

Chapter 17

Appendix

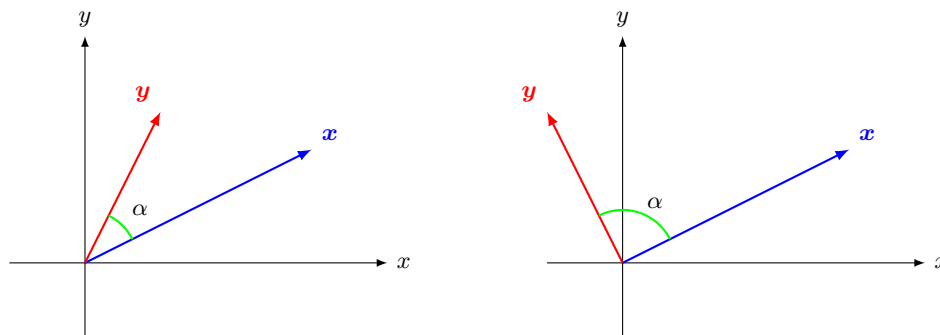
17.1 Review of Linear Algebra

For two vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, the **inner product** $\langle \mathbf{x}, \mathbf{y} \rangle$ of \mathbf{x} and \mathbf{y} is

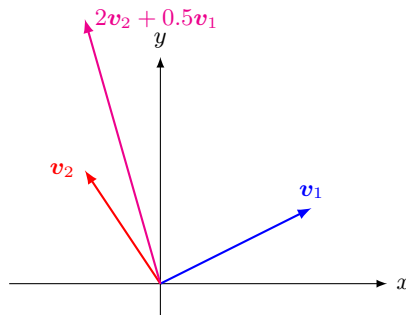
$$\mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \quad \langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y} = \sum_{i=1}^n x_i y_i. \quad (17.1)$$

where \mathbf{x}^T is the transpose of \mathbf{x} .

The **length** or the ℓ_2 norm of a vector \mathbf{x} is $\|\mathbf{x}\| = \|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^T \mathbf{x}}$ and we have $\|\mathbf{x}\|_2^2 = \mathbf{x}^T \mathbf{x}$. Let α be the angle between \mathbf{x} and \mathbf{y} . Then $\mathbf{x}^T \mathbf{y} = \|\mathbf{x}\| \|\mathbf{y}\| \cos \alpha$. If $\mathbf{x}^T \mathbf{y} = 0$, then the two are called **orthogonal**.



For a collection of vectors $\mathbf{v}_1, \dots, \mathbf{v}_m$, a **linear combination** of these is any vector of the form $a_1 \mathbf{v}_1 + \dots + a_m \mathbf{v}_m$, $a_i \in \mathbb{R}$. The set of all linear combinations of $\mathbf{v}_1, \dots, \mathbf{v}_m$ is their **span** and denoted as $\text{Span}\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$. This is a **subspace** (think line, plane, or the whole space). For a matrix \mathbf{A} , the span of the columns of \mathbf{A} is the **column space** of \mathbf{A} .

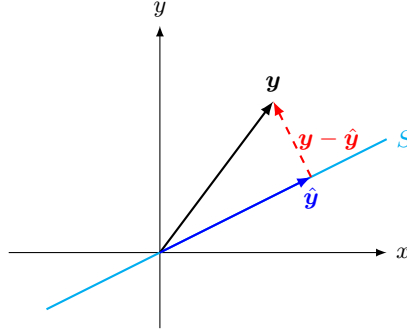


The vectors $\mathbf{v}_1, \dots, \mathbf{v}_m$ are **linearly independent** if there is no vector among them that can be written as a linear combination of the others, and linearly dependent otherwise. The vectors are linearly independent if and only if the only values for a_1, \dots, a_m satisfying $a_1\mathbf{v}_1 + \dots + a_m\mathbf{v}_m = \mathbf{0}$ are $a_1, \dots, a_m = 0$. In particular, the columns of a matrix \mathbf{A} are linearly independent if and only if the only vector \mathbf{a} satisfying $\mathbf{A}\mathbf{a} = \mathbf{0}$ is $\mathbf{a} = \mathbf{0}$.

The **inverse** of a square matrix \mathbf{A} is a matrix \mathbf{A}^{-1} such that $\mathbf{A}\mathbf{A}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$, where \mathbf{I} is the **identity matrix**, which has 1s on the diagonal and 0s elsewhere. A matrix that has an inverse is called **invertible**. For a square matrix \mathbf{A} , the following are equivalent:

- It is invertible.
- For all distinct vectors \mathbf{a} and \mathbf{b} , we have $\mathbf{A}\mathbf{a} \neq \mathbf{A}\mathbf{b}$.
- The only solution to $\mathbf{A}\mathbf{x} = \mathbf{0}$ is $\mathbf{x} = \mathbf{0}$.
- Its columns are linearly independent.
- Its determinant $|\mathbf{A}|$ is nonzero. (We also have $|\mathbf{A}^{-1}| = \frac{1}{|\mathbf{A}|}$.)

Given a subspace S (e.g., a plane or the column space of a matrix) and a vector \mathbf{y} , let $\hat{\mathbf{y}}$ be the vector in the subspace that is closest to \mathbf{y} . That is, we find $\hat{\mathbf{y}} \in S$ such that $\|\mathbf{y} - \hat{\mathbf{y}}\|$ is minimized. Then $\hat{\mathbf{y}}$ is called the **projection** of \mathbf{y} onto the subspace S .



Lemma 17.1 (Projection Lemma). *Let $\hat{\mathbf{y}}$ be the projection of a vector \mathbf{y} onto a subspace S . Then $\mathbf{y} - \hat{\mathbf{y}}$ is orthogonal to every vector in S .*

Proof. Suppose that this is not the case. Then there is a nonzero vector $\mathbf{v} \in S$ such that $(\mathbf{y} - \hat{\mathbf{y}})^T \mathbf{v} \neq 0$. We will show that this contradicts the minimality of $\|\mathbf{y} - \hat{\mathbf{y}}\|$. For any $a \in \mathbb{R}$,

$$\|\mathbf{y} - \hat{\mathbf{y}} - a\mathbf{v}\|_2^2 = (\mathbf{y} - \hat{\mathbf{y}} - a\mathbf{v})^T (\mathbf{y} - \hat{\mathbf{y}} - a\mathbf{v}) \quad (17.2)$$

$$= \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 - 2a\mathbf{v}^T (\mathbf{y} - \hat{\mathbf{y}}) + a^2\|\mathbf{v}\|_2^2. \quad (17.3)$$

This is a convex function in a . So setting the derivative to 0 gives the value of a that minimizes the error:

$$\frac{\partial}{\partial a} \|\mathbf{y} - \hat{\mathbf{y}} - a\mathbf{v}\|_2^2 = -2\mathbf{v}^T (\mathbf{y} - \hat{\mathbf{y}}) + 2a\|\mathbf{v}\|_2^2 = 0 \Rightarrow a = \frac{\mathbf{v}^T (\mathbf{y} - \hat{\mathbf{y}})}{\|\mathbf{v}\|_2^2} \neq 0. \quad (17.4)$$

Let

$$\hat{\mathbf{y}}' = \hat{\mathbf{y}} + \frac{\mathbf{v}^T (\mathbf{y} - \hat{\mathbf{y}})}{\mathbf{v}^T \mathbf{v}} \mathbf{v}, \quad (17.5)$$

and note that $\hat{\mathbf{y}}'$ is also in S but it is closer to \mathbf{y} , contradicting the optimality of $\hat{\mathbf{y}}$. \square

17.2 Vector and matrix differentiation

Definition 17.2 (The three derivatives). For a matrix \mathbf{A} , scalar z , and two vectors \mathbf{x}, \mathbf{y} (possibly one-dimensional), let

$$\frac{d\mathbf{A}}{dz} = \begin{pmatrix} \frac{\partial A_{11}}{\partial z} & \cdots & \frac{\partial A_{1n}}{\partial z} \\ \vdots & \ddots & \vdots \\ \frac{\partial A_{m1}}{\partial z} & \cdots & \frac{\partial A_{mn}}{\partial z} \end{pmatrix}, \quad \frac{dz}{d\mathbf{A}} = \begin{pmatrix} \frac{\partial z}{\partial A_{11}} & \cdots & \frac{\partial z}{\partial A_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial z}{\partial A_{m1}} & \cdots & \frac{\partial z}{\partial A_{mn}} \end{pmatrix}, \quad \frac{d\mathbf{y}}{d\mathbf{x}} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_n}{\partial x_1} & \cdots & \frac{\partial y_n}{\partial x_m} \end{pmatrix}$$

Lemma 17.3. For a scalar a , vectors $\mathbf{x}, \mathbf{y}, \mathbf{v}$, and constant matrices \mathbf{A} and \mathbf{S} ,

$$\begin{aligned} \frac{d\mathbf{y}}{d\mathbf{v}} &= \frac{d\mathbf{y}}{d\mathbf{x}} \frac{d\mathbf{x}}{d\mathbf{v}}, \\ \frac{d}{d\mathbf{v}}(a\mathbf{x}) &= a \frac{d\mathbf{x}}{d\mathbf{v}} + \mathbf{x} \frac{da}{d\mathbf{v}}, \\ \frac{d}{d\mathbf{v}}(\mathbf{y}^T \mathbf{A} \mathbf{x}) &= \mathbf{y}^T \mathbf{A} \frac{d\mathbf{x}}{d\mathbf{v}} + \mathbf{x}^T \mathbf{A}^T \frac{d\mathbf{y}}{d\mathbf{v}}, \\ \frac{d}{d\mathbf{v}}(\mathbf{y}^T \mathbf{S} \mathbf{y}) &= 2\mathbf{y}^T \mathbf{S} \frac{d\mathbf{y}}{d\mathbf{v}}, \quad (\mathbf{S} \text{ is symmetric}) \\ \frac{d}{d\mathbf{v}}(\mathbf{A} \mathbf{x}) &= \mathbf{A} \frac{d\mathbf{x}}{d\mathbf{v}}. \end{aligned}$$

Lemma 17.4. For matrix \mathbf{A} and constant vector \mathbf{x} ,

$$\begin{aligned} \frac{d}{d\mathbf{A}}(\mathbf{x}^T \mathbf{A} \mathbf{x}) &= \mathbf{x} \mathbf{x}^T \\ \frac{d}{d\mathbf{A}} \ln |\mathbf{A}| &= \mathbf{A}^{-T} \end{aligned}$$

Definition 17.5. Let $f : \mathbb{R}^m \rightarrow \mathbb{R}$. The gradient of $f(\mathbf{x})$ with respect to \mathbf{x} is defined as

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \left(\frac{df(\mathbf{x})}{d\mathbf{x}} \right)^T = \begin{pmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_m} \end{pmatrix}$$

and the Hessian of $f(\mathbf{x})$ with respect to \mathbf{x} is defined as

$$\mathbf{H}_{\mathbf{x}}(f(\mathbf{x})) = \frac{d\nabla_{\mathbf{x}} f(\mathbf{x})}{d\mathbf{x}} = \begin{pmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1 \partial x_1} & \cdots & \frac{\partial f(\mathbf{x})}{\partial x_m \partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_1 \partial x_m} & \cdots & \frac{\partial f(\mathbf{x})}{\partial x_m \partial x_m} \end{pmatrix}$$

Chain rule. Consider $h : \mathbb{R}^m \rightarrow \mathbb{R}$, $g : \mathbb{R} \rightarrow \mathbb{R}$, and $f(\mathbf{x}) = g(h(\mathbf{x}))$. From Lemma 17.3,

$$\begin{aligned} \nabla f(\mathbf{x}) &= g'(h(\mathbf{x})) \nabla h(\mathbf{x}), \\ \mathbf{H}f(\mathbf{x}) &= g'(h(\mathbf{x})) \mathbf{H}h(\mathbf{x}) + g''(h(\mathbf{x})) \nabla h(\mathbf{x}) \nabla^T h(\mathbf{x}) \end{aligned}$$

since

$$\begin{aligned} \mathbf{H}f(\mathbf{x}) &= \frac{d\nabla f}{d\mathbf{x}} \\ &= \frac{d(g'(h(\mathbf{x})) \nabla h(\mathbf{x}))}{d\mathbf{x}} \\ &= g'(h(\mathbf{x})) \frac{d\nabla h(\mathbf{x})}{d\mathbf{x}} + \nabla h(\mathbf{x}) \frac{d(g'(h(\mathbf{x})))}{d\mathbf{x}} \\ &= g'(h(\mathbf{x})) \mathbf{H}h(\mathbf{x}) + \nabla h(\mathbf{x}) \nabla^T h(\mathbf{x}) g''(h(\mathbf{x})) \end{aligned}$$

Example 17.6. Let us find the derivatives of $f(\mathbf{x}) = \log \sum_{i=1}^m e^{x_i}$. Let $\mathbf{z} = (\exp(x_i))_{i=1}^m$ so that $f(\mathbf{x}) = \log \mathbf{1}^T \mathbf{z}$.

$$\begin{aligned}\nabla f(\mathbf{x}) &= \frac{\mathbf{z}}{\mathbf{1}^T \mathbf{z}}, \\ \mathbf{H}f(\mathbf{x}) &= \frac{\text{diag}(\mathbf{z})}{\mathbf{1}^T \mathbf{z}} - \frac{\mathbf{z} \mathbf{z}^T}{(\mathbf{1}^T \mathbf{z})^2}.\end{aligned}$$

△

Chain rule. Let $\mathbf{h} = (h_1, \dots, h_n) : \mathbb{R}^m \rightarrow \mathbb{R}^n$, $g : \mathbb{R}^n \rightarrow \mathbb{R}$, and $f(\mathbf{x}) = g(\mathbf{h}(\mathbf{x}))$. Then

$$\begin{aligned}\frac{\partial f}{\partial x_i} &= \sum_{j=1}^n \frac{\partial g}{\partial h_j} \frac{\partial h_j}{\partial x_i} = \frac{dg}{d\mathbf{h}} \cdot \frac{d\mathbf{h}}{dx_i} = \nabla^T g \cdot \frac{d\mathbf{h}}{dx_i}, \\ \frac{df}{d\mathbf{x}} &= \frac{dg}{d\mathbf{h}} \frac{d\mathbf{h}}{d\mathbf{x}} = \nabla^T g \frac{d\mathbf{h}}{d\mathbf{x}}, \quad \nabla_{\mathbf{x}} f = \left(\frac{df}{d\mathbf{x}} \right)^T = \left(\frac{d\mathbf{h}}{d\mathbf{x}} \right)^T \nabla g\end{aligned}$$