

## Topic 2: Review of Linear Algebra

INSTRUCTOR: DANIEL L. PIMENTEL-ALARCÓN

© COPYRIGHT 2020

GO GREEN. AVOID PRINTING, OR PRINT 2-SIDED OR MULTIPAGE.

## 2.1 Fundamental Concepts

Linear algebra lies at the heart of machine learning, studying linear equations, mappings, and their representations in feature spaces. One of the most elemental vector manipulations are linear combinations, which essentially comprise scaling and addition.

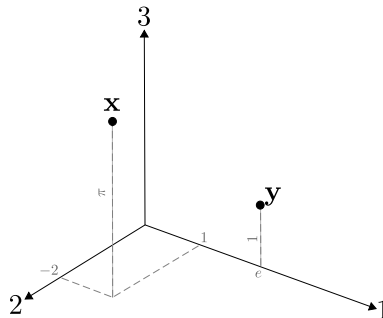
**Definition 2.1** (Linear combination, coefficients). A vector  $\mathbf{z}$  is a *linear combination* of vectors  $\{\mathbf{x}_1, \dots, \mathbf{x}_R\}$  if it can be written as

$$\mathbf{z} = \sum_{r=1}^R c_r \mathbf{x}_r \quad (2.1)$$

for some  $c_1, \dots, c_R \in \mathbb{R}$ . The scalars  $\{c_1, \dots, c_R\}$  are called the *coefficients* of  $\mathbf{z}$  with respect to (w.r.t.)  $\{\mathbf{x}_1, \dots, \mathbf{x}_R\}$ .

**Example 2.1.** Consider the following vectors in  $\mathbb{R}^3$ :

$$\mathbf{x} = \begin{bmatrix} 1 \\ -2 \\ \pi \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} e \\ 0 \\ 1 \end{bmatrix}.$$



Then

$$\mathbf{z} = -3\mathbf{x} + 2\mathbf{y} = -3 \begin{bmatrix} 1 \\ -2 \\ \pi \end{bmatrix} + 2 \begin{bmatrix} e \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} -3 \\ 6 \\ -3\pi \end{bmatrix} + \begin{bmatrix} 2e \\ 0 \\ 2 \end{bmatrix} = \begin{bmatrix} -3 + 2e \\ 6 \\ -3\pi + 2 \end{bmatrix}$$

is a linear combination of  $\mathbf{x}$  and  $\mathbf{y}$ , with coefficients  $-3$  and  $2$ .

Another fundamental concept is linear independence.

**Definition 2.2** (Linear independence). A set of vectors  $\{\mathbf{x}_1, \dots, \mathbf{x}_R\}$  is *linearly independent* if

$$\sum_{r=1}^R c_r \mathbf{x}_r = \mathbf{0}$$

implies  $c_r = 0$  for every  $r = 1, \dots, R$ . Otherwise we say it is *linearly dependent*.

Intuitively, a set of vectors is linearly independent if none of them can be written as a linear combination of the others.

**Example 2.2.** The following vectors are linearly independent:

$$\mathbf{x}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad \mathbf{x}_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \quad \mathbf{x}_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix},$$

because we cannot write either of them as a linear combination of the others. In contrast, the following vectors are linearly dependent:

$$\mathbf{y}_1 = \begin{bmatrix} 5 \\ 5 \\ 2 \end{bmatrix}, \quad \mathbf{y}_2 = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \quad \mathbf{y}_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix},$$

because we can write  $\mathbf{y}_1$  as a linear combination of  $\mathbf{y}_2$  and  $\mathbf{y}_3$ , namely,  $\mathbf{y}_1 = 5\mathbf{y}_2 + 2\mathbf{y}_3$ .

## 2.2 Matrices

Matrices are very handy structures to arrange and manipulate vectors.

**Example 2.3.** We can arrange the vectors in the Example 2.2 in the following matrices:

$$\mathbf{X} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{Y} = \begin{bmatrix} 5 & 5 & 2 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

### 2.2.1 Review of Basic Matrix Operations

- **Matrix Multiplication.** Given matrices  $\mathbf{A} \in \mathbb{R}^{m \times k}$  and  $\mathbf{B} \in \mathbb{R}^{k \times n}$ , their product, denoted by  $\mathbf{AB}$ , is another matrix of size  $m \times n$  whose  $(i, j)^{\text{th}}$  entry is given by:

$$[\mathbf{AB}]_{ij} = \sum_{\ell=1}^k \mathbf{A}_{i\ell} \mathbf{B}_{\ell j}.$$

Intuitively, the  $(i, j)^{\text{th}}$  entry of  $\mathbf{AB}$  is given by the multiplication of the  $i^{\text{th}}$  row of  $\mathbf{A}$  and the  $j^{\text{th}}$  column of  $\mathbf{B}$ . For example, if

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix}, \quad (2.2)$$

Then:

$$\mathbf{AB} = \begin{bmatrix} (1 \cdot 1 + 2 \cdot 3 + 3 \cdot 5) & (1 \cdot 2 + 2 \cdot 4 + 3 \cdot 6) \\ (4 \cdot 1 + 5 \cdot 3 + 6 \cdot 5) & (4 \cdot 2 + 5 \cdot 4 + 6 \cdot 6) \end{bmatrix} = \begin{bmatrix} 22 & 28 \\ 49 & 64 \end{bmatrix}.$$

- **Scalar Multiplication** Given a scalar  $c$  and a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , their product, denoted by  $c\mathbf{A}$ , is an  $m \times n$  matrix whose  $(i, j)^{\text{th}}$  entry is given by  $c$  times the  $(i, j)^{\text{th}}$  entry of  $\mathbf{A}$ . For example, with  $\mathbf{A}$  as in (2.2), and  $c = 7$ ,

$$c\mathbf{A} = 7\mathbf{A} = \begin{bmatrix} 7 & 14 & 21 \\ 28 & 35 & 42 \end{bmatrix},$$

- **Transposition.** The transpose of a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , denoted by  $\mathbf{A}^T$ , is an  $n \times m$  matrix whose  $(i, j)^{\text{th}}$  entry is given by the  $(j, i)^{\text{th}}$  entry of  $\mathbf{A}$ . Intuitively, transposing a matrix is like flipping its rows and columns along the diagonal. For example, with  $\mathbf{A}$  as in (2.2),

$$\mathbf{A}^T = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix}.$$

- **Trace.** Given a squared matrix  $\mathbf{X} \in \mathbb{R}^{m \times m}$ , its trace, denoted by  $\text{tr}(\mathbf{X})$  is the sum of its diagonal entries:

$$\text{tr}(\mathbf{X}) = \sum_{i=1}^m \mathbf{X}_{ii}.$$

For example, with  $\mathbf{X}$  as in Example 2.3,  $\text{tr}(\mathbf{X}) = 3$ .

- **Identity Matrix.** The identity matrix of size  $m \times m$ , denoted by  $\mathbf{I}_m$ , is a squared matrix whose diagonal entries are all ones, and off-diagonal entries are all zeros. For example, the  $3 \times 3$  identity matrix is

$$\mathbf{I}_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

Whenever there is no possible confusion about the size of the identity matrix, people often drop the  $m$  subindex, and simply denote it as  $\mathbf{I}$ .

- **Inverse.** Given a squared matrix  $\mathbf{X} \in \mathbb{R}^{m \times m}$ , its inverse, denoted by  $\mathbf{X}^{-1}$  is a matrix such that  $\mathbf{X}\mathbf{X}^{-1} = \mathbf{X}^{-1}\mathbf{X} = \mathbf{I}$ . For example, with  $\mathbf{X}$  as in Example 2.3,  $\mathbf{X}^{-1} = \mathbf{X}$ . As another example, consider

$$\mathbf{X} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \\ 1 & 3 & 9 \end{bmatrix}.$$

Then

$$\mathbf{X}^{-1} = \begin{bmatrix} 3 & -3 & 1 \\ -2.5 & 4 & -1.5 \\ 0.5 & -1 & 0.5 \end{bmatrix}.$$

You can verify that  $\mathbf{X}\mathbf{X}^{-1} = \mathbf{X}^{-1}\mathbf{X} = \mathbf{I}$ . Here is a special trick to invert  $2 \times 2$  matrices:

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}.$$

Of course, this requires that  $ad - bc \neq 0$ .

- **Hadamard Product.** Given two matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$ , their Hadamard product (aka point-wise product) is a matrix of size  $m \times n$ , denoted as  $\mathbf{A} \odot \mathbf{B}$ , whose  $(i,j)^{\text{th}}$  entry is given by the product of the  $(i,j)$  entries of  $\mathbf{A}$  and  $\mathbf{B}$ , i.e.,

$$[\mathbf{A} \odot \mathbf{B}]_{ij} = \mathbf{A}_{ij} \mathbf{B}_{ij}.$$

For example, with  $\mathbf{A}, \mathbf{B}$  as in (2.2),

$$\mathbf{A} \odot \mathbf{B}^T = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \odot \begin{bmatrix} 1 & 3 & 5 \\ 2 & 4 & 6 \end{bmatrix} = \begin{bmatrix} 1 & 6 & 15 \\ 8 & 20 & 36 \end{bmatrix}.$$

- **Vector Operations.** Notice that vectors are 1-column matrices, so all matrix operators that apply to non-squared matrices also apply to vectors. For instance, with the same setup as in Example 2.2,

$$\mathbf{x}_1^T \mathbf{x}_2 = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = 0.$$

Notice that  $\mathbf{x}_1^T \mathbf{x}_2 \neq \mathbf{x}_1 \mathbf{x}_2^T$ :

$$\mathbf{x}_1 \mathbf{x}_2^T = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}.$$

- **Norms.** Given a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , its Frobenius norm is defined as:

$$\|\mathbf{A}\|_F := \sqrt{\sum_{i=1}^m \sum_{j=1}^n \mathbf{A}_{ij}^2}.$$

In the particular case of vectors, this is often called the  $\ell_2$ -norm or Euclidean norm, and is denoted by  $\|\mathbf{x}\|_2$ , or simply  $\|\mathbf{x}\|$ . Norms are important because they essentially quantify the *size* of a matrix or vector. Just as there are several ways to quantify the size of a person (e.g., age, height, weight), there

are also several ways to quantify the *size* of a matrix or vector, for which we can use different norms. Another example is the  $\ell_1$ -norm, which for matrices is defined as:

$$\|\mathbf{A}\|_1 := \max_j \sum_{i=1} |\mathbf{A}_{ij}|,$$

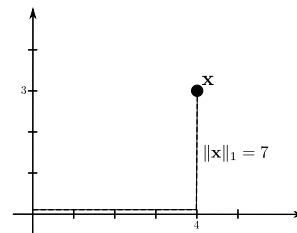
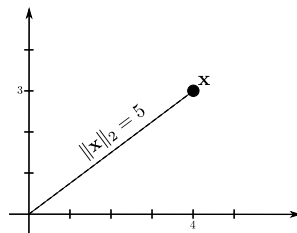
and for vectors is defined as

$$\|\mathbf{x}\|_1 := \sum_j |x_j|,$$

also known as the taxi-cab or Manhattan norm. Intuitively, the  $\ell_2$ -norm measures the *point-to-point* distance, while the  $\ell_1$ -norm measures the *taxi-cab* distance. For example, for the same vector  $\mathbf{x} = [4 \ 3]^\top$ , here are two different notions to quantify its size:

$$\|\mathbf{x}\|_2 = \sqrt{4^2 + 3^2} = 5$$

$$\|\mathbf{x}\|_1 = |4| + |3| = 7.$$



The norm  $\|\mathbf{x}\|$  of a vector  $\mathbf{x}$  is essentially its size. Norms are also useful because they allow us to measure distance between vectors (through their difference). For example, consider the following images:



and vectorize them to produce vectors  $\mathbf{x}, \mathbf{y}, \mathbf{z}$ . We want to do face clustering, i.e., we want to know which images correspond to the same person. If  $\|\mathbf{x} - \mathbf{y}\|$  is small (i.e.,  $\mathbf{x}$  is similar to  $\mathbf{y}$ ), it is reasonable to conclude that the first two images correspond to the same person. If  $\|\mathbf{x} - \mathbf{z}\|$  is large (i.e.,  $\mathbf{x}$  is very different from  $\mathbf{z}$ ), it is reasonable to conclude that the first and second images corresponds to different persons.

Norms satisfy the so-called *triangle inequality*:  $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ . This allows you to draw useful conclusions. For example, knowing that  $\|\mathbf{x} - \mathbf{y}\|$  is small and that  $\|\mathbf{x} - \mathbf{z}\|$  is large allows us to conclude that  $\|\mathbf{y} - \mathbf{z}\|$  is also large. Intuitively, this allows us to conclude that if  $\mathbf{x}, \mathbf{y}$  correspond to the same person, and  $\mathbf{x}$  and  $\mathbf{z}$  corresponds to different persons, then  $\mathbf{y}$  and  $\mathbf{z}$  also correspond to different persons. In other words, nothing weird will happen.

**Inner Product.** Given two matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$ , their inner product is defined as:

$$\langle \mathbf{A}, \mathbf{B} \rangle := \text{tr}(\mathbf{A}^\top \mathbf{B}).$$

For example, with  $\mathbf{A}$  as in (2.2),

$$\langle \mathbf{A}, \mathbf{A} \rangle = \text{tr}(\mathbf{A}^\top \mathbf{A}) = \text{tr} \left( \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \right) = \text{tr} \left( \begin{bmatrix} 17 & 22 & 27 \\ 22 & 29 & 36 \\ 27 & 36 & 45 \end{bmatrix} \right) = 91.$$

Notice that for vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^D$ ,  $\mathbf{x}^\top \mathbf{y}$  will always be a scalar, so we can drop the trace, and simply write  $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^\top \mathbf{y}$ . For example, with the same setup as in Example 2.2:

$$\begin{aligned} \langle \mathbf{x}_1, \mathbf{x}_2 \rangle &= \mathbf{x}_1^\top \mathbf{x}_2 = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = 0, \\ \langle \mathbf{y}_1, \mathbf{y}_2 \rangle &= \mathbf{y}_1^\top \mathbf{y}_2 = \begin{bmatrix} 5 & 5 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = 10. \end{aligned}$$

Inner products are of particular importance because they measure the similarity between matrices and vectors. In particular, recall from your kinder garden classes that the angle  $\theta$  between two vectors is given by:

$$\cos \theta = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|}.$$

More generally, the larger the inner product between two objects (in absolute value), the more similar they are.

### 2.2.2 Why do I care about Matrices?

Matrix operators are useful because they allow us to write otherwise complex burdensome operations in simple and concise matrix form.

**Example 2.4.** With the same settings as in Examples 2.2 and 2.3, we can write linear combinations using a simple matrix multiplication. For instance, instead of writing

$$\mathbf{y}_1 = \begin{bmatrix} 5 \\ 5 \\ 2 \end{bmatrix} = 5 \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + 5 \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} + 2 \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} = 5\mathbf{x}_1 + 5\mathbf{x}_2 + 2\mathbf{x}_3,$$

we can let  $\mathbf{c} = [5 \ 5 \ 2]^\top$  be  $\mathbf{y}_1$ 's coefficient vector, and equivalently write in simpler form:

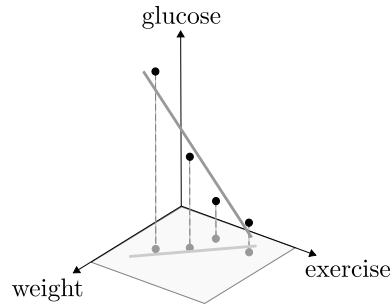
$$\mathbf{y}_1 = \begin{bmatrix} 5 \\ 5 \\ 2 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 5 \\ 5 \\ 2 \end{bmatrix} = \mathbf{X}^\top \mathbf{c}.$$

Similar matrix expressions will be ubiquitous throughout this course, so you should start getting familiar and feeling comfortable using matrix operations.

## 2.3 Subspaces

Subspaces are essentially high-dimensional lines. A 1-dimensional subspace is a line, a 2-dimensional subspace is a plane, and so on. Subspaces are useful because data often lies near subspaces.

**Example 2.5.** The following data lie near a 1-dimensional subspace (line):



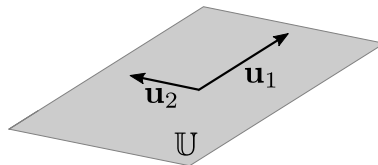
In higher dimensions subspaces may be harder to visualize, so you will have to use imagination to decide how a higher-dimensional subspace looks. Luckily, we have a precise and formal mathematical way to define them:

**Definition 2.3** (Subspace). A subset  $\mathbb{U} \subseteq \mathbb{R}^D$  is a *subspace* if for every  $a, b \in \mathbb{R}$  and every  $\mathbf{u}, \mathbf{v} \in \mathbb{U}$ ,  $a\mathbf{u} + b\mathbf{v} \in \mathbb{U}$ .

**Definition 2.4** (Span).  $\text{span}[\mathbf{u}_1, \dots, \mathbf{u}_R]$  is the set of all linear combinations of  $\{\mathbf{u}_1, \dots, \mathbf{u}_R\}$ . More formally,

$$\text{span}[\mathbf{u}_1, \dots, \mathbf{u}_R] := \left\{ \mathbf{x} \in \mathbb{R}^D : \mathbf{x} = \sum_{r=1}^R c_r \mathbf{u}_r \text{ for some } c_1, \dots, c_R \in \mathbb{R} \right\}.$$

For any vectors  $\mathbf{u}_1, \dots, \mathbf{u}_R \in \mathbb{R}^D$ ,  $\text{span}[\mathbf{u}_1, \dots, \mathbf{u}_R]$  is a subspace. For example, here is a subspace  $\mathbb{U}$  (plane) spanned by two vectors,  $\mathbf{u}_1$  and  $\mathbf{u}_2$ :



**Definition 2.5** (Basis). A set of linearly independent vectors  $\{\mathbf{u}_1, \dots, \mathbf{u}_R\}$  is a *basis* of a subspace  $\mathbb{U}$  if each  $\mathbf{v} \in \mathbb{U}$  can be written as

$$\mathbf{v} = \sum_{r=1}^R c_r \mathbf{u}_r$$

for a *unique* set of coefficients  $\{c_1, \dots, c_R\}$ .

**Definition 2.6** (Orthogonal). A collection of vectors  $\{\mathbf{x}_1, \dots, \mathbf{x}_R\}$  is *orthogonal* if  $\langle \mathbf{x}_i, \mathbf{x}_j \rangle = 0$  for every  $i \neq j$ .

If  $\|\mathbf{x}\| = 1$ , we say  $\mathbf{x}$  is a *unit* vector, or that it is *normalized*. Similarly, a collection of normalized, orthogonal vectors is called *orthonormal*. There is a tight relation between inner products and norms. The following is one of the most important and useful inequalities that describe this relationship.

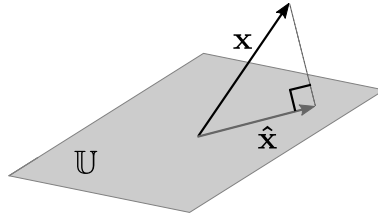
**Proposition 2.1** (Cauchy-Schwartz inequality). For every  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^D$ ,

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|.$$

Furthermore, if  $\mathbf{y} \neq \mathbf{0}$ , then equality holds if and only if  $\mathbf{x} = a\mathbf{y}$  for some  $a \in \mathbb{R}$ .

## 2.4 Projections

In words, the projection  $\hat{\mathbf{x}}$  of a vector  $\mathbf{x}$  onto a subspace  $\mathbb{U}$  is the vector in  $\mathbb{U}$  that is closest to  $\mathbf{x}$ :



More formally,

**Definition 2.7** (Projection). The *projection* of  $\mathbf{x} \in \mathbb{R}^D$  onto subspace  $\mathbb{U}$  is the vector  $\hat{\mathbf{x}} \in \mathbb{U}$  satisfying

$$\|\mathbf{x} - \hat{\mathbf{x}}\| \leq \|\mathbf{x} - \mathbf{u}\| \quad \text{for every } \mathbf{u} \in \mathbb{U}.$$

Notice that if  $\mathbf{x} \in \mathbb{U}$ , then  $\hat{\mathbf{x}} = \mathbf{x}$ . The following proposition tells us exactly how to compute projections.

**Proposition 2.2.** Let  $\{\mathbf{u}_1, \dots, \mathbf{u}_R\}$  be an orthonormal basis of  $\mathbb{U}$ . The projection of  $\mathbf{x} \in \mathbb{R}^D$  onto  $\mathbb{U}$  is given by

$$\hat{\mathbf{x}} = \sum_{r=1}^R \langle \mathbf{x}, \mathbf{u}_r \rangle \mathbf{u}_r.$$

In other words, the coefficient of  $\mathbf{x}$  w.r.t.  $\mathbf{u}_r$  is given by  $\langle \mathbf{x}, \mathbf{u}_r \rangle$ .



Furthermore, the following proposition tells us that we can compute projections very efficiently: just using a simple matrix multiplication! This makes projections very attractive in practice. For example, as we saw before, data often lies near subspaces. We can measure how close using the norm of the *residual*  $\mathbf{x} - \hat{\mathbf{x}}$ .

**Proposition 2.3** (Projector operator). Let  $\mathbf{U} \in \mathbb{R}^{D \times R}$  be a basis of  $\mathbb{U}$ . The *projection operator*  $\mathbf{P}_{\mathbf{U}} : \mathbb{R}^D \rightarrow \mathbb{U}$  that maps any vector  $\mathbf{x} \in \mathbb{R}^D$  to its projection  $\hat{\mathbf{x}} \in \mathbb{U}$  is given by:

$$\mathbf{P}_{\mathbf{U}} = \mathbf{U}(\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T.$$

Notice that if  $\mathbf{U}$  is orthonormal, then  $\mathbf{P}_{\mathbf{U}} = \mathbf{U} \mathbf{U}^T$ .

*Proof.* Since  $\hat{\mathbf{x}} \in \mathbb{U}$ , that means we can write  $\hat{\mathbf{x}}$  as  $\mathbf{U}\mathbf{c}$  for some  $\mathbf{c} \in \mathbb{R}^R$ . We thus want to find the  $\mathbf{c}$  that minimizes:

$$\|\mathbf{x} - \mathbf{U}\mathbf{c}\|_2^2 = (\mathbf{x} - \mathbf{U}\mathbf{c})^T (\mathbf{x} - \mathbf{U}\mathbf{c}) = \mathbf{x}^T - 2\mathbf{c}^T \mathbf{U}^T \mathbf{x} + \mathbf{c}^T \mathbf{U}^T \mathbf{U} \mathbf{c}.$$

Since this is convex in  $\mathbf{c}$ , we can use elemental optimization to find the desired minimizer, i.e., we will take derivative w.r.t.  $\mathbf{c}$ , set to zero and solve for  $\mathbf{c}$ . To learn more about how to take derivatives w.r.t. vectors and matrices see *Old and new matrix algebra useful for statistics* by Thomas P. Minka. The derivative w.r.t.  $\mathbf{c}$  is given by:

$$-2\mathbf{U}^T \mathbf{x} + 2\mathbf{U}^T \mathbf{U} \mathbf{c}.$$

Setting to zero and solving for  $\mathbf{c}$  we obtain:

$$\hat{\mathbf{c}} := \arg \min_{\mathbf{c} \in \mathbb{R}^R} \|\mathbf{x} - \mathbf{U}\mathbf{c}\|_2^2 = (\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T \mathbf{x},$$

where we know  $\mathbf{U}^T \mathbf{U}$  is invertible because  $\mathbf{U}$  is a basis by assumption, so its columns are linearly independent. It follows that

$$\hat{\mathbf{x}} = \mathbf{U}\hat{\mathbf{c}} = \underbrace{\mathbf{U}(\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T}_{\mathbf{P}_{\mathbf{U}}} \mathbf{x},$$

as claimed. If  $\mathbf{U}$  is orthonormal, then  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ , and hence  $\mathbf{P}_{\mathbf{U}}$  simplifies to  $\mathbf{U} \mathbf{U}^T$ . Notice that  $\hat{\mathbf{c}}$  are the coefficients of  $\hat{\mathbf{x}}$  w.r.t. the basis  $\mathbf{U}$ . □

## 2.5 Gram-Schmidt Orthogonalization

Orthonormal bases have very nice and useful properties. For example, in Proposition 2.3, if the basis  $\mathbf{U}$  is orthonormal, then the projection operator is simplified into  $\mathbf{U} \mathbf{U}^T$ , which requires much less computations than  $\mathbf{U}(\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T$ . The following procedure tells us exactly how to transform an arbitrary basis into an orthonormal basis.

**Proposition 2.4** (Gram-Schmidt procedure). Let  $\{\mathbf{u}_1, \dots, \mathbf{u}_R\}$  be a basis of  $\mathbb{U}$ . Let

$$\begin{aligned} \mathbf{v}'_r &= \begin{cases} \mathbf{u}_1 & r = 1, \\ \mathbf{u}_r - \sum_{k=1}^{r-1} \langle \mathbf{u}_r, \mathbf{v}_k \rangle \mathbf{v}_k & r = 2, \dots, R, \end{cases} \\ \mathbf{v}_r &= \mathbf{v}'_r / \|\mathbf{v}'_r\|. \end{aligned}$$

Then  $\{\mathbf{v}_1, \dots, \mathbf{v}_R\}$  are an orthonormal basis of  $\mathbb{U}$ .

*Proof.* We know from Proposition 2.2 that  $\sum_{k=1}^{r-1} \langle \mathbf{u}_r, \mathbf{v}_k \rangle \mathbf{v}_k$  is the projection of  $\mathbf{u}_r$  onto  $\text{span}[\mathbf{v}_1, \dots, \mathbf{v}_{r-1}]$ . This implies  $\mathbf{v}'_r$  is the orthogonal residual of  $\mathbf{u}_r$  onto  $\text{span}[\mathbf{v}_1, \dots, \mathbf{v}_{r-1}]$ , and hence it is orthogonal to  $\{\mathbf{v}_1, \dots, \mathbf{v}_{r-1}\}$ , as desired.  $\mathbf{v}_r$  is simply the normalized version of  $\mathbf{v}'_r$ .  $\square$

## 2.6 Conclusions

These notes review several fundamental concepts of linear algebra that lie at the heart of machine learning, and will be crucial in *every* topic to be studied in this course.