

Chapter 2

Notation and Norms

2.1 Introduction

This chapter recalls the usual convention for distinguishing scalars, vectors, and matrices. Vetter's notation for matrix derivatives is then explained, as well as the meaning of the expressions *little o* and *big O* employed for comparing the local or asymptotic behaviors of functions. The most important vector and matrix norms are finally described. Norms find a first application in the definition of types of convergence speeds for iterative algorithms.

2.2 Scalars, Vectors, and Matrices

Unless stated otherwise, scalar variables are real valued, as are the entries of vectors and matrices.

Italics are for *scalar variables* (v or V), bold lower-case letters for *column vectors* (\mathbf{v}), and bold upper-case letters for *matrices* (\mathbf{M}). *Transposition*, the transformation of columns into rows in a vector or matrix, is denoted by the superscript T . It applies to what is to its *left*, so \mathbf{v}^T is a *row vector* and, in $\mathbf{A}^T \mathbf{B}$, \mathbf{A} is transposed, not \mathbf{B} .

The *identity matrix* is \mathbf{I} , with \mathbf{I}_n the $(n \times n)$ identity matrix. The i th column vector of \mathbf{I} is the *canonical vector* \mathbf{e}^i .

The entry at the intersection of the i th row and j th column of \mathbf{M} is $m_{i,j}$. The product of matrices

$$\mathbf{C} = \mathbf{A}\mathbf{B} \quad (2.1)$$

thus implies that

$$c_{i,j} = \sum_k a_{i,k} b_{k,j}, \quad (2.2)$$

and the number of columns in \mathbf{A} must be equal to the number of rows in \mathbf{B} . Recall that the product of matrices (or vectors) is *not commutative*, in general. Thus, for instance, when \mathbf{v} and \mathbf{w} are columns vectors with the same dimension, $\mathbf{v}^T \mathbf{w}$ is a scalar whereas $\mathbf{w} \mathbf{v}^T$ is a (rank-one) square matrix.

Useful relations are

$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T, \quad (2.3)$$

and, provided that \mathbf{A} and \mathbf{B} are invertible,

$$(\mathbf{AB})^{-1} = \mathbf{B}^{-1} \mathbf{A}^{-1}. \quad (2.4)$$

If \mathbf{M} is square and symmetric, then all of its eigenvalues are real. $\mathbf{M} \succ 0$ then means that each of these eigenvalues is strictly positive (\mathbf{M} is *positive definite*), while $\mathbf{M} \succeq 0$ allows some of them to be zero (\mathbf{M} is *non-negative definite*).

2.3 Derivatives

Provided that $f(\cdot)$ is a sufficiently differentiable function from \mathbb{R} to \mathbb{R} ,

$$\dot{f}(x) = \frac{df}{dx}(x), \quad (2.5)$$

$$\ddot{f}(x) = \frac{d^2 f}{dx^2}(x), \quad (2.6)$$

$$f^{(k)}(x) = \frac{d^k f}{dx^k}(x). \quad (2.7)$$

Vetter's notation [1] will be used for derivatives of matrices with respect to matrices. (A word of caution is in order: there are other, incompatible notations, and one should be cautious about mixing formulas from different sources.)

If \mathbf{A} is $(n_A \times m_A)$ and \mathbf{B} $(n_B \times m_B)$, then

$$\mathbf{M} = \frac{\partial \mathbf{A}}{\partial \mathbf{B}} \quad (2.8)$$

is an $(n_A n_B \times m_A m_B)$ matrix, such that the $(n_A \times m_A)$ submatrix in position (i, j) is

$$\mathbf{M}_{i,j} = \frac{\partial \mathbf{A}}{\partial b_{i,j}}. \quad (2.9)$$

Remark 2.1 \mathbf{A} and \mathbf{B} in (2.8) may be row or column vectors. □

Example 2.1 If \mathbf{v} is a generic column vector of \mathbb{R}^n , then

$$\frac{\partial \mathbf{v}}{\partial \mathbf{v}^T} = \frac{\partial \mathbf{v}^T}{\partial \mathbf{v}} = \mathbf{I}_n. \quad (2.10)$$

□

Example 2.2 If $J(\cdot)$ is a differentiable function from \mathbb{R}^n to \mathbb{R} , and \mathbf{x} a vector of \mathbb{R}^n , then

$$\frac{\partial J}{\partial \mathbf{x}}(\mathbf{x}) = \begin{bmatrix} \frac{\partial J}{\partial x_1} \\ \frac{\partial J}{\partial x_2} \\ \vdots \\ \frac{\partial J}{\partial x_n} \end{bmatrix}(\mathbf{x}) \quad (2.11)$$

is a column vector, called the *gradient* of $J(\cdot)$ at \mathbf{x} . □

Example 2.3 If $J(\cdot)$ is a twice differentiable function from \mathbb{R}^n to \mathbb{R} , and \mathbf{x} a vector of \mathbb{R}^n , then

$$\frac{\partial^2 J}{\partial \mathbf{x} \partial \mathbf{x}^T}(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 J}{\partial x_1^2} & \frac{\partial^2 J}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 J}{\partial x_1 \partial x_n} \\ \frac{\partial^2 J}{\partial x_2 \partial x_1} & \frac{\partial^2 J}{\partial x_2^2} & & \vdots \\ \vdots & & \ddots & \vdots \\ \frac{\partial^2 J}{\partial x_n \partial x_1} & \cdots & \cdots & \frac{\partial^2 J}{\partial x_n^2} \end{bmatrix}(\mathbf{x}) \quad (2.12)$$

is an $(n \times n)$ matrix, called the *Hessian* of $J(\cdot)$ at \mathbf{x} . Schwarz's theorem ensures that

$$\frac{\partial^2 J}{\partial x_i \partial x_j}(\mathbf{x}) = \frac{\partial^2 J}{\partial x_j \partial x_i}(\mathbf{x}), \quad (2.13)$$

provided that both are continuous at \mathbf{x} and \mathbf{x} belongs to an open set in which both are defined. Hessians are thus symmetric, except in pathological cases not considered here. □

Example 2.4 If $\mathbf{f}(\cdot)$ is a differentiable function from \mathbb{R}^n to \mathbb{R}^p , and \mathbf{x} a vector of \mathbb{R}^n , then

$$\mathbf{J}(\mathbf{x}) = \frac{\partial \mathbf{f}}{\partial \mathbf{x}^T}(\mathbf{x}) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & & \vdots \\ \vdots & & \ddots & \vdots \\ \frac{\partial f_p}{\partial x_1} & \cdots & \cdots & \frac{\partial f_p}{\partial x_n} \end{bmatrix} \quad (2.14)$$

is the $(p \times n)$ *Jacobian matrix* of $\mathbf{f}(\cdot)$ at \mathbf{x} . When $p = n$, the Jacobian matrix is square and its determinant is the *Jacobian*. □

Remark 2.2 The last three examples show that the Hessian of $J(\cdot)$ at \mathbf{x} is the Jacobian matrix of its gradient function evaluated at \mathbf{x} . \square

Remark 2.3 Gradients and Hessians are frequently used in the context of optimization, and Jacobian matrices when solving systems of nonlinear equations. \square

Remark 2.4 The *Nabla operator* ∇ , a vector of partial derivatives with respect to all the variables of the function on which it operates

$$\nabla = \left(\frac{\partial}{\partial x_1}, \dots, \frac{\partial}{\partial x_n} \right)^T, \quad (2.15)$$

is often used to make notation more concise, especially for partial differential equations. Applying ∇ to a scalar function J and evaluating the result at \mathbf{x} , one gets the gradient vector

$$\nabla J(\mathbf{x}) = \frac{\partial J}{\partial \mathbf{x}}(\mathbf{x}). \quad (2.16)$$

If the scalar function is replaced by a vector function \mathbf{f} , one gets the Jacobian matrix

$$\nabla \mathbf{f}(\mathbf{x}) = \frac{\partial \mathbf{f}}{\partial \mathbf{x}^T}(\mathbf{x}), \quad (2.17)$$

where $\nabla \mathbf{f}$ is interpreted as $(\nabla \mathbf{f}^T)^T$.

By applying ∇ twice to a scalar function J and evaluating the result at \mathbf{x} , one gets the Hessian matrix

$$\nabla^2 J(\mathbf{x}) = \frac{\partial^2 J}{\partial \mathbf{x} \partial \mathbf{x}^T}(\mathbf{x}). \quad (2.18)$$

(∇^2 is sometimes taken to mean the *Laplacian operator* Δ , such that

$$\Delta f(\mathbf{x}) = \sum_{i=1}^n \frac{\partial^2 f}{\partial x_i^2}(\mathbf{x}) \quad (2.19)$$

is a scalar. The context and dimensional considerations should make what is meant clear.) \square

Example 2.5 If \mathbf{v} , \mathbf{M} , and \mathbf{Q} do not depend on \mathbf{x} and \mathbf{Q} is symmetric, then

$$\frac{\partial}{\partial \mathbf{x}}(\mathbf{v}^T \mathbf{x}) = \mathbf{v}, \quad (2.20)$$

$$\frac{\partial}{\partial \mathbf{x}^T}(\mathbf{M} \mathbf{x}) = \mathbf{M}, \quad (2.21)$$

$$\frac{\partial}{\partial \mathbf{x}}(\mathbf{x}^T \mathbf{M} \mathbf{x}) = (\mathbf{M} + \mathbf{M}^T) \mathbf{x} \quad (2.22)$$

and

$$\frac{\partial}{\partial \mathbf{x}}(\mathbf{x}^T \mathbf{Q} \mathbf{x}) = 2\mathbf{Q} \mathbf{x}. \quad (2.23)$$

These formulas will be used quite frequently. \square

2.4 Little o and Big O

The function $f(x)$ is $o(g(x))$ as x tends to x_0 if

$$\lim_{x \rightarrow x_0} \frac{f(x)}{g(x)} = 0, \quad (2.24)$$

so $f(x)$ gets negligible compared to $g(x)$ for x sufficiently close to x_0 . In what follows, x_0 is always taken equal to zero, so this need not be specified, and we just write $f(x) = o(g(x))$.

The function $f(x)$ is $O(g(x))$ as x tends to *infinity* if there exists real numbers x_0 and M such that

$$x > x_0 \Rightarrow |f(x)| \leq M|g(x)|. \quad (2.25)$$

The function $f(x)$ is $O(g(x))$ as x tends to *zero* if there exists real numbers δ and M such that

$$|x| < \delta \Rightarrow |f(x)| \leq M|g(x)|. \quad (2.26)$$

The notation $O(x)$ or $O(n)$ will be used in two contexts:

- when dealing with Taylor expansions, x is a real number tending to zero,
- when analyzing algorithmic complexity, n is a positive integer tending to infinity.

Example 2.6 The function

$$f(x) = \sum_{i=2}^m a_i x^i,$$

with $m \geq 2$, is such that

$$\lim_{x \rightarrow 0} \frac{f(x)}{x} = \lim_{x \rightarrow 0} \left(\sum_{i=2}^m a_i x^{i-1} \right) = 0,$$

so $f(x) = o(x)$ when x tends to zero. Now, if $|x| < 1$, then

$$\frac{|f(x)|}{x^2} < \sum_{i=2}^m |a_i|,$$

so $f(x) = O(x^2)$ when x tends to zero. If, on the other hand, x is taken equal to the (large) positive integer n , then

$$\begin{aligned} f(n) &= \sum_{i=2}^m a_i n^i \leq \sum_{i=2}^m |a_i n^i| \\ &\leq \left(\sum_{i=2}^m |a_i| \right) \cdot n^m, \end{aligned}$$

so $f(n) = O(n^m)$ when n tends to infinity. □

2.5 Norms

A function $f(\cdot)$ from a vector space \mathbb{V} to \mathbb{R} is a *norm* if it satisfies the following three properties:

1. $f(\mathbf{v}) \geq 0$ for all $\mathbf{v} \in \mathbb{V}$ (*positivity*),
2. $f(\alpha \mathbf{v}) = |\alpha| \cdot f(\mathbf{v})$ for all $\alpha \in \mathbb{R}$ and $\mathbf{v} \in \mathbb{V}$ (*positive scalability*),
3. $f(\mathbf{v}^1 \pm \mathbf{v}^2) \leq f(\mathbf{v}^1) + f(\mathbf{v}^2)$ for all $\mathbf{v}^1 \in \mathbb{V}$ and $\mathbf{v}^2 \in \mathbb{V}$ (*triangle inequality*).

These properties imply that $f(\mathbf{v}) = 0 \Rightarrow \mathbf{v} = \mathbf{0}$ (*non-degeneracy*). Another useful relation is

$$|f(\mathbf{v}^1) - f(\mathbf{v}^2)| \leq f(\mathbf{v}^1 \pm \mathbf{v}^2). \quad (2.27)$$

Norms are used to quantify distances between vectors. They play an essential role, for instance, in the characterization of the intrinsic difficulty of numerical problems via the notion of condition number (see Sect. 3.3) or in the definition of cost functions for optimization.

2.5.1 Vector Norms

The most commonly used norms in \mathbb{R}^n are the l_p norms

$$\|\mathbf{v}\|_p = \left(\sum_{i=1}^n |v_i|^p \right)^{\frac{1}{p}}, \quad (2.28)$$

with $p \geq 1$. They include

- the *Euclidean norm* (or l_2 norm)

$$\|\mathbf{v}\|_2 = \sqrt{\sum_{i=1}^n v_i^2} = \sqrt{\mathbf{v}^T \mathbf{v}}, \quad (2.29)$$

- the *taxicab norm* (or Manhattan norm, or grid norm, or l_1 norm)

$$\|\mathbf{v}\|_1 = \sum_{i=1}^n |v_i|, \quad (2.30)$$

- the *maximum norm* (or l_∞ norm, or Chebyshev norm, or uniform norm)

$$\|\mathbf{v}\|_\infty = \max_{1 \leq i \leq n} |v_i|. \quad (2.31)$$

They are such that

$$\|\mathbf{v}\|_2 \leq \|\mathbf{v}\|_1 \leq n \|\mathbf{v}\|_\infty, \quad (2.32)$$

and

$$\mathbf{v}^T \mathbf{w} \leq \|\mathbf{v}\|_2 \cdot \|\mathbf{w}\|_2. \quad (2.33)$$

The latter result is known as the *Cauchy-Schwarz inequality*.

Remark 2.5 If the entries of \mathbf{v} were complex, norms would be defined differently. The Euclidean norm, for instance, would become

$$\|\mathbf{v}\|_2 = \sqrt{\mathbf{v}^H \mathbf{v}}, \quad (2.34)$$

where \mathbf{v}^H is the *transconjugate* of \mathbf{v} , i.e., the row vector obtained by transposing the column vector \mathbf{v} and replacing each of its entries by its complex conjugate. \square

Example 2.7 For the complex vector

$$\mathbf{v} = \begin{bmatrix} a \\ ai \end{bmatrix},$$

where a is some nonzero real number and i is the imaginary unit (such that $i^2 = -1$), $\mathbf{v}^T \mathbf{v} = 0$. This proves that $\sqrt{\mathbf{v}^T \mathbf{v}}$ is not a norm. The value of the Euclidean norm of \mathbf{v} is $\sqrt{\mathbf{v}^H \mathbf{v}} = \sqrt{2}|a|$. \square

Remark 2.6 The so-called l_0 norm of a vector is the number of its nonzero entries. Used in the context of sparse estimation, where one is looking for an estimated parameter vector with as few nonzero entries as possible, it is *not* a norm, as it does not satisfy the property of positive scalability. \square

2.5.2 Matrix Norms

Each vector norm *induces* a *matrix norm*, defined as

$$||\mathbf{M}|| = \max_{||\mathbf{v}||=1} ||\mathbf{M}\mathbf{v}||, \quad (2.35)$$

so

$$||\mathbf{M}\mathbf{v}|| \leq ||\mathbf{M}|| \cdot ||\mathbf{v}|| \quad (2.36)$$

for any \mathbf{M} and \mathbf{v} for which the product $\mathbf{M}\mathbf{v}$ makes sense. This matrix norm is *sub-ordinate* to the vector norm inducing it. The matrix and vector norms are then said to be *compatible*, an important property for the study of products of matrices and vectors.

- The matrix norm induced by the vector norm l_2 is the *spectral norm*, or *2-norm*,

$$||\mathbf{M}||_2 = \sqrt{\rho(\mathbf{M}^T\mathbf{M})}, \quad (2.37)$$

where $\rho(\cdot)$ is the function that computes the *spectral radius* of its argument, i.e., the modulus of the eigenvalue(s) with the largest modulus. Since all the eigenvalues of $\mathbf{M}^T\mathbf{M}$ are real and non-negative, $\rho(\mathbf{M}^T\mathbf{M})$ is the largest of these eigenvalues. Its square root is the largest *singular value* of \mathbf{M} , denoted by $\sigma_{\max}(\mathbf{M})$. So

$$||\mathbf{M}||_2 = \sigma_{\max}(\mathbf{M}). \quad (2.38)$$

- The matrix norm induced by the vector norm l_1 is the *1-norm*

$$||\mathbf{M}||_1 = \max_j \sum_i |m_{i,j}|, \quad (2.39)$$

which amounts to summing the absolute values of the entries of each *column* in turn and keeping the largest result.

- The matrix norm induced by the vector norm l_∞ is the *infinity norm*

$$||\mathbf{M}||_\infty = \max_i \sum_j |m_{i,j}|, \quad (2.40)$$

which amounts to summing the absolute values of the entries of each *row* in turn and keeping the largest result. Thus

$$||\mathbf{M}||_1 = ||\mathbf{M}^T||_\infty. \quad (2.41)$$

Since each subordinate matrix norm is compatible with its inducing vector norm,

$$\|\mathbf{v}\|_1 \text{ is compatible with } \|\mathbf{M}\|_1, \quad (2.42)$$

$$\|\mathbf{v}\|_2 \text{ is compatible with } \|\mathbf{M}\|_2, \quad (2.43)$$

$$\|\mathbf{v}\|_\infty \text{ is compatible with } \|\mathbf{M}\|_\infty. \quad (2.44)$$

The *Frobenius norm*

$$\|\mathbf{M}\|_F = \sqrt{\sum_{i,j} m_{i,j}^2} = \sqrt{\text{trace}(\mathbf{M}^T \mathbf{M})} \quad (2.45)$$

deserves a special mention, as it is not induced by any vector norm yet

$$\|\mathbf{v}\|_2 \text{ is compatible with } \|\mathbf{M}\|_F. \quad (2.46)$$

Remark 2.7 To evaluate a vector or matrix norm with MATLAB (or any other interpreted language based on matrices), it is much more efficient to use the corresponding dedicated function than to access the entries of the vector or matrix individually to implement the norm definition. Thus, `norm(X, p)` returns the p -norm of \mathbf{X} , which may be a vector or a matrix, while `norm(M, 'fro')` returns the Frobenius norm of the matrix \mathbf{M} . \square

2.5.3 Convergence Speeds

Norms can be used to study how quickly an iterative method would converge to the solution \mathbf{x}^* if computation were exact. Define the error at iteration k as

$$\mathbf{e}^k = \mathbf{x}^k - \mathbf{x}^*, \quad (2.47)$$

where \mathbf{x}^k is the estimate of \mathbf{x}^* at iteration k . The *asymptotic* convergence speed is *linear* if

$$\limsup_{k \rightarrow \infty} \frac{\|\mathbf{e}^{k+1}\|}{\|\mathbf{e}^k\|} = \alpha < 1, \quad (2.48)$$

with α the rate of convergence.

It is *superlinear* if

$$\limsup_{k \rightarrow \infty} \frac{\|\mathbf{e}^{k+1}\|}{\|\mathbf{e}^k\|} = 0, \quad (2.49)$$

and *quadratic* if

$$\limsup_{k \rightarrow \infty} \frac{\|\mathbf{e}^{k+1}\|}{\|\mathbf{e}^k\|^2} = \alpha < \infty. \quad (2.50)$$

A method with quadratic convergence thus also has superlinear and linear convergence. It is customary, however, to qualify a method with the best convergence it achieves. Quadratic convergence is better than superlinear convergence, which is better than linear convergence.

Remember that these convergence speeds are asymptotic, valid when the error has become small enough, and that they do not take the effect of rounding into account. They are meaningless if the initial vector \mathbf{x}^0 was too badly chosen for the method to converge to \mathbf{x}^* . When the method does converge to \mathbf{x}^* , they may not describe accurately its initial behavior and will no longer be true when rounding errors become predominant. They are nevertheless an interesting indication of what can be expected at best.

Reference

1. Vetter, W.: Derivative operations on matrices. IEEE Trans. Autom. Control **15**, 241–244 (1970)

<http://www.springer.com/978-3-319-07670-6>

Numerical Methods and Optimization

A Consumer Guide

Walter, É.

2014, XV, 476 p. 67 illus., 23 illus. in color., Hardcover

ISBN: 978-3-319-07670-6