

Nonlinear Optimization

(Com S 477/577 Notes)

Yan-Bin Jia

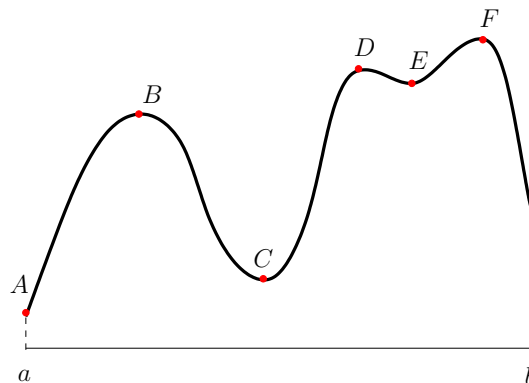
Oct 29, 2020

1 Introduction

Given a single function f that depends on one or more independent variable, we want to find the values of those variables where f is maximized or minimized. Often the computational cost is dominated by the cost of evaluating f (and also perhaps its partial derivatives with respect to all variables).

Finding a global extremum is, in general, a very difficult problem. Two standard heuristics are widely used: i) find local extrema starting from widely varying values of the independent variables, and then pick the most extreme of these; ii) perturb a local extremum by taking a finite amplitude step away from it, and then see if your routine can get to a better point, or “always” to the same one. Recently, “simulated annealing” methods have demonstrated important successes on a variety of global optimization problems.

The diagram below describes a function in an interval $[a, b]$. The derivative vanishes at the points B, C, D, E, F . The points B and D are local but not global maxima. The points C and E are local but not global minima. The global maximum occurs at F where the derivative vanishes. The global minimum occurs at the left endpoint A of the interval so that the derivative need not vanish.



Nonlinear optimization, also referred to as *nonlinear programming* when nonlinear constraints are involved, have many applications. Just to name a few: minimization of energy/power consumption, optimal control, robot path planning in potential fields, thermal and fluid model calibration.

Recall how the optimization process works in one dimension when f is a function from \mathbb{R} to \mathbb{R} . First we might compute the critical set of f ,

$$c_f = \{x \mid f'(x) = 0\}$$

By examining this set we can determine those x that are global minima or maxima. Notice that the computation of c_f seems to entail finding the zeros of the derivative f' . In other words, we have reduced the optimization problem to the root-finding problem. So why do we need to study nonlinear programming? Well, in higher dimensions, it is often easier to find (local) minimum than one would think. Intuitively, this is because f' is not an arbitrary function but rather a derivative whose integral (f) is given. We will thus have a lot more to say about optimization in higher dimension than we did about root finding.

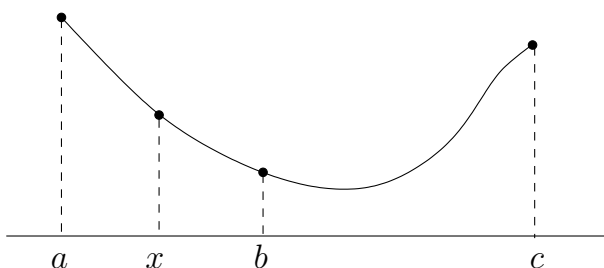
Since a maximization problem can be turned into a minimization problem simply by negating the objective function, we will deal with minimization only from now on.

2 Golden Section Search

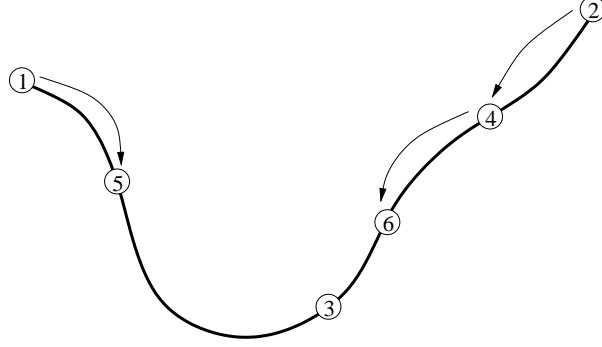
This method bears some analogy to bisection used in root finding. The basic idea is to bracket a minimum and then shrink the bracket. One problem is that we do not know the value of the function at the minimum, so we cannot know whether we have bracketed a minimum. Fortunately, if we are willing to settle for a local minimum, then we really just want to bracket a zero of f' . That we can do, by a numerical approximation to the derivative of f . The inner part of the loop works as follows. We start with three points a, b, c with $a < b < c$ such that

$$f(b) < \min\{f(a), f(c)\}.$$

Now choose a point x , say, half way between a and b . If $f(x) > f(b)$, as shown in the figure below, then the new bracketing triple becomes $[x, b, c]$. If $f(x) < f(b)$, then the new bracketing triple becomes $[a, x, b]$. To ensure that the interval $[a, c]$ will shrink toward a point, *one should alternate the bracketing tuples*, at least after a few rounds. For instance, halve $[a, b]$ in this round while $[b, c]$ in the next round.



The figure on the next page shows a more complete example about how golden section works. The minimum is originally bracketed by points 1, 3, 2. The function is evaluated at 4, which replaces 2; then at 5, which replaces 1; and then at 6, which replaces 4. Note that the center point is always lower than the two outside points. The minimum is bracketed by points 5, 3, 6 after the three steps.



Golden section search is applicable to extremizing functions of one variable only, just like bisection is to finding the roots of functions of this type only. Now we consider functions of more than one variables. Let the function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be twice continuous differentiable, that is, $f \in C^2$. A point $\mathbf{x}^* \in \mathbb{R}^n$ is said to be a *relative minimum point* or a *local minimum point* if there is an $\epsilon > 0$ such that $f(\mathbf{x}) \geq f(\mathbf{x}^*)$ for all $\mathbf{x} \in \mathbb{R}^n$ with $\|\mathbf{x} - \mathbf{x}^*\| < \epsilon$. If $f(\mathbf{x}) > f(\mathbf{x}^*)$ for all $\mathbf{x} \neq \mathbf{x}^*$ with $\|\mathbf{x} - \mathbf{x}^*\| < \epsilon$, then \mathbf{x}^* is said to be a *strict relative minimum point* of f .

A point \mathbf{x}^* is said to be a *global minimum point* of f if $f(\mathbf{x}) \geq f(\mathbf{x}^*)$ for all \mathbf{x} . It is said to be a *strict global minimum point* if $f(\mathbf{x}) > f(\mathbf{x}^*)$ for all $\mathbf{x} \neq \mathbf{x}^*$.

Let $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$. Recall that the *gradient* of f is a vector

$$\nabla f(\mathbf{x}) = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right). \quad (1)$$

It gives the direction in which the value of f increases the fastest. The *Hessian* H of f is defined as an $n \times n$ matrix:

$$H(\mathbf{x}) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix}$$

The Taylor series at $\mathbf{x}^* = (x_1^*, x_2^*, \dots, x_n^*)^T$ has the form

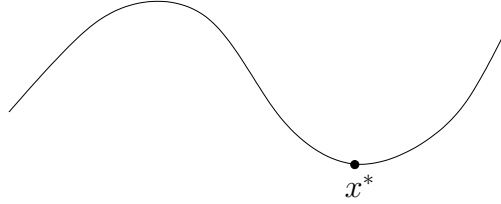
$$\begin{aligned} f(\mathbf{x}) &= f(\mathbf{x}^*) + \left(\frac{\partial f}{\partial x_1}(\mathbf{x}^*)(x_1 - x_1^*) + \cdots + \frac{\partial f}{\partial x_n}(\mathbf{x}^*)(x_n - x_n^*) \right) \\ &\quad + \frac{1}{2} \left(\frac{\partial^2 f}{\partial x_1^2}(\mathbf{x}^*)(x_1 - x_1^*)^2 + \frac{\partial^2 f}{\partial x_1 \partial x_2}(\mathbf{x}^*)(x_1 - x_1^*)(x_2 - x_2^*) + \cdots + \frac{\partial^2 f}{\partial x_n^2}(\mathbf{x}^*)(x_n - x_n^*)^2 \right) \\ &\quad + \cdots \\ &= f(\mathbf{x}^*) + \nabla f(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) + \frac{1}{2}(\mathbf{x} - \mathbf{x}^*)^T H(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) + O(\|\mathbf{x} - \mathbf{x}^*\|^3). \end{aligned} \quad (2)$$

If \mathbf{x}^* is a relative minimum, then the following conditions hold:

- i) $\nabla f(\mathbf{x}^*) = \mathbf{0}$,
- ii) $\mathbf{d}^T H(\mathbf{x}^*) \mathbf{d} \geq 0$ for every $\mathbf{d} \in \mathbb{R}^n$.

In one dimension, the above necessary conditions are familiar to us:

$$f'(x^*) = 0 \quad \text{and} \quad f''(x^*) \geq 0.$$



The Hessian $H(\mathbf{x}^*)$ at a relative minimum \mathbf{x}^* is symmetric *positive semi-definite*, that is, $\mathbf{x}^T H(\mathbf{x}^*) \mathbf{x} \geq 0$ for any \mathbf{x} . If \mathbf{x}^* is a strict relative minimum, then $H(\mathbf{x}^*)$ is *positive definite*, that is, $\mathbf{x}^T H(\mathbf{x}^*) \mathbf{x} > 0$ for any $\mathbf{x} \neq \mathbf{0}$. We can now derive sufficient conditions for a relative minimum.

If $\nabla f(\mathbf{x}^) = \mathbf{0}$ and $H(\mathbf{x}^*)$ is positive definite, then \mathbf{x}^* is a strict relative minimum.*

EXAMPLE 1 Suppose c is a real number, \mathbf{b} is a $n \times 1$ vector, and A is an $n \times n$ symmetric positive definite matrix. Consider the function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ given by

$$f(\mathbf{x}) = c + \mathbf{b}^T \mathbf{x} + \frac{1}{2} \mathbf{x}^T A \mathbf{x}. \quad (3)$$

We can easily see that

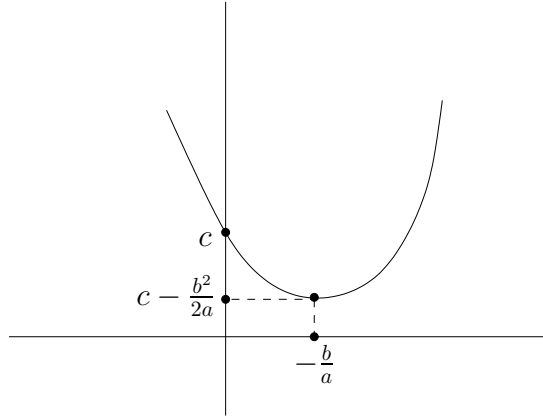
$$\begin{aligned} \nabla f(\mathbf{x}) &= \mathbf{b}^T + \mathbf{x}^T A, \\ H(\mathbf{x}) &= A. \end{aligned}$$

So there is a single extremum, located at \mathbf{x}^* , which is the solution to the system $A\mathbf{x} = -\mathbf{b}$. Since A is positive definite, this extremum is a strict local minimum. Since it is the only one, it is in fact a global minimum.

If we neglect the higher order terms inside the Big- O in the Taylor series (2), every function in one dimension behaves like this (locally near a minimum):

$$\begin{aligned} y &= c + bx + \frac{1}{2}ax^2, \\ y' &= b + ax, \\ y'' &= a > 0. \end{aligned}$$

So $y' = 0$ when $x = -\frac{b}{a}$. This is illustrated in the next figure.



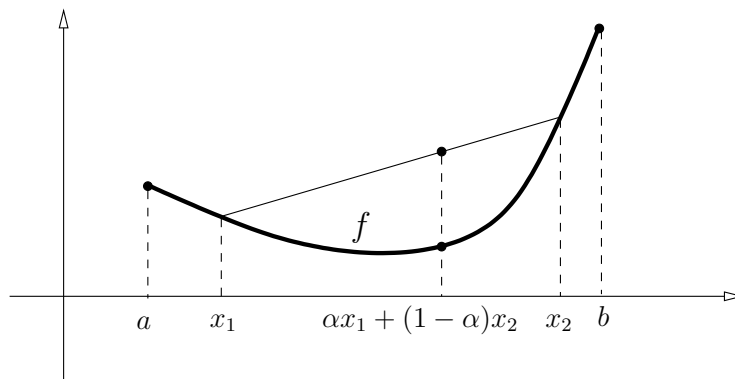
3 Convex Function

We just saw that H is positive definite at and near a strict local minimum. By Taylor's theorem every function looks like a quadratic near a strict local minimum. Furthermore, if f happens to be quadratic globally, formed from a symmetric positive definite matrix A as in (3), then it has a unique local minimum. This local minimum is therefore a global minimum.

Can we say more about functions whose local minima are global minima? A broad class of such functions are the *convex functions*.

A function $f : \Omega \rightarrow \mathbb{R}$ defined on a convex domain Ω is said to be *convex* if for every pair of points $\mathbf{x}_1, \mathbf{x}_2 \in \Omega$ and any α with $0 \leq \alpha \leq 1$, the following holds:

$$f(\alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2) \leq \alpha f(\mathbf{x}_1) + (1 - \alpha) f(\mathbf{x}_2).$$



Convex functions describe functions near local minima as studied in the following proposition.

Proposition 1 *Let $f \in C^2$. Then f is convex over a convex set Ω containing an interior point if and only if the Hessian matrix H is positive semi-definite in Ω .*

The minima of convex functions are global minima, as shown by the following theorem.

Theorem 2 *Let f be a convex function defined on a convex set Ω . Then the set Γ where f achieves its minimum value is convex. Furthermore, any relative minimum is a global minimum.*

Proof If f has no relative minima then the theorem is valid by default. Assume therefore that c_0 is the minimum value of f on Ω . Define the set

$$\Gamma = \{ \mathbf{x} \in \Omega \mid f(\mathbf{x}) = c_0 \}.$$

Suppose $\mathbf{x}_1, \mathbf{x}_2 \in \Gamma$ are two values that minimize f . We have that, for $0 \leq \alpha \leq 1$,

$$\begin{aligned} f(\alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2) &\leq \alpha f(\mathbf{x}_1) + (1 - \alpha) f(\mathbf{x}_2) \\ &= \alpha c_0 + (1 - \alpha) c_0 \\ &= c_0. \end{aligned}$$

But c_0 is the minimum, hence the above inequality must be an equality:

$$f(\alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2) = c_0.$$

In other words, f is also minimized at the point $\alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2$. Thus all the points on the line segment connecting \mathbf{x}_1 and \mathbf{x}_2 are in Γ . Since \mathbf{x}_1 and \mathbf{x}_2 are arbitrarily chosen from Γ , the set must be convex.

Suppose now that $\mathbf{x}^* \in \Omega$ is a relative minimum point of f but not a global minimum. Then there exists some $\mathbf{y} \in \Omega$ such that $f(\mathbf{y}) < f(\mathbf{x}^*)$. On the line $\{ \alpha \mathbf{y} + (1 - \alpha) \mathbf{x}^* \mid 0 \leq \alpha \leq 1 \}$, we have, for $0 < \alpha \leq 1$

$$\begin{aligned} f(\alpha \mathbf{y} + (1 - \alpha) \mathbf{x}^*) &\leq \alpha f(\mathbf{y}) + (1 - \alpha) f(\mathbf{x}^*) \\ &< \alpha f(\mathbf{x}^*) + (1 - \alpha) f(\mathbf{x}^*) \\ &= f(\mathbf{x}^*), \end{aligned}$$

contradicting the fact that \mathbf{x}^* is a relative minimum. □

4 Steepest Descent

Now let us return to the general problem of minimizing a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$. We want to find the critical values where the gradient $\nabla f = \mathbf{0}$. This is a system of n equations:

$$\begin{aligned} \frac{\partial f}{\partial x_1} &= 0, \\ &\vdots \\ \frac{\partial f}{\partial x_n} &= 0. \end{aligned}$$

We might expect to encounter the usual difficulties associated with higher-dimensional root finding. Fortunately, the derivative nature of the equations imposes some helpful structure to the problem.

In one dimension, to find a local minimum, we might employ the following rule:

if $f'(x) < 0$ then move to the right;
 if $f'(x) > 0$ then move to the left;
 if $f'(x) = 0$ then stop.

In higher dimensions we use the negative gradient $-\nabla f$ to point us toward a minimum. This is called *steepest descent*. In particular, the algorithm repeatedly performs one-dimensional minimizations along the direction of steepest descent.

In the algorithm, we start with $\mathbf{x}^{(0)}$ as an approximation to a local minimum of $f : \mathbb{R}^n \rightarrow \mathbb{R}$.

```

for  $m = 0, 1, 2, \dots$  until satisfied do
   $\mathbf{u} = \nabla f(\mathbf{x}^{(m)})$ 
  if  $\mathbf{u} = \mathbf{0}$  then stop
  else minimize the function  $g(t) = f(\mathbf{x}^{(m)} - t\mathbf{u})$ 
    let  $t^* > 0$  be the closest such minimum to zero
     $\mathbf{x}^{(m+1)} \leftarrow \mathbf{x}^{(m)} - t^*\mathbf{u}$ 
  
```

The method is also referred to as the *line search* strategy since during each iteration it moves on the line $\mathbf{x}^{(m)} - t\mathbf{u}$ away from $\mathbf{x}^{(m)}$ until encountering a local minimum of $g(t)$. How to carry out the line minimization of $g(t)$? Anyway you want. For instance, solve $g'(t) = 0$ directly. Or, step along the line until you produce a bracket, and then refine it.

EXAMPLE 2. The function

$$f(x_1, x_2) = x_1^3 + x_2^3 - 2x_1^2 + 3x_2^2 - 8$$

has the gradient

$$\nabla f = (3x_1^2 - 4x_1, 3x_2^2 + 6x_2).$$

Given the guess $\mathbf{x}^{(0)} = (1, -1)^T$ for a local minimum of f , we find

$$\nabla f(\mathbf{x}^{(0)}) = (-1, -3).$$

Thus, in the first step of steepest descent, we look for a minimum of the function

$$\begin{aligned} g(t) &= f(\mathbf{x}^{(0)} - t\nabla f(\mathbf{x}^{(0)})) \\ &= f(1+t, -1+3t) \\ &= (1+t)^3 + (-1+3t)^3 - 2(1+t)^2 + 3(-1+3t)^2 - 8. \end{aligned}$$

Setting $g'(t) = 0$ yields the equation

$$\begin{aligned} 0 &= 3(1+t)^2 + 3(3t-1)^2 - 4(1+t) + 3 \cdot 2(3t-1) \\ &= 84t^2 + 2t - 10, \end{aligned}$$

which has two solutions

$$t_1 = \frac{1}{3} \quad \text{and} \quad t_2 = -\frac{5}{14}.$$

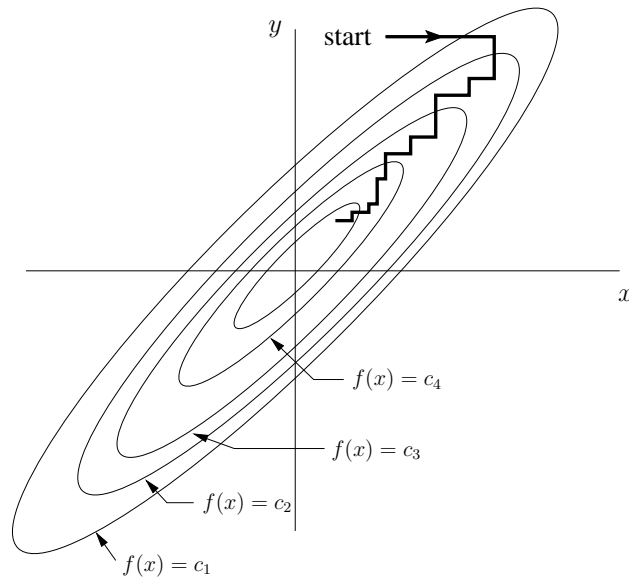
We choose the positive root, $t_1 = \frac{1}{3}$, since we intend to walk from $\mathbf{x}^{(0)}$ in the direction of $-\nabla f(\mathbf{x}^{(0)})$. This gives $\mathbf{x}^{(1)} = (\frac{4}{3}, 0)^T$. The gradient ∇f vanishes at $\mathbf{x}^{(1)}$. Therefore f achieves at least a local minimum at $(\frac{4}{3}, 0)^T$.

It turns out that steepest descent converges globally to a relative minimum. And the convergence rate is linear. Let A and a be the largest and smallest eigenvalues, respectively, of the Hessian H at the local minimum. Then the following holds regarding the ratio between the errors at two adjacent steps

$$\frac{|e_{m+1}|}{|e_m|} \sim \left(\frac{A - a}{A + a} \right)^2.$$

The steepest descent method can take many steps. The problem is that the method repeatedly moves in a steepest direction to a minimum along that direction. Consequently, consecutive steps are perpendicular to each other (this behavior is illustrated in the figure below). To see why, consider moving at $\mathbf{x}^{(k)}$ along $\mathbf{u} = -\nabla f(\mathbf{x}^{(k)})$ to reach $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t^* \mathbf{u}$ where $f(\mathbf{x})$ no longer decreases. The rate at which the value of f changes, after a movement of $t\mathbf{u}$ from $\mathbf{x}^{(k)}$, is measured by the directional derivative, $\nabla f(\mathbf{x}^{(k)} + t\mathbf{u}) \cdot \mathbf{u}$. This derivative is negative at $\mathbf{x}^{(k)}$ (when $t = 0$), and has not changed its sign before $\mathbf{x}^{(k+1)}$ (when $t = t^*$). Suppose $\nabla f(\mathbf{x}^{(k+1)}) \not\perp \nabla f(\mathbf{x}^{(k)})$. Then, $\nabla f(\mathbf{x}^{(k+1)}) \cdot \mathbf{u} < 0$ must hold. This means that the value of f would further decrease if continuing the movement in the direction \mathbf{u} from $\mathbf{x}^{(k+1)}$. Hence a contradiction.

So there are a number of back and forth steps that only slowly converge to a minimum. This situation can get very bad in a narrow valley, where successive steps undo some of their previous progress. Ideally, in \mathbb{R}^n we would like to take n perpendicular steps, each of which attains a minimum. This idea will lead to the conjugate gradient method.



A Matrix Calculus

This appendix presents some basic rules of differentiations of scalars and vectors with respect to vectors and matrices. These rules will be used later on in the course.

A.1 Differentiation With Respect to a Vector

The derivative of a vector function $\mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_m(\mathbf{x}))^T$ with respect to \mathbf{x} is an $m \times n$ matrix:

$$\frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \dots & \frac{\partial f_m}{\partial x_n} \end{pmatrix}. \quad (4)$$

The above matrix is referred to as the *Jacobian* of \mathbf{f} .

Let \mathbf{c} be an n -vector and A an $m \times n$ matrix. Then $\mathbf{c}^T \mathbf{x}$ and $A\mathbf{x}$ are scalar and vector functions of \mathbf{x} , respectively. Applying (1) and (4), respectively, we can easily verify the following:

$$\frac{\partial(\mathbf{c}^T \mathbf{x})}{\partial \mathbf{x}} = \mathbf{c}^T, \quad (5)$$

$$\frac{\partial(A\mathbf{x})}{\partial \mathbf{x}} = A. \quad (6)$$

A.2 Differentiation of Inner and Cross Products

Suppose that \mathbf{u} and \mathbf{v} are both three-dimensional vectors. Then, equation (5) is applied in differentiation of their dot product:

$$\begin{aligned} \frac{\partial(\mathbf{u} \cdot \mathbf{v})}{\partial \mathbf{v}} &= \frac{\partial(\mathbf{u}^T \mathbf{v})}{\partial \mathbf{v}} = \mathbf{u}^T, \\ \frac{\partial(\mathbf{u} \cdot \mathbf{v})}{\partial \mathbf{u}} &= \frac{\partial(\mathbf{v}^T \mathbf{u})}{\partial \mathbf{u}} = \mathbf{v}^T. \end{aligned}$$

To differentiate the cross product $\mathbf{u} \times \mathbf{v}$, for any vector $\mathbf{w} = (w_1, w_2, w_3)^T$ we denote by $\mathbf{w} \times$ the following 3×3 anti-symmetric matrix:

$$\mathbf{w} \times = \begin{pmatrix} 0 & -w_3 & w_2 \\ w_3 & 0 & -w_1 \\ -w_2 & w_1 & 0 \end{pmatrix}.$$

Apparently, the product of the matrix $\mathbf{w} \times$ with \mathbf{v} is the cross product $\mathbf{w} \times \mathbf{v}$. It then follows that

$$\begin{aligned} \frac{\partial(\mathbf{u} \times \mathbf{v})}{\partial \mathbf{v}} &= \mathbf{u} \times, \\ \frac{\partial(\mathbf{u} \times \mathbf{v})}{\partial \mathbf{u}} &= -\frac{\partial(\mathbf{v} \times \mathbf{u})}{\partial \mathbf{u}} \\ &= -\mathbf{v} \times. \end{aligned}$$

Now let us look at how to differentiate the scalar $\mathbf{x}^T A \mathbf{x}$, where \mathbf{x} is an n -vector and A an $n \times n$ matrix, with respect to \mathbf{x} . We have

$$\frac{\partial \mathbf{x}^T A \mathbf{x}}{\partial \mathbf{x}} = \frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T (A\mathbf{y})) \Big|_{\mathbf{y}=\mathbf{x}} + \frac{\partial}{\partial \mathbf{x}} ((\mathbf{y}^T A)\mathbf{x}) \Big|_{\mathbf{y}=\mathbf{x}}$$

$$\begin{aligned}
&= \frac{\partial}{\partial \mathbf{x}} \left((A\mathbf{y})^T \mathbf{x} \right) \Big|_{y=x} + \mathbf{y}^T A \Big|_{y=x} \\
&= (A\mathbf{y})^T \Big|_{y=x} + \mathbf{y}^T A \Big|_{y=x} \\
&= \mathbf{x}^T A^T + \mathbf{x}^T A.
\end{aligned}$$

The above reduces to $2\mathbf{x}^T A$ when A is symmetric.

A.3 Trace of a Product Matrix

The derivative of a scalar function $f(A)$, where $A = (a_{ij})_{m \times n}$ is a matrix, is defined as

$$\frac{\partial f}{\partial A} = \begin{pmatrix} \frac{\partial f}{\partial a_{11}} & \frac{\partial f}{\partial a_{12}} & \cdots & \frac{\partial f}{\partial a_{1n}} \\ \frac{\partial f}{\partial a_{21}} & \frac{\partial f}{\partial a_{22}} & \cdots & \frac{\partial f}{\partial a_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial a_{m1}} & \frac{\partial f}{\partial a_{m2}} & \cdots & \frac{\partial f}{\partial a_{mn}} \end{pmatrix}. \quad (7)$$

Let $C = (c_{ij})_{m \times n}$ and $X = (x_{ij})_{m \times n}$ be two matrices. The trace of the product matrix CX^T is the sum of its diagonal entries:

$$\text{Tr}(CX^T) = \sum_{r=1}^m \left(\sum_{s=1}^n c_{rs} x_{rs} \right).$$

Immediately, we have

$$\frac{\partial}{\partial x_{ij}} \text{Tr}(CX^T) = c_{ij},$$

which implies that

$$\frac{\partial}{\partial X} \text{Tr}(CX^T) = C.$$

Next, we differentiate the trace of the product matrix XCX^T as follows:

$$\begin{aligned}
\frac{\partial}{\partial X} \text{Tr}(XCX^T) &= \frac{\partial}{\partial X} \text{Tr}(XCX^T) \Big|_{Y=X} + \frac{\partial}{\partial X} \text{Tr}((YX)X^T) \Big|_{Y=X} \\
&= \frac{\partial}{\partial X} \text{Tr}((YX^T)X^T) \Big|_{Y=X} + YX \Big|_{Y=X} \\
&= YX^T \Big|_{Y=X} + YX \\
&= XC^T + XC.
\end{aligned}$$

The above reduces to $2XC$ when C is symmetric.

A.4 Matrix in One Variable

Consider an $m \times n$ matrix A in which every element is a function of some variable t . Write the matrix as $A(t) = (a_{ij}(t))_{m \times n}$. Let the overdot denote differentiation with respect to t . Differentiation and

integration of the matrix operate element-wise:

$$\begin{aligned}\dot{A}(t) &= (\dot{a}_{ij}), \\ \int A(t) dt &= \left(\int a_{ij}(t) dt \right).\end{aligned}$$

Suppose $A(t)$ is $n \times n$ and non-singular. Then we have $AA^{-1} = I_n$, the $n \times n$ identity matrix. Thus,

$$\begin{aligned}0 &= \frac{d}{dt}(AA^{-1}) \\ &= \dot{A}A^{-1} + A\frac{d}{dt}(A^{-1}),\end{aligned}$$

which yields the derivative of the inverse matrix:

$$\frac{d}{dt}(A^{-1}) = -A^{-1}\dot{A}A^{-1}. \quad (8)$$

An interesting case is with the rotation matrix R , which is also orthogonal, i.e., $RR^T = R^T R = I_n$. We obtain

$$\begin{aligned}RR^T = I &\Rightarrow \dot{R}R^T + R\dot{R}^T = 0 \\ &\Rightarrow \dot{R}R^T + (\dot{R}R^T)^T = 0 \\ &\Rightarrow \dot{R}R^T = -(\dot{R}R^T)^T.\end{aligned}$$

The above implies that the matrix $\dot{R}R^T$ is anti-symmetric. Therefore, it can be written as

$$\dot{R}R^T = \begin{pmatrix} 0 & -\omega_z & \omega_y \\ \omega_z & 0 & -\omega_x \\ -\omega_y & \omega_x & 0 \end{pmatrix}.$$

The vector $\boldsymbol{\omega} = (\omega_x, \omega_y, \omega_z)^T$ is the angular velocity, where the cross product $\boldsymbol{\omega} \times \boldsymbol{v} = \dot{R}R^T \boldsymbol{v}$ describes the change rate of the vector \boldsymbol{v} (i.e., the velocity of the destination point of the vector) due to the body rotation described by R .

Using the Taylor expansion, we define the exponential function:

$$e^{At} = \sum_{j=0}^{\infty} \frac{(At)^j}{j!},$$

where A is an $n \times n$ matrix. The function's importance comes from that it is the solution to the linear system $\dot{\boldsymbol{x}} = A\boldsymbol{x} + \boldsymbol{b}\boldsymbol{u}$, where \boldsymbol{u} is the control vector. It has the derivative

$$\frac{d}{dt}e^{At} = Ae^{At} = e^{At}A.$$

References

- [1] M. Erdmann. Lecture notes for *16-811 Mathematical Fundamentals for Robotics*. The Robotics Institute, Carnegie Mellon University, 1998.
- [2] D. G. Luenberger. *Linear and Nonlinear Programming*. Addison-Wesley, 2nd edition, 1984.
- [3] W. H. Press, *et al.* *Numerical Recipes in C++: The Art of Scientific Computing*. Cambridge University Press, 2nd edition, 2002.