

# Matrix Calculus

Yangang Cao

June 19, 2019

While the topics in the previous sections are typically covered in a standard course on linear algebra, one topic that does not seem to be covered very often (and which we will use extensively) is the extension of calculus to the vector setting. Despite the fact that all the actual calculus we use is relatively trivial, the notation can often make things look much more difficult than they are. In this section we present some basic definitions of matrix calculus and provide a few examples.

## 1 The Gradient

Suppose that  $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$  is a function that takes as input a matrix of size  $m \times n$  and returns a real value. Then the **gradient** of  $f$  (with respect to  $A \in \mathbb{R}^{m \times n}$ ) is the matrix of partial derivatives, defined as:

$$\nabla_A f(A) \in \mathbb{R}^{m \times n} = \begin{bmatrix} \frac{\partial f(A)}{\partial A_{11}} & \frac{\partial f(A)}{\partial A_{12}} & \cdots & \frac{\partial f(A)}{\partial A_{1n}} \\ \frac{\partial f(A)}{\partial A_{21}} & \frac{\partial f(A)}{\partial A_{22}} & \cdots & \frac{\partial f(A)}{\partial A_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f(A)}{\partial A_{m1}} & \frac{\partial f(A)}{\partial A_{m2}} & \cdots & \frac{\partial f(A)}{\partial A_{mn}} \end{bmatrix}$$

i.e., an  $m \times n$  matrix with

$$(\nabla_A f(A))_{ij} = \frac{\partial f(A)}{\partial A_{ij}}.$$

Note that the size of  $\nabla_A f(A)$  is always the same as the size of  $A$ . So if, in particular,  $A$  is just a vector  $x \in \mathbb{R}^n$ ,

$$\nabla_x f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix}.$$

It is very important to remember that the gradient of a function is only defined if the function is real-valued, that is, it returns a scalar value. We can not, for example, take the gradient of  $Ax$ ,  $A \in \mathbb{R}^{n \times n}$  with respect to  $x$ , since this quantity is vector-valued.

It follows directly from the equivalent properties of partial derivatives that:

- $\nabla_x(f(x) + g(x)) = \nabla_x f(x) + \nabla_x g(x)$
- For  $t \in \mathbb{R}$ ,  $\nabla_x(tf(x)) = t\nabla_x f(x)$

In principle, gradients are a natural extension of partial derivatives to functions of multiple variables. In practice, however, working with gradients can sometimes be tricky for notational reasons. For example, suppose that  $A \in \mathbb{R}^{m \times n}$  is a matrix of fixed coefficients and suppose that  $b \in \mathbb{R}^m$  is a vector of fixed coefficients. Let  $f : \mathbb{R}^m \rightarrow \mathbb{R}$  be the function defined by  $f(z) = z^T z$ , such that  $\nabla_z f(z) = 2z$ . But now, consider the expression,

$$\nabla f(Ax).$$

How should this expression be interpreted? There are at least two possibilities:

1. In the first interpretation, recall that  $\nabla_z f(z) = 2z$ . Here, we interpret  $\nabla f(Ax)$  as evaluating the gradient at the point  $Ax$ , hence,

$$\nabla f(Ax) = 2(Ax) = 2Ax \in \mathbb{R}^m$$

2. In the second interpretation, we consider the quantity  $f(Ax)$  as a function of the input variables  $x$ . More formally, let  $g(x) = f(Ax)$ . Then in this interpretation,

$$\nabla f(Ax) = \nabla_x g(x) \in \mathbb{R}^n$$

Here, we can see that these two interpretations are indeed different. One interpretation yields an  $m$ -dimensional vector as a result, while the other interpretation yields an  $n$ -dimensional vector as a result! How can we resolve this?

Here, the key is to make explicit the variables which we are differentiating with respect to. In the first case, we are differentiating the function  $f$  with respect to its arguments  $z$  and then substituting the argument  $Ax$ . In the second case, we are differentiating the composite function  $g(x) = f(Ax)$  with respect to  $x$  directly. We denote the first case as  $\nabla_z f(Ax)$  and the second case as  $\nabla_x f(Ax)$ . Keeping the notation clear is extremely important (as you'll find out in your homework, in fact!).

## 2 The Hessian

Suppose that  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is a function that takes a vector in  $\mathbb{R}^n$  and returns a real number. Then the **Hessian** matrix with respect to  $x$ , written  $\nabla_x^2 f(x)$  or simply as  $H$  is the  $n \times n$  matrix of partial derivatives,

$$\nabla_x^2 f(x) \in \mathbb{R}^{n \times n} = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2^2} & \cdots & \frac{\partial^2 f(x)}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_n \partial x_1} & \frac{\partial^2 f(x)}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_n^2} \end{bmatrix}$$

In other words,  $\nabla_x^2 f(x) \in \mathbb{R}^{n \times n}$ , with

$$(\nabla_x^2 f(x))_{ij} = \frac{\partial^2 f(x)}{\partial x_i \partial x_j}$$

Note that the Hessian is always symmetric, since

$$\frac{\partial^2 f(x)}{\partial x_i \partial x_j} = \frac{\partial^2 f(x)}{\partial x_j \partial x_i}$$

Similar to the gradient, the Hessian is defined only when  $f(x)$  is real-valued.

It is natural to think of the gradient as the analogue of the first derivative for functions of vectors, and the Hessian as the analogue of the second derivative (and the symbols we use also suggest this relation). This intuition is generally correct, but there are a few caveats to keep in mind.

First, for real-valued functions of one variable  $f : \mathbb{R} \rightarrow \mathbb{R}$ , it is a basic definition that the second derivative is the derivative of the first derivative, i.e.,

$$\frac{\partial^2 f(x)}{\partial x^2} = \frac{\partial}{\partial x} \frac{\partial}{\partial x} f(x).$$

However, for functions of a vector, the gradient of the function is a vector, and we cannot take the gradient of a vector – i.e.,

$$\nabla_x \nabla_x f(x) = \nabla_x \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix}$$

and this expression is not defined. Therefore, it is not the case that the Hessian is the gradient of the gradient. However, this is almost true, in the following sense: If we look at the  $i$ th entry of the gradient  $(\nabla_x f(x))_i = \partial f(x)/\partial x_i$ , and take the gradient with respect to  $x$  we get

$$\nabla_x \frac{\partial f(x)}{\partial x_i} = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_i \partial x_1} \\ \frac{\partial^2 f(x)}{\partial x_i \partial x_2} \\ \vdots \\ \frac{\partial^2 f(x)}{\partial x_i \partial x_n} \end{bmatrix}$$

which is the  $i$ th column (or row) of the Hessian. Therefore,

$$\nabla_x^2 f(x) = \begin{bmatrix} \nabla_x (\nabla_x f(x))_1 & \nabla_x (\nabla_x f(x))_2 & \cdots & \nabla_x (\nabla_x f(x))_n \end{bmatrix}$$

If we don't mind being a little bit sloppy we can say that (essentially)  $\nabla_x^2 f(x) = \nabla_x (\nabla_x f(x))^T$ , so long as we understand that this really means taking the gradient of each entry of  $(\nabla_x f(x))^T$ , not the gradient of the whole vector.

Finally, note that while we can take the gradient with respect to a matrix  $A \in \mathbb{R}^n$ , for the purposes of this class we will only consider taking the Hessian with respect to a vector  $x \in \mathbb{R}^n$ . This is simply a matter of convenience (and the fact that none of the calculations we do require us to find the Hessian with respect to a matrix), since the Hessian with respect to a matrix would have to represent all the partial derivatives  $\partial^2 f(A)/(\partial A_{ij} \partial A_{kl})$ , and it is rather cumbersome to represent this as a matrix.

### 3 Gradients and Hessians of Quadratic and Linear Functions

Now let's try to determine the gradient and Hessian matrices for a few simple functions. It should be noted that all the gradients given here are special cases of the gradients given in the CS229 lecture notes.

For  $x \in \mathbb{R}^n$ , let  $f(x) = b^T x$  for some known vector  $b \in \mathbb{R}^n$ . Then

$$f(x) = \sum_{i=1}^n b_i x_i$$

so

$$\frac{\partial f(x)}{\partial x_k} = \frac{\partial}{\partial x_k} \sum_{i=1}^n b_i x_i = b_k$$

From this we can easily see that  $\nabla_x b^T x = b$ . This should be compared to the analogous situation in single variable calculus, where  $\partial/(\partial x) ax = a$ .

Now consider the quadratic function  $f(x) = x^T A x$  for  $A \in \mathbb{S}^n$ . Remember that

$$f(x) = \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j$$

To take the partial derivative, we'll consider the terms including  $x_k$  and  $x_k^2$  factors separately:

$$\begin{aligned} \frac{\partial f(x)}{\partial x_k} &= \frac{\partial}{\partial x_k} \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j \\ &= \frac{\partial}{\partial x_k} \left[ \sum_{i \neq k} \sum_{j \neq k} A_{ij} x_i x_j + \sum_{i \neq k} A_{ik} x_i x_k + \sum_{j \neq k} A_{kj} x_k x_j + A_{kk} x_k^2 \right] \\ &= \sum_{i \neq k} A_{ik} x_i + \sum_{j \neq k} A_{kj} x_j + 2A_{kk} x_k \\ &= \sum_{i=1}^n A_{ik} x_i + \sum_{j=1}^n A_{kj} x_j = 2 \sum_{i=1}^n A_{ki} x_i \end{aligned}$$

where the last equality follows since  $A$  is symmetric (which we can safely assume, since it is appearing in a quadratic form). Note that the  $k$ th entry of  $\nabla_x f(x)$  is just the inner product of the  $k$ th row of  $A$  and  $x$ . Therefore,  $\nabla_x x^T A x = 2Ax$ . Again, this should remind you of the analogous fact in single-variable calculus, that  $\partial/(\partial x) ax^2 = 2ax$ .

Finally, let's look at the Hessian of the quadratic function  $f(x) = x^T A x$  (it should be obvious that the Hessian of a linear function  $b^T x$  is zero). In this case,

$$\frac{\partial^2 f(x)}{\partial x_k \partial x_\ell} = \frac{\partial}{\partial x_k} \left[ \frac{\partial f(x)}{\partial x_\ell} \right] = \frac{\partial}{\partial x_k} \left[ 2 \sum_{i=1}^n A_{\ell i} x_i \right] = 2A_{\ell k} = 2A_{k\ell}$$

Therefore, it should be clear that  $\nabla_x^2 x^T A x = 2A$ , which should be entirely expected (and again analogous to the single-variable fact that  $\partial^2/(\partial x^2) ax^2 = 2a$ ).

To recap,

- $\nabla_x b^T x = b$
- $\nabla_x x^T A x = 2Ax$  (if  $A$  symmetric)
- $\nabla_x^2 x^T A x = 2A$  (if  $A$  symmetric)

## 4 Least Square

Let's apply the equations we obtained in the last section to derive the least squares equations. Suppose we are given matrices  $A \in \mathbb{R}^{m \times n}$  (for simplicity we assume  $A$  is full rank) and a vector  $b \in \mathbb{R}^m$  such that  $b \notin \mathcal{R}(A)$ . In this situation we will not be able to find a vector  $x \in \mathbb{R}^n$ , such that  $Ax = b$ , so instead we want to find a vector  $x$  such that  $Ax$  is as close as possible to  $b$ , as measured by the square of the Euclidean norm  $\|Ax - b\|_2^2$ .

Using the fact that  $\|x\|_2^2 = x^T x$ , we have

$$\begin{aligned}\|Ax - b\|_2^2 &= (Ax - b)^T (Ax - b) \\ &= x^T A^T A x - 2b^T A x + b^T b\end{aligned}$$

Taking the gradient with respect to  $x$  we have, and using the properties we derived in the previous section

$$\begin{aligned}\nabla_x (x^T A^T A x - 2b^T A x + b^T b) &= \nabla_x x^T A^T A x - \nabla_x 2b^T A x + \nabla_x b^T b \\ &= 2A^T A x - 2A^T b\end{aligned}$$

Setting this last expression equal to zero and solving for  $x$  gives the normal equations

$$x = (A^T A)^{-1} A^T b$$

which is the same as what we derived in class.

## 5 Gradients of the Determinant

Now let's consider a situation where we find the gradient of a function with respect to a matrix, namely for  $A \in \mathbb{R}^{n \times n}$ , we want to find  $\nabla_A |A|$ . Recall from our discussion of determinants that

$$|A| = \sum_{i=1}^n (-1)^{i+j} A_{ij} |A_{\setminus i, \setminus j}| \quad (\text{for any } j \in 1, \dots, n)$$

so

$$\frac{\partial}{\partial A_{k\ell}} |A| = \frac{\partial}{\partial A_{k\ell}} \sum_{i=1}^n (-1)^{i+j} A_{ij} |A_{\setminus i, \setminus j}| = (-1)^{k+\ell} |A_{\setminus k, \setminus \ell}| = (\text{adj}(A))_{\ell k}$$

From this it immediately follows from the properties of the adjoint that

$$\nabla_A |A| = (\text{adj}(A))^T = |A| A^{-T}$$

Now let's consider the function  $f : \mathbb{S}_{++}^n \rightarrow \mathbb{R}, f(A) = \log|A|$ . Note that we have to restrict the domain of  $f$  to be the positive definite matrices, since this ensures that  $|A| > 0$ , so that the log of  $|A|$  is a real number. In this case we can use the chain rule (nothing fancy, just the ordinary chain rule from single-variable calculus) to see that

$$\frac{\partial \log|A|}{\partial A_{ij}} = \frac{\partial \log|A|}{\partial |A|} \frac{\partial |A|}{\partial A_{ij}} = \frac{1}{|A|} \frac{\partial |A|}{\partial A_{ij}}$$

From this it should be obvious that

$$\nabla_A \log|A| = \frac{1}{|A|} \nabla_A |A| = A^{-1}$$

where we can drop the transpose in the last expression because  $A$  is symmetric. Note the similarity to the single-valued case, where  $\partial/(\partial x) \log x = 1/x$ .

## 6 Eigenvalues as Optimization

Finally, we use matrix calculus to solve an optimization problem in a way that leads directly to eigenvalue/eigenvector analysis. Consider the following, equality constrained optimization problem:

$$\max_{x \in \mathbb{R}^n} x^T A x \quad \text{subject to } \|x\|_2^2 = 1$$

for a symmetric matrix  $A \in \mathbb{S}^n$ . A standard way of solving optimization problems with equality constraints is by forming the **Lagrangian**, an objective function that includes the equality constraints. The Lagrangian in this case can be given by

$$\mathcal{L}(x, \lambda) = x^T A x - \lambda x^T x$$

where  $\lambda$  is called the Lagrange multiplier associated with the equality constraint. It can be established that for  $x^*$  to be a optimal point to the problem, the gradient of the Lagrangian has to be zero at  $x^*$  (this is not the only condition, but it is required). That is,

$$\nabla_x \mathcal{L}(x, \lambda) = \nabla_x (x^T A x - \lambda x^T x) = 2A^T x - 2\lambda x = 0$$

Notice that this is just the linear equation  $Ax = \lambda x$ . This shows that the only points which can possibly maximize (or minimize)  $x^T A x$  assuming  $x^T x = 1$  are the eigenvectors of  $A$ .