

## 1 Vector Calculus

Below,  $\mathbf{x} \in \mathbb{R}^d$  means that  $\mathbf{x}$  is a  $d \times 1$  column vector with real-valued entries. Likewise,  $\mathbf{A} \in \mathbb{R}^{d \times d}$  means that  $\mathbf{A}$  is a  $d \times d$  matrix with real-valued entries. In this course, we will by convention consider vectors to be column vectors.

Consider  $\mathbf{x}, \mathbf{w} \in \mathbb{R}^d$  and  $\mathbf{A} \in \mathbb{R}^{d \times d}$ . In the following questions,  $\nabla_{\mathbf{x}}$  denotes the gradient with respect to  $\mathbf{x}$ , which, by convention, is a column vector. See the appendix for more details on definitions for vector calculus.

Calculate the following derivatives.

(a)  $\nabla_{\mathbf{x}}(\mathbf{w}^T \mathbf{x})$   
 $\nabla_{\mathbf{x}}(\mathbf{w}^T \mathbf{x}) = \left( \frac{\partial(\mathbf{w}^T \mathbf{x})}{\partial \mathbf{x}} \right)^T = \mathbf{w}$

(b)  $\nabla_{\mathbf{x}}(\mathbf{w}^T \mathbf{A} \mathbf{x})$   
 $\nabla_{\mathbf{x}}(\mathbf{w}^T \mathbf{A} \mathbf{x}) = \mathbf{A}^T \mathbf{w}$

(c)  $\nabla_{\mathbf{A}}(\mathbf{w}^T \mathbf{A} \mathbf{x})$   
 $\mathbf{w}^T \mathbf{A} \mathbf{x} = \sum_{i=1}^n \sum_{j=1}^m \mathbf{w}_i \mathbf{A}_{ij} \cdot \mathbf{x}_j$   
 $\nabla_{\mathbf{A}}(\mathbf{w}^T \mathbf{A} \mathbf{x}) = \mathbf{w} \mathbf{x}^T$

(d)  $\nabla_{\mathbf{x}}(\mathbf{x}^T \mathbf{A} \mathbf{x})$

$$\nabla_{\mathbf{x}}(\mathbf{x}^T \mathbf{A} \mathbf{x}) = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$$

(e)  $\nabla_{\mathbf{x}}^2(\mathbf{x}^T \mathbf{A} \mathbf{x})$

$$\nabla_{\mathbf{x}}(\mathbf{x}^T \mathbf{A} \mathbf{x}) = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$$

$$\begin{aligned} \nabla_{\mathbf{x}}^2(\mathbf{x}^T \mathbf{A} \mathbf{x}) &= \left( \frac{\partial (\mathbf{A} + \mathbf{A}^T) \mathbf{x}}{\partial \mathbf{x}} \right)^T \\ &= \mathbf{A} + \mathbf{A}^T \end{aligned}$$

(f) Now let's apply our identities derived above to a practical problem. Given a design matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and a label vector  $\mathbf{Y} \in \mathbb{R}^n$ , the ordinary least squares regression problem is

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{Y}\|_2^2$$

Using parts (a)–(e), derive a necessary condition for  $\mathbf{w}^*$ . *Note: We do not necessarily assume  $\mathbf{X}$  is full rank! Hint: A necessary condition for a minimum point of a function is that it is a critical point, i.e. where the gradient is 0.*

$$\begin{aligned} g(\mathbf{w}) &= \nabla_{\mathbf{w}} \left( \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{Y}\|_2^2 \right) = \left( \frac{\partial \left( \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{Y}\|_2^2 \right)}{\partial \mathbf{w}} \right)^T \\ &= \left( \frac{\partial \left( \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{Y}\|_2^2 \right)}{\partial (\mathbf{X}\mathbf{w} - \mathbf{Y})} \cdot \frac{\partial (\mathbf{X}\mathbf{w} - \mathbf{Y})}{\partial \mathbf{w}} \right)^T \\ &= \left( (\mathbf{X}\mathbf{w} - \mathbf{Y})^T \cdot \mathbf{X} \right)^T \\ &= \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{Y}) \\ g(\mathbf{w}^*) &= \mathbf{X}^T (\mathbf{X}\mathbf{w}^* - \mathbf{Y}) = 0 \quad \mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \quad \text{if } \mathbf{X}^T \mathbf{X} \text{ is invertible} \end{aligned}$$

## 2 Back to Basics: Linear Algebra

Let  $X \in \mathbb{R}^{n \times m}$ . When we write  $\subseteq$ , it means “is a subspace of.” We study a few important subspaces in the theory of linear maps:

- The **columnspace**, also called the range or span, of  $X$  is  $\text{Range}(X) := \{Xv : v \in \mathbb{R}^m\} \subseteq \mathbb{R}^n$ . Consists of all vectors in the span (the set of all linear combinations) of the columns of  $X$ .
- The **rowspace** is  $\text{Row}(X) := \{X^T v : v \in \mathbb{R}^n\} \subseteq \mathbb{R}^m$ . Consists of all vectors in the span of the rows of  $X$ .
- The **nullspace**, also called the kernel, of  $X$  is  $\mathcal{N}(X) := \{v \in \mathbb{R}^m : Xv = 0\} \subseteq \mathbb{R}^m$ .
- The **orthogonal complement** of a subspace  $U$  in some vector space  $V$  is a subspace, denoted  $U^\perp$ , such that  $u \in U, v \in U^\perp \implies u \cdot v = 0$  and  $U$  and  $U^\perp$  together span  $V$ . (These facts imply that  $\dim U + \dim U^\perp = \dim V$ . It also implies that  $U^{\perp\perp} = U$ .) For example, in the three-dimensional Euclidean space  $V = \mathbb{R}^3$ , if  $U$  is a plane through the origin, then  $U^\perp$  is a line through the origin perpendicular to  $U$ .

For this problem we do not assume that  $X$  has full rank.

(a) Show that the following facts are true.

(i)  $\text{Row}(X) = \text{Range}(X^T)$

$$\text{Range}(X^T) := \{X^T v : v \in \mathbb{R}^n\} \subseteq \mathbb{R}^m$$

$$\text{Row}(X) = \text{Range}(X^T)$$

(ii)  $\mathcal{N}(X)^\perp = \text{Row}(X)$ .

$$\forall v_1 \in \text{Row}(X), v_2 \in \mathcal{N}(X), v_1 \cdot v_2 = v_1^T v_2 = (Xv_0)^T v_2 = v_0^T (Xv_2) = 0 \implies v_1 \in \mathcal{N}(X)^\perp \implies \text{Row}(X) \subseteq \mathcal{N}(X)^\perp$$

According to the definition of Null Space,  $\text{Row}(X)^\perp = \mathcal{N}(X)$

$$\forall u \in \mathbb{R}^m, u = u_x + u_y, u_x \in \text{Row}(X), u_y \in \text{Row}(X)^\perp$$

$$\forall v \in \mathcal{N}(X)^\perp, v = v_x + v_y, v_x \in \text{Row}(X), v_y \in \text{Row}(X)^\perp$$

$$v \in \mathcal{N}(X)^\perp, v_y \in \mathcal{N}(X) \implies v \cdot v_y = 0 \implies (v_x + v_y) \cdot v_y = 0 \implies \|v_y\|_2^2 = 0 \implies v_y = 0$$

$$v_x \in \text{Row}(X), v_y \in \text{Row}(X)^\perp \implies v_x \cdot v_y = 0$$

$$v = v_x + v_y = v_x + 0 = v_x \in \text{Row}(X) \implies \mathcal{N}(X)^\perp \subseteq \text{Row}(X) \quad \text{Therefore, } \mathcal{N}(X)^\perp = \text{Row}(X)$$

(iii)  $\mathcal{N}(X^T X) = \mathcal{N}(X)$ . Hint: if  $v \in \mathcal{N}(X^T X)$ , then  $v^T X^T X v = 0$ .

$$\forall v_1 \in \mathcal{N}(X), Xv_1 = 0, X^T X v_1 = 0, v_1 \in \mathcal{N}(X^T X)$$

$$\forall v_2 \in \mathcal{N}(X^T X), X^T X v_2 = 0, v_2^T X^T X v_2 = 0, \|Xv_2\|_2^2 = 0, Xv_2 = 0$$

$$v_2 \in \mathcal{N}(X)$$

$$\mathcal{N}(X^T X) = \mathcal{N}(X)$$

- (b) We now prove an important result of linear algebra, the rank-nullity theorem. Let  $\text{Rank}(X) = \dim \text{Range}(X) = \dim \text{Row}(X)$  and  $\text{Nullity}(X) = \dim \mathcal{N}(X)$ . (The fact that  $\dim \text{Range}(X) = \dim \text{Row}(X)$ —that is, the dimension spanned by the rows equals the dimension spanned by the columns—is itself a pretty important result, which you should always remember when you hear the word “rank.”) The rank-nullity theorem says that for any  $X \in \mathbb{R}^{n \times m}$ ,

$$\text{Rank}(X) + \text{Nullity}(X) = m.$$

Use the above results to prove this theorem. *Hint: Use the orthogonal complement of the nullspace to connect the rank to the nullity.*

$$\dim \mathcal{N}(X) + \dim \mathcal{N}(X)^\perp = m$$

$$\text{Rank}(X) + \text{Nullity}(X) = m$$

Gilbert Strang has proposed that a collection of four facts be called “fundamental theorem of linear algebra.” Two of these facts are the rank-nullity theorem, part (b), and the fact that the row space is the orthogonal complement of the nullspace, part (a)(ii). The other two facts are related to the singular value decomposition, which we’ll learn late in the semester.

### 3 Eigenvalues

- (a) Let  $\mathbf{A}$  be an invertible matrix. Show that if  $\mathbf{v}$  is an eigenvector of  $\mathbf{A}$  with eigenvalue  $\lambda$ , then it is also an eigenvector of  $\mathbf{A}^{-1}$  with eigenvalue  $\lambda^{-1}$ .

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v} \Rightarrow \mathbf{A}^{-1}\mathbf{A}\mathbf{v} = \mathbf{A}^{-1}\lambda\mathbf{v} \Rightarrow \mathbf{v} = \lambda\mathbf{A}^{-1}\mathbf{v} \Rightarrow \lambda^{-1}\mathbf{v} = \mathbf{A}^{-1}\mathbf{v}$$

- (b) A symmetric matrix  $\mathbf{A}$  is said to be positive semidefinite (PSD) ( $\mathbf{A} \geq 0$ ) if  $\forall \mathbf{v} \neq 0, \mathbf{v}^T \mathbf{A} \mathbf{v} \geq 0$ . Show that  $\mathbf{A}$  is PSD if and only if all of its eigenvalues are nonnegative.

Hint: Use the eigendecomposition of the matrix  $\mathbf{A}$ .

$$\begin{aligned} \mathbf{A} &= \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T & \vdots & \mathbf{z}^T \mathbf{\Lambda} \mathbf{z} = \sum \lambda_i z_i^2 \geq 0 \\ \mathbf{v}^T \mathbf{A} \mathbf{v} &= \mathbf{v}^T \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T \mathbf{v} & \langle & \text{All } \lambda_i \geq 0 \\ \mathbf{z} &= \mathbf{Q}^T \mathbf{v} & \langle & \\ \mathbf{v}^T \mathbf{A} \mathbf{v} &= \mathbf{z}^T \mathbf{\Lambda} \mathbf{z} \geq 0 & \langle & \end{aligned}$$

- (c) Let  $\mathbf{A}$  be a PSD matrix. Show that its eigenvalues are equal to its singular values.

$$\begin{aligned} \mathbf{A} &= \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{-1} \\ \mathbf{A}^2 &= \mathbf{Q} \mathbf{\Lambda}^2 \mathbf{Q}^{-1} \\ \therefore \lambda_{\mathbf{A}} &\geq 0 \\ \therefore \lambda_{\mathbf{A}} &= \sqrt{\lambda_{\mathbf{A}^2}} \end{aligned}$$

## 4 Probability Review

There are  $n$  archers all shooting at the same target (bulls-eye) of radius 1. Let the score for a particular archer be defined to be the distance away from the center (the lower the score, the better, and 0 is the optimal score). Each archer's score is independent of the others, and is distributed uniformly between 0 and 1. What is the expected value of the worst (highest) score?

- (a) Define a random variable  $Z$  equal to the worst (highest) score, in terms of random variables that indicate each archer's score.

$$Z = \max\{X_i \mid i \in [1, n], i \in \mathbb{Z}\}$$

- (b) Derive the Cumulative Distribution Function (CDF) of  $Z$ . *Hint: Recall the CDF of a random variable  $Z$  is given by  $F(z) = P(Z \leq z)$*

$$F(z) = P(Z \leq z) = \underline{z^n}$$
$$\frac{dF(z)}{dz} = n z^{n-1}$$

- (c) Let  $X$  be a non-negative random variable. The Tail-Sum formula states that

$$\mathbb{E}[X] = \int_0^\infty \mathbb{P}(X \geq t) dt$$

Using both the Tail-Sum formula and the CDF of  $Z$  you derived, calculate the expected value of  $Z$ . *Hint: Write  $\mathbb{P}(X \geq t)$  in terms of the CDF of  $X$ .*

$$\begin{aligned} \mathbb{E}[X] &= \int_0^1 \mathbb{P}(X \geq t) dt = \int_0^1 (1 - t^n) dt \\ &= \left( t - \frac{1}{n+1} t^{n+1} \right) \Big|_0^1 = \frac{n}{n+1} \end{aligned}$$

(d) Consider what happens to  $\mathbb{E}[Z]$  as  $n \rightarrow \infty$ . Does this match your intuition?

Yes,  $n \rightarrow \infty$ , the max will more be  
more close to 1.

## 5 Vector Calculus Appendix <sup>1</sup>

Let us first understand the definition of the derivative. Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  denote a scalar function. Then the derivative  $\frac{\partial f}{\partial \mathbf{x}}$  is an operator that can help find the change in function value at  $\mathbf{x}$ , up to first order, when we add a little perturbation  $\Delta \in \mathbb{R}^d$  to  $\mathbf{x}$ . That is,

$$f(\mathbf{x} + \Delta) = f(\mathbf{x}) + \frac{\partial f}{\partial \mathbf{x}} \Delta + o(\|\Delta\|) \quad (1)$$

where  $o(\|\Delta\|)$  stands for any term  $r(\Delta)$  such that  $r(\Delta)/\|\Delta\| \rightarrow 0$  as  $\|\Delta\| \rightarrow 0$ . An example of such a term is a quadratic term like  $\|\Delta\|^2$ . Let us quickly verify that  $r(\Delta) = \|\Delta\|^2$  is indeed an  $o(\|\Delta\|)$  term. As  $\|\Delta\| \rightarrow 0$ , we have

$$\frac{r(\Delta)}{\|\Delta\|} = \frac{\|\Delta\|^2}{\|\Delta\|} = \|\Delta\| \rightarrow 0,$$

thereby verifying our claim. As a rule of thumb, any term that has a higher-order dependence on  $\|\Delta\|$  than linear is  $o(\|\Delta\|)$  and is ignored to compute the derivative.<sup>2</sup>

We call  $\frac{\partial f}{\partial \mathbf{x}}$  the *derivative of  $f$  at  $\mathbf{x}$* . Sometimes we use  $\frac{df}{dx}$  but we use  $\partial$  to indicate that  $f$  may depend on some other variable too. (But to define  $\frac{\partial f}{\partial \mathbf{x}}$ , we study changes in  $f$  with respect to changes in  $\mathbf{x}$  only.)

Since  $\Delta$  is a column vector, the vector  $\frac{\partial f}{\partial \mathbf{x}}$  should be a row vector so that  $\frac{\partial f}{\partial \mathbf{x}} \Delta$  is a scalar. So one way to compute the derivative is to expand out  $f(\mathbf{x} + \Delta)$  and guess the expression for the derivative. We call this method *computation via first principle*.

The gradient of  $f$  at  $\mathbf{x}$  is defined to be the transpose of this derivative. That is  $\nabla f(\mathbf{x}) = \left(\frac{\partial f}{\partial \mathbf{x}}\right)^\top$ .

We now write down some formulas that would be helpful to compute different derivatives in various settings where a solution via first principle might be hard to compute. We will also distinguish between the derivative, gradient, Jacobian, and Hessian in our notation.

1. Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  denote a scalar function. Let  $\mathbf{x} \in \mathbb{R}^d$  denote the vector input to  $f$ . We have

$$\frac{\partial f}{\partial \mathbf{x}} \in \mathbb{R}^{1 \times d} \quad \text{such that} \quad \frac{\partial f}{\partial \mathbf{x}} = \left[ \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_d} \right] \quad (2)$$

$$\nabla_{\mathbf{x}} f = \left( \frac{\partial f}{\partial \mathbf{x}} \right)^\top = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix}. \quad (3)$$

---

<sup>1</sup>Good resources for matrix calculus are:

- The Matrix Cookbook: <https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>
- Wikipedia: [https://en.wikipedia.org/wiki/Matrix\\_calculus](https://en.wikipedia.org/wiki/Matrix_calculus)
- Khan Academy: <https://www.khanacademy.org/math/multivariable-calculus/multivariable-derivatives>
- YouTube: <https://www.youtube.com/playlist?list=PLSQL0a2vh4HC5feHa6Rc5c0wbRTx56nF7>.

<sup>2</sup>Note that  $r(\Delta) = \sqrt{\|\Delta\|}$  is not an  $o(\|\Delta\|)$  term. Since for this case,  $r(\Delta)/\|\Delta\| = 1/\sqrt{\|\Delta\|} \rightarrow \infty$  as  $\|\Delta\| \rightarrow 0$ .



2. Let  $y : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$  be a scalar function defined on the space of  $m \times n$  matrices. Let  $\mathbf{B}$  denote the matrix input to  $f$ . Then, the derivative of  $f$  with respect to  $\mathbf{B}$  is an  $n \times m$  matrix given by

$$\frac{\partial y}{\partial \mathbf{B}} \in \mathbb{R}^{n \times m} \quad \text{such that} \quad \left[ \frac{\partial y}{\partial \mathbf{B}} \right]_{ij} = \frac{\partial y}{\partial B_{ji}}, \quad (4)$$

As in the vector case above, the gradient and derivative with respect to a matrix are also transposes of each other. So, the gradient of  $f$  with respect to  $\mathbf{B}$  is an  $m \times n$  matrix given by

$$\nabla_{\mathbf{B}} y = \left( \frac{\partial y}{\partial \mathbf{B}} \right)^\top \in \mathbb{R}^{m \times n} \quad \text{such that} \quad [\nabla_{\mathbf{B}} y]_{ij} = \frac{\partial y}{\partial B_{ij}}. \quad (5)$$

An argument via first principle follows as:

$$y(\mathbf{B} + \Delta) = y(\mathbf{B}) + \text{trace} \left( \frac{\partial y}{\partial \mathbf{B}} \Delta \right) + o(\|\Delta\|). \quad (6)$$

3. For  $\mathbf{z} : \mathbb{R}^d \rightarrow \mathbb{R}^k$ , a vector-valued function, its derivative  $\frac{\partial \mathbf{z}}{\partial \mathbf{x}}$  is an operator that can help find the change in function value at  $\mathbf{x}$ , up to first order, when we add a little perturbation  $\Delta$  to  $\mathbf{x}$ :

$$\mathbf{z}(\mathbf{x} + \Delta) = \mathbf{z}(\mathbf{x}) + \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \Delta + o(\|\Delta\|). \quad (7)$$

A formula for the same can be derived as

$$J(\mathbf{z}) = \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \in \mathbb{R}^{k \times d} = \begin{bmatrix} \frac{\partial z_1}{\partial \mathbf{x}} \\ \frac{\partial z_2}{\partial \mathbf{x}} \\ \vdots \\ \frac{\partial z_k}{\partial \mathbf{x}} \end{bmatrix} = \begin{bmatrix} \frac{\partial z_1}{\partial x_1} & \frac{\partial z_1}{\partial x_2} & \cdots & \frac{\partial z_1}{\partial x_d} \\ \frac{\partial z_2}{\partial x_1} & \frac{\partial z_2}{\partial x_2} & \cdots & \frac{\partial z_2}{\partial x_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial z_k}{\partial x_1} & \frac{\partial z_k}{\partial x_2} & \cdots & \frac{\partial z_k}{\partial x_d} \end{bmatrix} \quad (8)$$

That is,

$$[J(\mathbf{z})]_{ij} = \left[ \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right]_{ij} = \frac{\partial z_i}{\partial x_j}. \quad (9)$$

4. The Hessian of  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is the generalization of a second derivative and is defined as

$$H(f) = \nabla^2 f(\mathbf{x}) = J(\nabla f)^\top = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_d} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_d \partial x_1} & \frac{\partial^2 f}{\partial x_d \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_d^2} \end{bmatrix} \quad (10)$$

A first principle definition is

$$\nabla f(\mathbf{x} + \Delta) \approx \nabla f(\mathbf{x}) + \nabla^2 f(\mathbf{x}) \Delta \quad (11)$$

or equivalently

$$\nabla f(\mathbf{x} + \Delta) = \nabla f(\mathbf{x}) + \nabla^2 f(\mathbf{x}) \Delta + o(\|\Delta\|).$$

For sufficiently smooth functions (when the mixed derivatives are equal), the Hessian is a symmetric matrix. Most of the functions we cover in this class will have symmetric Hessians.