

Machine Learning

Section 9: Matrix Differential Calculus

第9节：矩阵微分计算

Stefan Harmeling

8. November 2021

Linear regression summary

Model assumption

$$p(w) = \mathcal{N}(w|0, \tau^2 I) \quad \text{prior}$$

$$p(y|X, w) = \mathcal{N}(y|Xw, \sigma^2 I) \quad \text{likelihood}$$

1. Ordinary least squares (MLE):

$$w_{\text{MLE}} = (X^T X)^{-1} X^T y \quad \text{maximizer of likelihood, no prior}$$

2. Ridge regression (MAP):

$$p(w|X, y) = \mathcal{N}(w|w_n, V_n) \quad \text{posterior 后方}$$

$$w_{\text{ridge}} = (\lambda I + X^T X)^{-1} X^T y \quad \text{maximizer of posterior}$$

3. Bayesian linear regression:

$$p(w|X, y) = \mathcal{N}(w|w_n, V_n) \quad \text{posterior}$$

$$w_n = (\lambda I + X^T X)^{-1} X^T y \quad \text{posterior mean}$$

$$V_n = \sigma^2 (\lambda I + X^T X)^{-1} \quad \text{posterior covariance}$$

Today:

**How to calculate derivatives of
scalar-/vector-/matrix-valued functions of
scalar-/vector-/matrix-valued variables!**

如何计算
标量-/向量-/矩阵值变量的标量-/向量-/矩阵
值函数!

Matrix differential calculus

A tool to find complicated derivatives involving vectors and matrices.

Sources

- ▶ Download bn142.pdf from public sciebo folder of this lecture
- ▶ Magnus/Neudecker: Matrix differential calculus, 2007, PDF was available at <http://www.janmagnus.nl/misc/mdc2007-3rdedition.pdf>, the link doesn't work anymore, it looks like a new version is planned so the old PDF is removed. Instead much of the stuff can be found in a paper: <http://www.janmagnus.nl/papers/JRM012.pdf>
- ▶ <http://www.janmagnus.nl/misc/mdc-ch18.pdf>
- ▶ Lütkepohl: Handbook of matrices, 1996

What is a differential?

Definition 9.1 (source: https://en.wikipedia.org/wiki/Differential_of_a_function)

A differential is an infinitesimal change in some varying quantity, e.g.

- ▶ *for variable x the differential is denoted as dx*
- ▶ *for variable $y = f(x)$ being the image of x under function f , the differential is*

$$dy = f'(x) dx$$

where $f'(x)$ is the derivative of f .

Note

- ▶ writing the derivative $f'(x) = dy/dx$ in Leibniz notation:

$$dy = \frac{dy}{dx} dx$$

(looks reasonable!)

- ▶ an integral is a summation over infinitesimal rectangles

$$\int f(x) dx = \sum_x f(x) dx$$

Differential calculus微分计算

Notation for derivative (not for the differential) 导数的符号

$$Df(x) = f'(x)$$

Identification formula: (relates derivative and differential) 识别公式

$$df = Df(x)dx$$

How to find a derivative with the differential calculus:

如何用微积分找导数。

1. write the letter d in front of the expression. 在表达式前面写上字母d
2. identity the constants and variables. 确定常数和变量的身份
3. transform the expression 转换表达式
4. read off the derivative using the identification table 使用识别表读出导数

Notation!

For the rest of the slides:

- ▶ **Scalars**: small greek letters: $\phi, \psi, \alpha, \dots$ 疤痕
- ▶ **Vectors**: small latin letters: u, v, x, f, a, \dots 媒介物
- ▶ **Matrices**: capital latin letters: U, V, F, A, \dots

Rules for the calculus (1)

Only with scalars:

$$d\alpha = 0$$

$$d(\alpha\phi) = \alpha d\phi$$

$$d(\phi + \psi) = d\phi + d\psi$$

$$d(\phi\psi) = (d\phi)\psi + \phi d\psi$$

$$d(\phi/\psi) = ((d\phi)\psi - \phi d\psi) / \psi^2$$

where α is constant.

Example:

Find the derivative of $\phi(\xi) = \xi^2$.

$$d\phi = d\xi^2 = 2\xi d\xi$$

Thus using $d\phi = D\phi(\xi)d\xi$ we read off the derivative:

$$D\phi(\xi) = 2\xi$$

Rules for the calculus (2)

With scalars and vectors:

$$da = 0_n$$

$$d(\alpha u) = \alpha du$$

$$d(u + v) = du + dv$$

$$d(u^T v) = (du)^T v + u^T dv$$

$$d(u^T) = (du)^T$$

where α and a are constant.

Example:

Find the derivative of $\phi(x) = x^T x$.

$$d\phi = \dots$$

Rules for the calculus (3)

With scalars, vectors and matrices:

$$dA = 0_{m \times n}$$

$$d(\alpha U) = \alpha dU$$

$$d(U + V) = dU + dV$$

$$d(UV) = (dU)V + UdV$$

$$d(U^T) = (dU)^T$$

$$d \operatorname{tr} U = \operatorname{tr} dU$$

where α and A are constant.

Example:

Find the derivative of $\phi(X) = \operatorname{tr}(X^T X)$.

$$d\phi = d \operatorname{tr}(X^T X) = \operatorname{tr}((dX)^T X + X^T dX) = \operatorname{tr}(2X^T dX)$$

Thus

$$D\phi(X) = 2(\operatorname{vec} X)^T$$

Rules for the calculus (4)

Many more rules:

$$d\alpha = 0$$

$$d(\alpha\phi) = \alpha d\phi$$

$$d(\phi + \psi) = d\phi + d\psi$$

$$d(\phi\psi) = (d\phi)\psi + \phi d\psi$$

$$d(\phi/\psi) = ((d\phi)\psi - \phi d\psi) / \psi^2$$

$$da = 0_n$$

$$d(\alpha u) = \alpha du$$

$$d(u + v) = du + dv$$

$$d(u^T v) = (du)^T v + u^T dv$$

$$d(u^T) = (du)^T$$

$$dA = 0_{m \times n}$$

$$d(\alpha U) = \alpha dU$$

$$d(U + V) = dU + dV$$

$$d(UV) = (dU)V + UdV$$

$$d(U^T) = (dU)^T$$

$$d \operatorname{vec} U = \operatorname{vec} dU$$

$$d(U \otimes V) = (dU) \otimes V + U \otimes dV$$

$$d(\phi^\alpha) = \alpha \phi^{\alpha-1} d\phi$$

$$d \det U = \det(U) \operatorname{tr}(U^{-1} dU)$$

$$d \exp \phi = \exp(\phi) d\phi$$

$$d \operatorname{tr} U = \operatorname{tr} dU$$

$$d(U \odot V) = (dU) \odot V + U \odot dV$$

$$d(U^{-1}) = -U^{-1} (dU) U^{-1}$$

$$d \log(\det U) = \operatorname{tr}(U^{-1} dU)$$

$$\operatorname{tr}(d \exp U) = \operatorname{tr}(\exp(U) dU)$$

- ▶ α, a, A be constants
- ▶ $\phi, \psi, u, v, x, f, U, V, F$ be variables/functions.

Identification table (more identification formulas)

	function	differential	derivative	shape of derivative
$\phi(\xi)$	$\mathbb{R} \rightarrow \mathbb{R}$	$d\phi = \alpha(\xi)d\xi$	$D\phi(\xi) = \alpha(\xi)$	1×1
$\phi(x)$	$\mathbb{R}^n \rightarrow \mathbb{R}$	$d\phi = a(x)^T dx$	$D\phi(x) = a(x)^T$	$1 \times n$
$\phi(X)$	$\mathbb{R}^{n \times q} \rightarrow \mathbb{R}$	$d\phi = \text{tr}(A(X)^T dX)$	$D\phi(X) = (\text{vec } A(X))^T$	$1 \times nq$
$f(\xi)$	$\mathbb{R} \rightarrow \mathbb{R}^m$	$df = a(\xi) d\xi$	$Df(\xi) = a(\xi)$	$m \times 1$
$f(x)$	$\mathbb{R}^n \rightarrow \mathbb{R}^m$	$df = A(x) dx$	$Df(x) = A(x)$	$m \times n$
$f(X)$	$\mathbb{R}^{n \times q} \rightarrow \mathbb{R}^m$	$df = A(X) d\text{vec } X$	$Df(X) = A(X)$	$m \times nq$
$F(\xi)$	$\mathbb{R} \rightarrow \mathbb{R}^{m \times p}$	$dF = A(\xi) d\xi$	$DF(\xi) = \text{vec } A(\xi)$	$mp \times 1$
$F(x)$	$\mathbb{R}^n \rightarrow \mathbb{R}^{m \times p}$	$d\text{vec } F = A(x) dx$	$DF(x) = A(x)$	$mp \times n$
$F(X)$	$\mathbb{R}^{n \times q} \rightarrow \mathbb{R}^{m \times p}$	$d\text{vec } F = A(X) d\text{vec } X$	$DF(X) = A(X)$	$mp \times nq$

Note: The differential has always the same shape as the function.

Many more matrix tricks

$$\text{tr}(AB) = \text{tr}(BA)$$

$$\text{diag}(UV^T) = (U \odot V)1_n$$

$$A \otimes 1_l = (I_m \otimes 1_l)A$$

$$\text{Diag } a = a1_n^T \odot I_n$$

$$\text{Diag}(\text{diag } A) = I_n \odot A$$

$$\text{vec}(a) = \text{vec}(a^T) = a$$

$$\text{tr}(u^T v) = \text{tr}(v^T u) = v^T u$$

$$\det(\exp A) = \exp(\text{tr } A)$$

$$\text{tr } A = 1_m^T (A \odot I_{m \times n}) 1_n = 1^T \text{diag}(A) = \text{tr } A^T$$

$$\text{tr}(U^T (V \odot C)) = \text{tr}((U^T \odot V^T)C)$$

$$1_l \otimes A = (1_l \otimes I_m)A$$

$$\text{diag } A = (A \odot I_n)1_n$$

$$\|U\|_{\text{Fro}}^2 = \text{tr}(U^T U) = \text{vec}(U)^T \text{vec}(U)$$

$$\text{vec}(ABC) = (C^T \otimes A) \text{vec}(B)$$

$$ABc = (c^T \otimes A) \text{vec } B = (A \otimes c^T) \text{vec } B^T = \text{vec}(c^T B^T A^T)$$

$$\text{vec } ab^T = b \otimes a$$

Notation

$0_n, 0_{m \times n}$

$1_n, 1_{m \times n}$

$I_n, I_{m \times n}$

$\text{vec } A$

$\text{diag } A$

$\text{Diag } a$

$\exp \alpha, \exp A$

$A \odot B$

$A \oslash B$

$A \otimes B$

n vector of zeros, $m \times n$ matrix of zeros

n vector of ones, $m \times n$ matrix of ones

$n \times n$ identity matrix, $m \times n$ identity matrix

vector containing the stacked columns of A

vector containing the diagonal of A

diagonal matrix with a along the diagonal

scalar exponential function, matrix exponential function

Hadamard product (component-wise product)

component-wise division

Kronecker product

Interlude: finite differencing (Python/Numpy version)

You should always check your derivatives with finite differencing which is an alternative way to calculate a derivative! Here is some Python/Numpy code:你应该经常用有限差分法来检查你的导数，这是计算导数的另一种方法! 这里有一些Python/Numpy代码。

```
import numpy as np
def finite_diff(f, x, delta):

    """estimate the gradient by finite-differencing method"""
    grad_f, dx = np.zeros_like(x), np.zeros_like(x)
    for i in range(x.size):
        dx.flat[i]      = delta
        grad_f.flat[i] = f(x+dx) - f(x-dx)
        dx.flat[i]      = 0.0
    return grad_f / (2*delta)
```

Interlude: finite differencing (MATLAB version)

You should always check your derivatives with finite differencing which is an alternative way to calculate a derivative! Here is some matlab **code**:你应该经常用有限差分法来检查你的导数，这是计算导数的另一种方法。是计算导数的另一种方法!下面是一些matlab

代码。

```
function df = finitediff(fun, x, d, varargin)
% FINITEDIFF estimates a gradient by finite-differencing method.
% (c) Stefan Harmeling, 2012-07-10.
sx = size(x);
nx = numel(x);
df = zeros(sx);
dx = zeros(sx);
for i = 1:nx
    dx(i) = d;
    df(i) = (fun(x+dx, varargin{:})-fun(x-dx, varargin{:}))/(2*d);
    dx(i) = 0;
end
```

Matrix differential calculus 矩阵微分计算

General recipe and examples

1. write the letter d in front of the expression
2. identity the constants and variables
3. transform the expression
4. read off the derivative using the identification table

Some pros and cons

- + clean notation
- + vectorized function leads to vectorized derivative (good for coding)
- + powerful: The book of Magnus/Neudecker shows how to take the derivative of eigenvalues and eigenvectors
- complicated formulas
- requires tricks and practice

More examples

Example: least squares

Find the derivative of $\phi(x) = (y - Ax)^2$.

$$d\phi = d((y - Ax)^T(y - Ax)) = -2(y - Ax)^T A dx$$

Thus:

$$D\phi(x) = -2(y - Ax)^T A = (-2A^T(y - Ax))^T$$

Example: vector-valued function of matrix

Find the derivative of $f(X) = (X^T X)^{-1} X^T y$.

We write $A = (X^T X)^{-1}$.

$$\begin{aligned} df &= d(X^T X)^{-1} X^T y \\ &= -A((dX)^T X + X^T(dX))AX^T y \\ &= -A(dX)^T XAX^T y - AX^T(dX)AX^T y - A(dX)^T y \\ &= -(A \otimes (XAX^T y)^T) \operatorname{dvec} X - ((XAX^T y)^T \otimes A) \operatorname{dvec} X - (A \otimes y^T) \operatorname{dvec} X \\ &= \underbrace{-(A \otimes (XAX^T y)^T + (XAX^T y)^T \otimes A + A \otimes y^T)}_{Df(X)} \operatorname{dvec} X \end{aligned}$$

Example: scalar-valued function of matrix 例子：矩阵的标量值函数

Often it is easier to find the differential of a scalar function $\phi(X) = c^T (X^T X)^{-1} X^T y$ where we again write $A = (X^T X)^{-1}$.

$$\begin{aligned} d\phi &= -c^T A(dX)^T X A X^T y - c^T A X^T (dX) A X^T y + c^T A (dX)^T y \\ &= -\text{tr}(c^T A (dX)^T X A X^T y) - \text{tr}(c^T A X^T (dX) A X^T y) + \text{tr}(c^T A (dX)^T y) \\ &= -\text{tr}(y^T X A^T X^T (dX) A^T c) - \text{tr}(c^T A X^T (dX) A X^T y) + \text{tr}(y^T (dX) A^T c) \\ &= -\text{tr}(A^T c y^T X A^T X^T dX) - \text{tr}(A X^T y c^T A X^T dX) + \text{tr}(A^T c y^T dX) \\ &= \text{tr}((-A^T c y^T X A^T X^T - A X^T y c^T A X^T + A^T c y^T) dX) \\ &= \text{tr}(\underbrace{(-A c y^T X A X^T - A X^T y c^T A X^T + A c y^T)}_{\text{call this } C(X)^T} dX) \end{aligned}$$

Using the identification formula: $D\phi(X) = (\text{vec } C(X))^T$ with shape $1 \times nq$ we get:

$$\begin{aligned} D\phi(X) &= (\text{vec}(-A c y^T X A X^T - A X^T y c^T A X^T + A c y^T))^T \\ &= (\text{vec}(-X A X^T y c^T A - X A c y^T X A + y c^T A))^T \end{aligned}$$

Example: use indices or not 例如：是否使用指数

Always avoiding indices can be painful:总是回避指数会让人感到痛苦。

$$\begin{aligned} d \operatorname{tr}(A \operatorname{Diag} v) &= d \operatorname{tr}(A(I_n \odot v \mathbf{1}_n^T)) = d \operatorname{tr}((A \odot I_n) v \mathbf{1}_n^T) \\ &= d \mathbf{1}_n^T \operatorname{Diag}(A) v = d \operatorname{diag}(A)^T v = \operatorname{diag}(A)^T dv \end{aligned}$$

Thus better rewrite with indices in this case:因此，在这种情况下，最好用指数重写。

$$d \operatorname{tr}(A \operatorname{Diag} v) = d \sum_i A_{ii} v_i = d \operatorname{tr} \operatorname{diag}(A)^T v = \operatorname{tr} \operatorname{diag}(A)^T dv$$

Example: Rayleigh coefficient

Find the derivative of Rayleigh coefficient $\phi(x) = x^T A x / (x^T x)$ for symmetric A . 求对称 A 的瑞利系数 $\phi(x) = x^T A x / (x^T x)$ 的导数

$$\begin{aligned} d\phi &= \frac{2x^T A(dx)(x^T x) - 2x^T A x x^T dx}{(x^T x)^2} \\ &= \frac{2(x^T x)x^T A(dx) - 2x^T A x x^T dx}{(x^T x)^2} \\ &= \frac{2(x^T x)x^T A - 2x^T A x x^T}{(x^T x)^2} dx \\ &= \frac{2}{(x^T x)^2} x^T (x x^T A - A x x^T) dx \end{aligned}$$

Thus the derivative is:

$$\begin{aligned} D\phi(x) &= \left(\frac{2}{(x^T x)^2} (A x x^T - x x^T A) x \right)^T \\ &= \left(2 \frac{A x}{x^T x} - 2 \frac{x^T A x}{(x^T x)^2} x \right)^T \end{aligned}$$

Example: derivative of steepest descent例子：最陡峭下降的导数

steepest descent:

$$x^{(k+1)} = x^{(k)} - \xi A^T (y - Ax^{(k)})$$

We can derive the following differentials:我们可以推导出以下差值。

- Find the derivative of $x^{(k)}(x^{(0)})$.

$$dx^{(k)} = (I_n + \xi A^T A)^k dx^{(0)}$$

- Find the derivative of $x^{(k+1)}(\xi)$.

$$dx^{(k+1)} = - \left(\sum_{i=0}^k (\xi A^T A + I_n)^i A^T (y - Ax^{(k-i)}) \right) d\xi$$

- Find the derivative of $x^{(1)}(A)$.

$$dx^{(1)} = -\xi (I_n \otimes (y - Ax^{(0)})^T + (x^{(0)})^T \otimes A^T) \text{dvec } A$$

solution see bn142.pdf

End of Section 09