
- want more 'objective' approach: prob. represents 'variability in the world' rather than personal uncertainty
- focus on evaluating methods of drawing conclusions by evaluating performance under 'repeated use', given noisy observations/realisations of random processes
- a 'good' estimate is an output of a reliable method ... performs well in the 'long run' ... but did it perform well in this case?
- can think of in terms of a two player game between you & nature

Frequentist statistics : trust the process

Provides methods of evaluating the performance of estimation methods, under repeated use in a 'noisy' environment (noisy observations/realisations of random variables), regardless of where the methods come from... Eg can do:

'frequentist performance of Bayesian methods'!

Methods?

An estimator is a method (function) that provides estimates of an estimand (parameter/target) given observations

→ Eg the generalised inverse is an estimator:

$$\hat{\theta} = A^+ y$$

↑ estimate of θ A^+ estimator
data (realisation of random var)

where:
A ← forward model
parameter $\theta \rightarrow Y \leftarrow$ data (or sample)
 $A^+ \leftarrow$ estimator

Frequentist statistics as a two player game.

- nature chooses fixed θ & generates random data $Y \sim P(Y; \theta)$
- you come up with a 'guessing method' (estimator) that gives an estimate of θ given any realisation of Y (ie you see data, not θ)

$$\begin{array}{|c|} \hline \text{estimate} = T(Y) \\ \text{for } \theta \\ \hline \end{array}$$

↑ value ↑ method data realisation

→ The estimate is random because Y (data) is random, even if θ fixed (but unknown)

Here:

- Probability applies to
 - 'actual variation' in nature (aleatoric)
 - The estimates that the guessing method 'spits out' (since estimator is a function of the data, $T(Y)$)
- Does not apply to 'fixed but unknown' things (eg θ)
 - Not rep. of personal uncertainty (epistemic)

Types of estimate: Many!

Point estimate: 'single best guess' $T(Y) = \theta_{\text{est}} \in \text{parameter space}$ } eg max likelihood estimate

Interval/Set estimate: a set (eg interval) of parameter values indicating 'good', 'plausible' etc ranges of parameter values/estimates, ie $T(Y) = S_{\theta_{\text{est}}} \subseteq \text{param space}$
↳ eg confidence interval.

Function or distribution estimate:

(* more common for Bayes)

a distribution or function over parameter values, indicating 'good', 'plausible' etc ranges of parameter values/estimates.
→ likelihood function
→ posterior distribution
→ confidence distribution etc.

Performance:

Point estimator

- Want to guarantee that for all θ we have good expected performance in repeated use, eg want:
 $\mathbb{E}[d(T(Y), \theta)]$ for $Y \sim P(Y; \theta)$ to be small.

Example: want $\mathbb{E}[\|T(Y) - \theta\|^2]$ small for all θ & $Y \sim P(Y; \theta)$
(low mean squared error)

Interval estimator

- Want to guarantee that for all of nature's choices of θ , our (random) interval guesses will often (eg 95% of the time we get data) contain ('trap') the true but unknown value.

$C(Y)$ is 95% confidence interval procedure $\rightarrow \{ | \frac{P(C(Y) \ni \theta; \theta)}{P(\text{random fixed.})} \geq 0.95, \text{ for } Y \sim P(Y; \theta) \}$ — coverage (confidence) level

Confidence intervals?

A confidence interval procedure is a method of guessing intervals given data realisations:

eg $[L(Y), U(Y)]$ } random*
 ↓ ↓
 lower upper } random interval guess
 limit limit (+depends on random data)

The method has a given coverage probability (eg 95% prob) which is the prob. of containing the true (usually fixed but unknown) value θ , ie $P([L(Y), U(Y)] \ni \theta; \theta)$, under repeated samples/realisations of Y .

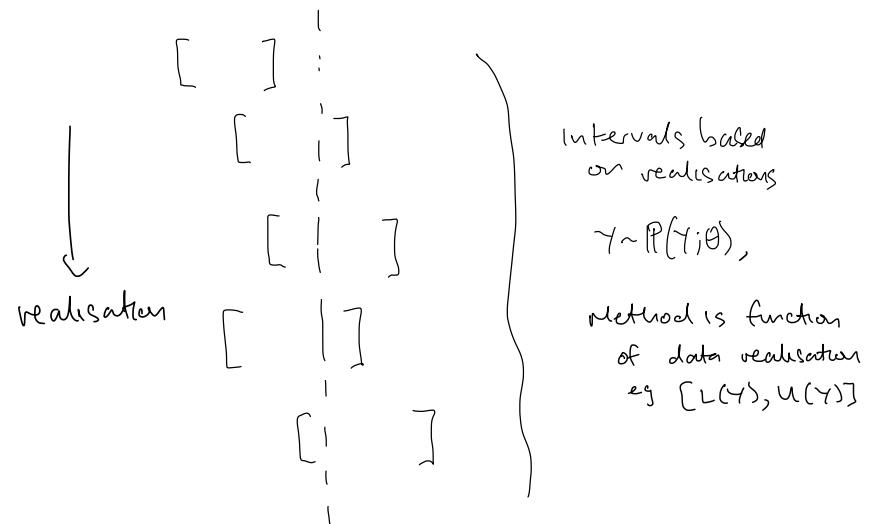
A given interval realisation either contains the true value or doesn't. Not 95% prob.

↳ compare $P(Y=3)$ to $P(4=3)$ } doesn't make sense

eg $P_Y([L(Y), U(Y)] \ni \theta; \theta) \checkmark$

$P_Y([3, 5] \ni \theta; \theta) \times$ (either 0 or 1.)

Coverage $\theta \leftarrow$ nature chooses



Above 4/5 intervals trap the true value (empirical coverage = 0.8)

Coverage checks

- given a method, you can play the role of nature to check performance.
 → choose a θ , generate data from $P(Y; \theta)$, see how often interval traps true value.

→ get coverage for each choice of θ

→ usually want uniform coverage, eg 95% coverage for all θ

↳ can do for selection or prove mathematically for all.

Other topics & further reading:

Intervals from point estimates, intervals from tests

Linearisation & uncertainty propagation

Estimation vs testing? p-values & NHST?

Bayesian vs Frequentist wars?

Machine learning vs statistics?

Objective vs subjective Bayes?

Confidence intervals/sets for ill-posed inverse problems?

Simulation-based inference?

Causal inference?

Randomisation, experimental design, computer
experiment design?

→ Will (maybe) put some on canvas.

Can always ask me
directly for suggestions!
