

EngSci 721

Inverse Problems and Learning From Data

Oliver Maclaren (oliver.maclaren@auckland.ac.nz)

1. Basic concepts [5 lectures + 1 Tutorial]

Forward vs inverse problems. Well-posed vs ill-posed problems. Algebra and calculus of inverse problems (left and right inverses, generalised and pseudo inverses, resolution operators, matrix calculus). Representing higher dimensional problems (image data etc).

2. Instability and regularisation in linear and nonlinear problems [6 lectures + 1 Tutorial]

Instability and related issues for generalised inverses. Introduction to regularisation and trade-offs. Tikhonov regularisation. Higher-order Tikhonov regularisation. Sparsity and regularisation using different norms. Truncated singular value decomposition. Iterative regularisation, including stochastic/mini-batch gradient descent.

3. Further topics [3 lectures + 1 Tutorial]

Regularisation parameter choice, including statistical and machine learning views of regularisation. Confidence sets for linear and nonlinear models. Physics-informed machine learning and neural networks.

Module overview

Inverse Problems and Learning From Data (*Oliver Maclaren*)

[~14 lectures/3 tutorials]

Lecture 5: Matrix Calculus

Topics:

- Motivation
- Three key rules for derivatives
- First principles
- Key rules for differentials
- Examples

Eng Sci 721 : Lecture 5: Matrix Calculus

- Motivation
- Three key rules: derivatives for scalar & vector-valued functions of a vector
- First principles: derivatives & differentials for functions between any vector spaces (e.g. matrices \xrightarrow{f} matrices)
- Other key rules: differentials
- Examples

Motivation:

Many inverse/learning from data / ML problems require derivatives (modern ML = 'differentiable programming') of scalar/vector/matrix functions of scalar/vector/matrix inputs

Recall

1. Tall / over-determined system: approximate solution

$$\rightarrow \text{Define } r = y - Ax \quad \left\{ \begin{array}{l} \text{residual 'error'} \\ \text{norm } \| \cdot \| \end{array} \right.$$

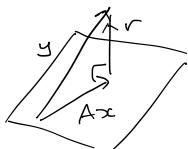
- Measure size of error with a norm $\| \cdot \|$ (see handout for diff. types)
- Typically assume $\| \cdot \|_2$

New problem:

$$\underset{x \in \mathbb{R}^n}{\text{minimise}} \quad \|y - Ax\|, \quad A \& y \text{ given}$$

- "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:

$$\underset{x}{\text{minimise}} \quad \|y - Ax\|_2^2 \quad (\text{equiv. problem})$$

First:

$\|y - Ax\|^2$ as a matrix-vector expression

$$As: \|z\|^2 = z^T z$$

we have:

$$\|y - Ax\|^2 = (y - Ax)^T (y - Ax)$$

$$= (y^T - x^T A^T)(y - Ax)$$

Since:

$$\left| \begin{array}{l} (A + B)^T = A^T + B^T \\ \& (AB)^T = B^T A^T \end{array} \right.$$

for matrices
 A, B

Also

$$\left| \begin{array}{l} (A + B)(C + D) \\ = AC + AD + BC + BD \end{array} \right.$$

$$So \quad (y^T - x^T A^T)(y - Ax)$$

$$= y^T y - y^T Ax - x^T A^T y + x^T A^T Ax$$

Also $y^T A x = \underline{\text{scalar}} \quad & \text{scalar}^T = \underline{\text{scalar}}$

$$\begin{aligned} & \& (y^T A x)^T = x^T A^T y \\ & & = y^T A x \end{aligned}$$

So :

$$\begin{aligned} & \|y - Ax\|^2 \\ & = y^T y - 2y^T A x + x^T A^T A x \end{aligned}$$

How to minimise wrt x ?

→ take derivative wrt x &

set = 0 (actually just a necessary condition but ignore...)

→ Need to be able to differentiate matrix expressions with respect to vectors etc

→ Also often want to differentiate eg matrix expressions wrt matrices! (or tensors etc ~)

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

- important for inverse prob. & machine learning etc but... requires:

(Matrix calculus, Tensor calculus) etc

→ multiple conventions/notation

(see e.g. Wikipedia page on Matrix calculus)

- I will give you (three key rules) you can use for many problems, instead of remembering the details!

- I will also show you how to derive these & more from first principles

- The approach is based on 'matrix calculus' rather than eg tensor calculus, tho is equiv.

↳ we work in 'flattened' (vec) linear world!

i.e. vec(M_{ij}) vs M_{ij} etc.

First: derivatives of scalar & vector-valued functions of a vector

$$\text{Eq: } f(\mathbf{x}) = \mathbf{x}^T A \mathbf{x}$$

↑
vector input

scalar output

Conventions

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

→ use 'Jacobian layout', e.g. vector f, \mathbf{x} ,
components:

$$[D_{\mathbf{x}} f]_{ij} = \frac{\partial f_i}{\partial x_j} = \begin{matrix} f_1 \\ \vdots \\ f_m \end{matrix} \left[\begin{matrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{matrix} \right]$$

Note, if f is scalar-valued, this implies
 $D_{\mathbf{x}} f$ is a row vector:

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

Hence this is sometimes written

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

to indicate the 'layout' direction is 'row like'
in \mathbf{x} .

We can define the gradient of f as

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

which again indicates how to 'lay out' results.

$$\boxed{\text{Note } \left(\frac{\partial f}{\partial x_i} \right)^T = \frac{\partial f^T}{\partial x_i}}$$

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

Derivatives

Constant

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

Linear

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

Quadratic

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

Know these !

We can also show e.g

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

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

o [Multivariable vector] chain rule:

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

$$\begin{aligned} &= (D_g f) \cdot (D_x g) + (D_h f) \cdot (D_x h) \\ &= \frac{\partial f}{\partial g^T} \frac{\partial g}{\partial x^T} + \frac{\partial f}{\partial h^T} \frac{\partial h}{\partial x^T} \end{aligned}$$

Imply e.g 'product rule' of form:

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

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

where do these come from??

Back to least squares!

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

→ soon!

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

→ first, consider \longrightarrow

$$(\text{or } \nabla_x f = 0 \dots)$$

3 Rules:

$$D_{x_i} (y^T y) = 0^T$$

$$D_{x_i} (-2y^T A x) = -2y^T A$$

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

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

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

$$\Rightarrow \boxed{A^T A x = A^T y} \quad \text{as before!}$$

Recall:

Least squares approximation

We have seen this leads to....

The normal equations:

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

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

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

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

First principles matrix calculus

→ Derivatives, linear approximation & differentials

This approach is heavily influenced by that of Magnus & Neudecker,
 ↗ A lot beyond scope! ↗

What is a derivative?

Recall the usual defⁿ for $f: \mathbb{R} \rightarrow \mathbb{R}$ at x_0 :

$$(1) \quad \boxed{\frac{df(x_0)}{dx} := \lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h}}$$

This is equivalent to, for some $a \in \mathbb{R}$

$$(2) \quad \boxed{f(x_0 + h) = f(x_0) + ah + o(h), h \rightarrow 0}$$

where ' $o(h)$ ' is 'little oh' & means, roughly, 'goes to zero faster than h ',

& we define: $\boxed{\frac{df(x_0)}{dx} := a \quad \rightarrow}$
from the above.

(a depends on x_0).

Linear approximations & differentials

The form:

$$f(x_0 + h) = f(x_0) + ah + o(h), h \rightarrow 0$$

emphasises that f can be approximated by the linear function ah at x_0 .

[note: ah should really be $a(x_0) \cdot h$ to allow different a at each point]

We further define the differential of f to be the linear part at

$$x_0 : \boxed{df = a \cdot h = \frac{df}{dx} \cdot h} \quad \text{via } \frac{df}{dx} := a \\ (\text{i.e., in full detail...})$$

$$\boxed{df(x_0; h) := a(x_0) \cdot h = \frac{df}{dx}(x_0)h}$$

We now allow this linear approximation to be evaluated for any 'increment'

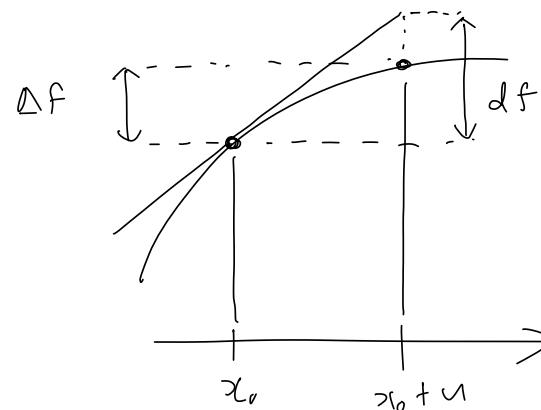
not just $h \rightarrow 0$, i.e.

$$\boxed{df(x_0; u) = a(x_0) \cdot u} \quad \text{for any } u.$$

eg finite

$$\boxed{\left| \frac{df(x_0)}{dx} \cdot u \right|} \quad \text{by defn.}$$

Geometric picture



$\Delta f = f(x_0 + u) - f(x_0)$ = actual change in f for increment u

$$df = a(x_0) \cdot u = \frac{df}{dx}(x_0) \cdot u \quad (\text{not nec. } u \rightarrow 0)$$

for $a(x_0)$ satisfying

$$\Delta f = a(x_0) \cdot h + o(h), \forall h \rightarrow 0$$

Generalisations to higher-dim

Both def'n's (1) & (2) allow us to generalise derivatives & differentials to functions having vector inputs & vector outputs

Here we mean 'vectors' in the general sense, so can apply to eg 'matrices as vectors'!]

These generalisations make different assumptions

(1) \rightarrow Gâteaux derivative & differential
↳ weak, exists more often

(2) \rightarrow Fréchet derivative & differential
↳ stronger, can have (1)
but not (2)

For us, both exist & are equal tho!

(1) Fréchet (linear approx., implicit)

$$\boxed{\text{if}} \quad f(x_0 + h) = f(x_0) + Ah + o(\|h\|)$$

where $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ $\forall h \in \mathbb{R}^n$

$$x_0, h \in \mathbb{R}^n$$

$$A \in \mathbb{R}^{m \times n} \text{ ie matrix}$$

[then] this defines (implicitly!)

- derivative: $\boxed{Df(x_0) = A}$

- differential: $\boxed{df(x_0, u) = Df(x_0)u}, u \in \mathbb{R}^n$

↳ also write $\boxed{df(x_0) = Df(x_0)dx}$ ie $u = dx$

or just $\boxed{df = Df dx}$



(2) Gâteaux | (directional derivative, explicit)

$$\boxed{\text{If}} \quad \lim_{t \rightarrow 0} \frac{f(x_0 + th) - f(x_0)}{t} = Ah, \quad \text{where}$$

where $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$

$t \in \mathbb{R}$ (scalar)

$x_0, h \in \mathbb{R}^n$

$A \in \mathbb{R}^{m \times n}$ is matrix

Then this defines (by 'direct calculation')

◦ derivative: $Df(x_0) = A$

◦ differential: $Df(x_0; u) = Df(x_0)u$, $u \in \mathbb{R}^n$

↪ write $Df(x_0) = Df(x_0)dx$ if $u = dx$

or just $Df = Df dx$

Subtleties

→ Gâteaux differentiability allows the def'n of partial derivatives:

$$\text{eg } df(x_0; e_j) = \lim_{t \rightarrow 0} \frac{f(x_0 + te_j) - f(x_0)}{t}$$

where $e_j = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} \leftarrow j$

= Ae_j (if exists)

= j^{th} col of A

If $f_i = i^{\text{th}}$ component of f ,

$$df_i(x_0; e_j) = \lim_{t \rightarrow 0} \frac{f_i(x_0 + te_j) - f_i(x_0)}{t}$$

= i^{th} row, j^{th} col
of A

$$= e_i Df(x_0; e_j)$$

→ defines $D_j f_i(x_0) = \frac{\partial f_i}{\partial x_j}(x_0)$ note order

& $Df(x_0) = [D_j f_i(x_0)]_{i=1:n, j=1:m}$
array of partial deriv.
= $\boxed{\text{Jacobian}}$

Subtleties

- We can have Gâteaux differentiability
 / existence of partial derivatives / Jacobian
 but no Fréchet derivative (linear approx.)
matrix
 - We will ignore these cases!
 - take 'differentiable' to mean
 'Fréchet exists' (& equals Gâteaux)

Leads to correspondence (identification thm)

If $d f(x_0; u) = A(x_0) u$ for some matrix A depending on x_0 , then $A(x_0) = \begin{bmatrix} D_j f_i(x_0) \end{bmatrix} = \underline{\underline{Df(x_0)}}$

(Derivative is uniquely determined

Via linearisation or partial

derivative / Jacobian matrix)

le

$\rightarrow \cancel{\boxed{\text{If } df = A dx \text{ then } Df = A}}$ \leftarrow
 (& vice-versa)

Derivatives & differentials of matrix-valued functions of matrices

Recall that matrices can be directly thought of as 'abstract' vectors in the vector-space $\mathbb{R}^{m \times n}$

$$\text{e.g. } M_3 = aM_1 + bM_2 \text{ etc}$$

→ these are not column vectors
in the usual sense

→ However, we can also explicitly map these to column vectors & use this familiar space for manipulations!

$$M \mapsto \text{vec}(M)$$

'vec' is linear so:

$$\overline{\text{vec}(\alpha M_1 + \beta M_2)} = \alpha \text{vec}(M_1) + \beta \text{vec}(M_2)$$

etc.



Implications for derivatives

$$\text{If } f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{m \times n}$$

is differentiable in the Frechet sense, we have a linear approx. in the vector space of matrices:

$$df(M_0, dM) = L(M_0) dM$$

where L is a 'linear operator' mapping matrices to matrices.

Importantly, if M is $m \times n$

then $[dM \text{ is also } m \times n]$

However $[L \text{ is not a matrix!}]$

→ it is a 'tensor', here represented by a 4D array of dimensions $m \times n \times m \times n$



Tensors in brief. Much more to it!

→ Represented by multidimensional arrays ($0, 1, 2, 3$ etc dimensions)
Aijkl...

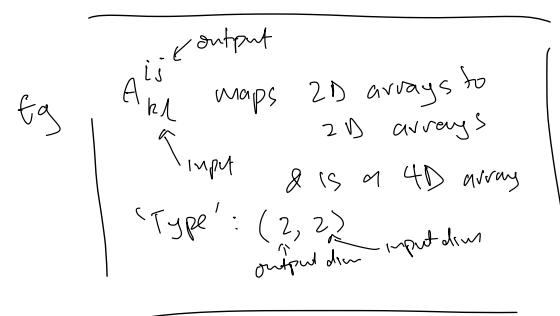
→ Group 'axes' into two sets:
inputs & outputs

e.g. ijkkl → $(i, j | k, l)$ or $A_{k|l}^{ij}$

→ write or just Aijkl if clear

→ Represent linear mappings

between tensors of dimensions of input & dimensions of output



→ we can re-group indices to change input-output specification:

A_{ij}^i : maps vector x^j to vector y^i

→ can also 'input from other side'
e.g. $x_i \rightarrow A_{ij}^i \rightarrow y_j$

Derivatives as tensors

So... if $f: \underbrace{\mathbb{R}^{m \times m}}_{\text{matrices}} \rightarrow \underbrace{\mathbb{R}^{m \times n}}_{\text{matrices}}$

Then $df = \underbrace{\mathcal{L}}_{\substack{\text{is a } 4D \\ \text{tensor!}}} \underbrace{d\mathbf{x}}_{\substack{\text{components} \quad \mathcal{L}_{ikl}^{ij}}} \quad$

Formula for 'applying' \mathcal{L} to B to get A :

$$\begin{aligned} A^{ij} &= \sum_{k,l} \mathcal{L}^{ij}_{kl} B^{kl} \\ &= \mathcal{L}^{ij}_{kl} B^{kl} \quad (\text{implicit sum: Einstein convention}) \end{aligned}$$

→ OK to do this... but
can also 'flatten'!



Derivative matrices (cf tensor)

We can also define the derivative matrix
eg a matrix function of a matrix
variable

→ we simply map matrices to col. vectors
using vec!

- Derivative is linear approx. to } matrix
vector → vector mapping } -

→ Define for $F(\mathbf{x}): \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{p \times q}$

$$\boxed{\begin{aligned} d\text{vec}(F(\mathbf{x})) &= A(\mathbf{x}) d\text{vec}(\mathbf{x}) \\ \Leftrightarrow D_{\mathbf{x}} F(\mathbf{x}) &= A(\mathbf{x}) \\ \Leftrightarrow \frac{\partial \text{vec} F(\mathbf{x})}{\partial (\text{vec} \mathbf{x})^T} &= A(\mathbf{x}) \end{aligned}}$$

i.e. $D_{\mathbf{x}} F$ represents the derivative matrix
(vectorised)

Differentials vs derivatives

In the original space, dM is same shape as M , while ' DM ' becomes a higher-order tensor unless vectorised

→ It is often easier to do some calculations with differentials in the original space, & then, if needed, vectorise to get a derivative matrix

This is helped by the following useful rules that apply to differentials directly in their natural space (in contrast, derivative expressions can become complex)

Key rules of differentials (original space)

For matrix functions F, G , constant matrix A , constant number α ,

- $dA = 0$
- $d(\alpha A) = \alpha dA$
- $d(A+B) = dA + dB$
- $d(A-B) = dA - dB$
- $d(uv) = (du)v + u(dv)$
- $d(U^T) = (dU)^T$
- $d(\text{vec}(u)) = \text{vec}(du)$
- $d(\text{linear function}(u)) = \text{linear } f(du)$
- Also, if $H = G(F(x))$
 $dH(x_0; dx) = dG(F(x_0); dF(x_0; dx))$
(chain rule / Cauchy invariance)

Examples

- Consider $f(x) = x^T A x$

→ quadratic (nonlinear)

$$f: \underbrace{\mathbb{R}^n}_{\text{vector}} \rightarrow \underbrace{\mathbb{R}}_{\text{scalar}} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{already 'vecteurised'}$$

$$df = df(x; dx)$$

$$= (dx)^T A x + x^T A dx$$

$$= x^T A^T dx + x^T A dx$$

$$= x^T (A^T + A) dx$$

$$\Rightarrow Df(x) = x^T (A^T + A) \quad \checkmark$$

Examples (harder!)

- Suppose A is invertible. $\underline{A^{-1}} = f(A)$

$$\text{Then } \underline{A^{-1} A} = \underline{I}$$

$$d(A^{-1} A) = dI = 0$$

$$\Rightarrow d(A^{-1}) A + A^{-1} dA = 0$$

$$d(A^{-1}) A A^{-1} + A^{-1} dA A^{-1} = 0$$

$$\Rightarrow \boxed{d(A^{-1}) = -A^{-1} dA A^{-1}}$$

$$\text{vec}(d(A^{-1})) = -(A^T)^{-1} \otimes A^{-1} \text{vec}(dA)$$

$$\Rightarrow \boxed{\frac{D A^{-1}}{A} = - (A^T)^{-1} \otimes A^{-1} \quad !}$$

✓ (Magnus & Neudecker)

matrix-val.
f of matrix.

examples (harder!)

$x \in \mathbb{R}^n$

- Consider $F(x) = \underbrace{x x^T}$

$F: \text{vector} \rightarrow \text{matrix}$

$$dF = (dx)x^T + x d(x^T)$$

$$= (dx)x^T + x (dx)^T$$

$$\begin{aligned}\text{vec}(dF) &= \text{vec}(I_n(dx)x^T + x(dx)^T I_n) \\ &= (x \otimes I)\text{vec}(dx) + (I \otimes x)\text{vec}(dx) \\ &= (x \otimes I)dx + (I \otimes x)dx \\ &= \left[x \otimes I + I \otimes x \right] \overbrace{\text{vec}(dx)}^{dx}\end{aligned}$$

$$\Rightarrow D_x F = x \otimes I + I \otimes x$$

✓ (see Magnus &
Neudecker)

Higher derivatives: second order

suppose $\phi: \mathbb{R}^n \rightarrow \mathbb{R}$ (scalar-valued)

- Differentials:

$$\rightarrow \left| \overline{d^2\phi := d(d\phi)} \right| \leftarrow$$

- Derivatives:

Define Hessian matrix as

$$\left| \begin{array}{c} \overline{H\phi = \nabla(D\phi)^T} \\ \overline{| = \frac{\partial^2\phi}{\partial x_i \partial x_j}} \end{array} \right|$$

$$\text{ie } \left| \overline{[H\phi]_{ij} = D_{ij}\phi} \right|$$

- Relation (identification):

$$\text{if } \left| \overline{d^2\phi = (dx)^T B dx} \right| \text{ for some } B,$$

$$\text{Then } \left| \overline{H\phi = \frac{B + B^T}{2}} \right| \left\{ \begin{array}{l} \text{note!} \\ \text{related to symmetry of Hessian} \end{array} \right.$$

& conversely.

Exercises

- Re-derive the 'normal equations' for least squares from memory
 - Derive the 3 key rules for derivatives from rules for differentials
 - If $\phi = \mathbf{x}^T \mathbf{A} \mathbf{x}$
find $d(d\phi)$
 - Look up automatic differentiation
in ML libraries
 - Tensorflow
 - Pytorch
 - Jax
 - ⋮
etc !
-