

# ENGSCI 721

## INVERSE PROBLEMS

*Oliver Maclaren*  
*oliver.maclaren@auckland.ac.nz*

## MODULE OVERVIEW

Inverse Problems (*Oliver Maclaren*) [~8 lectures/2 tutorials]

### 1. Basic concepts [4 lectures]

Forward vs inverse problems. Well-posed vs ill-posed problems. Algebra and calculus of inverse problems (matrix calculus, 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. Preview of the statistical view of inverse problems

[1 lectures]

Bayesians, Frequentists and all that. Basic frequentist analysis.

## LECTURE 8: INTRODUCTION TO STATISTICAL APPROACHES TO INVERSE PROBLEMS

Topics:

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

## EngSci 721 Lecture 8.

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 applications

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

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

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  $\theta$  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

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

using explicit overbar/arrow  
to distinguish from scalar  $y$ , but will drop

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

$$\text{" } y = \vec{1}_m \theta \text{"} \quad \begin{array}{l} \text{(vector eqn of obs., have} \\ \text{dropped overbars)} \end{array}$$

ie

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

for  $\theta$  given  $y$  (vector of obs.).

Note: Overdetermined system ( $m$  eqns, one unknown)

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

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

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

$\rightarrow$

so:

$$\begin{aligned} A^+ &= (A^T A)^{-1} A^T = (I_m^T I_m^{-1}) I_m^T \\ &= (m)^{-1} I_m^T \\ &= \frac{1}{m} I_m^T \end{aligned}$$

& our (pseudo-inverse) sol<sup>n</sup> is

$$\begin{aligned} \theta^+ &= A^+ y = \frac{1}{m} I_m^T y \\ &= \frac{1}{m} \sum_i y_i \\ &= \text{mean}(y) \quad (\text{See Ru's material}) \\ &= \underbrace{\langle E_m(y) \rangle}_{\substack{\text{sample based expected} \\ \text{value.}}} \quad \begin{matrix} \swarrow \text{for more} \\ \text{on } \langle E \rangle \end{matrix} \end{aligned}$$

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

$$\begin{aligned} \theta &= \text{sample mean}(y) \\ &= \langle E_m(y) \rangle. \end{aligned}$$

→ The data resol. operator  $R_D = A A^+ = I_m \langle E_m \rangle$  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:

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

↑      ↑      ↑      ↗  
 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

$$\theta \rightarrow [F] \rightarrow y$$

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

### New story

$$\theta \rightarrow [F] \rightarrow y = \theta + e$$

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

### More generally

$$\theta \rightarrow [F] \rightarrow y \sim P_Y(y)$$

↗ '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:

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

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

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

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

(1B.) Sample version of expected value:

$$\theta = E[Y] \approx E_m[Y] = \frac{1}{m} \sum_i Y_i$$

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

→ 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  
(eg over-determined), we  
can solve in 'average sense'  
↳ match expected value.

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

↳ see eg 'An introduction to the bootstrap'  
by Efron & Tibshirani  
'Analogy estimation methods in  
econometrics' by Manski

→ 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 'minimally' ]

## Estimate?

We used a sample-based estimate of  $\theta$ :

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

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  $\sigma^2$

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

→ see Ru's part for (much) more!



## Normal/Gaussian distribution

→ Ru's favorite 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  $\theta$  & constant variance  $\sigma^2$  is:

$$P_Y(y; \theta) = \underbrace{\frac{1}{(2\pi\sigma^2)^{m/2}}}_{\text{'given } \theta\text{'}, when \theta \text{ isn't vec. a random var. some, esp. Bayesians just use } P(y|\theta) \text{ instead of } P(y;\theta)}} \exp\left\{-\frac{1}{2} \left\| \frac{y - \theta}{\sigma} \right\|^2\right\}$$

(assuming  $\sigma$  known here)

( $\theta$  scalar,  $1\theta$  vector)

## Likelihood

Given an observation/realisation  $y$  of  $\gamma$ ,

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.
- key is relative values eg

$$\frac{L(\theta_1; y)}{L(\theta_2; y)} = \frac{P(Y=y; \theta_1)}{P(Y=y; \theta_2)}$$

→ constants & 'dy' cancel.

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

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

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

$$P(Y=y; \theta) \neq P(\overbrace{\theta}^{\text{cf.}} | Y)$$

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

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

→ doesn't require  $\theta$  to be random

→ isn't additive  
 etc.

'Maximum likelihood' estimation determines

a point estimate by solving:

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

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

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

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

Now:

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

$$\Rightarrow \begin{cases} \hat{\theta}_{ML} = \arg \max L(\theta; y) \\ = \arg \max \lambda(\theta; y) \\ = \arg \min -\lambda(\theta; y) \end{cases}$$

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  $\sigma$ . ]

## 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  $\theta$ ?

Bayesians

- want a 'posterior probability'  $P(\theta|y)$  that represents their personal belief/'state of information' etc about the true value of  $\theta$  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.

Frequentists

- 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:

$$\hat{\theta} = A^+ y \quad \begin{array}{l} \text{↑ data realisation of random var} \\ \text{estimate of } \theta \quad \text{estimator} \end{array}$$

Where:

$$\begin{array}{c} A \leftarrow \text{forward model (for mean)} \\ \theta \leftarrow Y \leftarrow \text{data (or sample)} \\ \text{parameter} \quad A^+ \leftarrow \text{estimator} \end{array}$$

## Frequentist statistics as a two player game.

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

$$\begin{array}{c} \text{estimate} = T(Y) \\ \text{for } \theta \\ \uparrow \quad \uparrow \\ \text{value} \quad \text{method} \\ \text{data realisation} \end{array}$$

→ The estimate is random because  $Y$  (data) is random, even if  $\theta$  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  $\theta$ )
  - 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  $\theta$  we have good expected performance in repeated use,  
eg want:

$$\mathbb{E}[d(T(Y), \theta)] \text{ for } Y \sim P(Y; \theta) \text{ to be small.}$$

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

Interval estimator

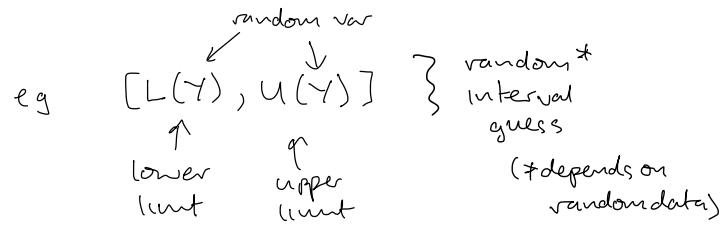
- Want to guarantee that for all of nature's choices of  $\theta$ , our (random) interval guesses will often (eg 95% of the time we get data) contain ('trap') the true but unknown value.

$C(Y)$  is 95% confidence interval procedure →  $\left\{ \begin{array}{l} \left| P(C(Y) \ni \theta; \theta) \geq 0.95, \text{ for } Y \sim P(Y; \theta) \right. \\ \left. \begin{array}{c} \uparrow \quad \uparrow \quad \uparrow \\ \text{random fixed.} \quad \text{interval} \quad \text{coverage (confidence level} \end{array} \right. \end{array} \right.$

since nature's choice!

## Confidence intervals?

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



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

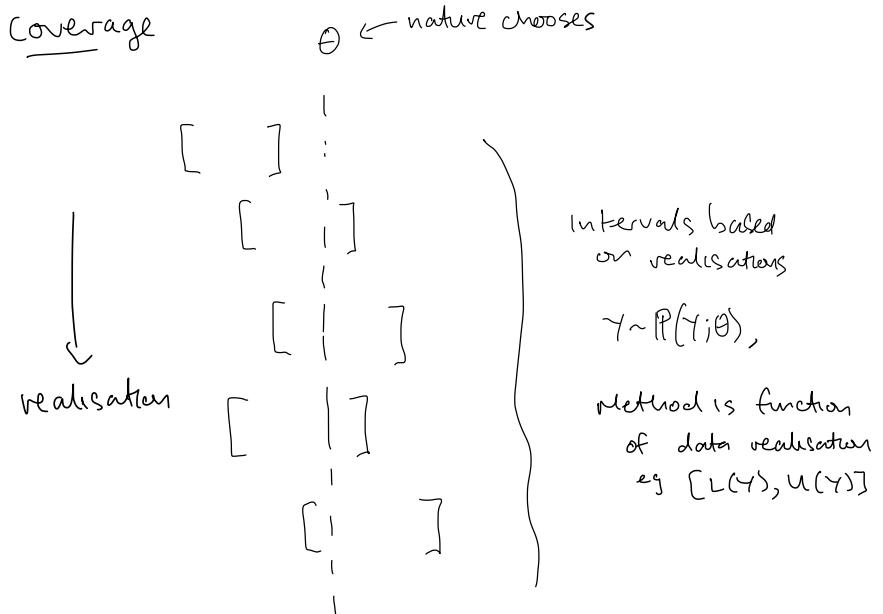
A given interval realisation either contains the true value or doesn't. Not 95% prob.

↪ 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



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  $\theta$ , generate data from  $P(Y; \theta)$ , see how often interval traps true value.
- get coverage for each choice of  $\theta$
- usually want uniform coverage, eg 95% coverage for all  $\theta$
- ↪ 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!

---