

## Linear Algebra Review

**Vectors** To begin, let us describe an element of the state space as a point with numerical coordinates, that is

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

Vectors of up to three dimensions are easy to diagram. For example,

$$\mathbf{x} = \begin{pmatrix} 3 \\ 2 \\ 5 \end{pmatrix}$$

can be drawn as follows.

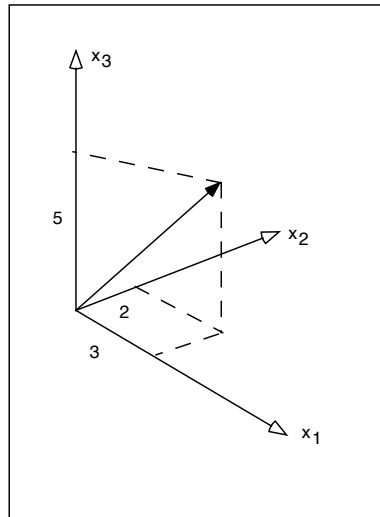


Figure 1: An example of a vector.

**Addition** To add two vectors together simply add their components. To multiply a vector by a scalar (number), multiply each of the components by the scalar. For example, if  $\mathbf{z} = \mathbf{x} + \mathbf{y}$ , then if

$$\mathbf{x} = \begin{pmatrix} 3 \\ 5 \\ 2 \end{pmatrix} \text{ and } \mathbf{y} = \begin{pmatrix} -4 \\ 1 \\ 3 \end{pmatrix} \text{ then } \mathbf{z} = \begin{pmatrix} -1 \\ 6 \\ 5 \end{pmatrix}$$

And if  $\mathbf{z} = \alpha \mathbf{x}$  for a scalar  $\alpha$ , then

$$\mathbf{z} = \alpha \begin{pmatrix} 3 \\ 5 \\ 2 \end{pmatrix} = \begin{pmatrix} 3\alpha \\ 5\alpha \\ 2\alpha \end{pmatrix}$$

**Dot Product** The *dot product* of two vectors, denoted  $\mathbf{x} \cdot \mathbf{y}$ , is defined as the sum of the product of their pairwise components; that is,

$$\mathbf{x} \cdot \mathbf{y} = \sum_{i=1}^n x_i y_i$$

For our example,  $\mathbf{x} \cdot \mathbf{y} = (3)(-4) + (5)(1) + (2)(3) = -1$ .

Two vectors are said to be *orthogonal* if their dot product is zero; that is,  $\mathbf{x} \cdot \mathbf{y} = 0$ .

The *length* of a vector, denoted  $\|\mathbf{x}\|$ , is simply  $\sqrt{\mathbf{x} \cdot \mathbf{x}}$ . The angle  $\theta$  between two vectors  $\mathbf{x}$  and  $\mathbf{y}$  is given by

$$\cos \theta = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}$$

From this equation it is seen that if vectors are orthogonal, then  $\cos \theta = 0$  or  $\theta = 90^\circ$ .

The vector  $(\mathbf{x} \cdot \mathbf{y}) \frac{\mathbf{y}}{\|\mathbf{y}\|}$  is termed the *projection* of  $\mathbf{x}$  onto  $\mathbf{y}$ .

The *transpose* of a matrix, denoted  $A^T$ , is simply described using the element notation as  $\{a_{ji}\}$ . In other words, the elements are “flipped” about the diagonal. A square  $n \times n$  matrix is *symmetric* if  $A^T = A$ .

**Linear Transformations** For any function  $f(x)$ , a linear transformation is such that

$$f(ax + by) = af(x) + bf(y)$$

An important linear transformation is *matrix multiplication*. Matrix multiplication  $A = BC$  is defined by

$$a_{ij} = \sum_{k=1}^N b_{ik}c_{kj}, i = 1, \dots, P, j = 1, \dots, Q$$

From this formula it is seen that the number of columns of  $B$  has to be the same as the number of rows of  $C$  for multiplication to be defined.

**Determinant** To define the *determinant* of a matrix first requires defining the number of inversions in a number sequence. Consider the sequence  $\{1, 3, 4, 2\}$ . The number of inversions in this sequence is 2 because 3 and 4 come after 2. Similarly the number of inversions in  $\{4, 2, 1, 3\}$  is 3. Denote the number of inversions of a sequence as  $n$ . The determinant of a matrix  $A$ , denoted  $|A|$ , is the sum of all  $n!$  possible different products that compose elements from columns of the matrix with a term that depends on the number of inversions in the row indices; that is,

$$|A| = \sum (-1)^{n(i_1, i_2, \dots, i_n)} a_{i_1, 1} a_{i_2, 2} \cdots a_{i_n, n}$$

Like the inverse of a matrix, the determinant is expensive to calculate for large matrices, and a standard text

should be referred to for an algorithm. For practice calculations, however, it is useful to remember that the determinant of the  $2 \times 2$  matrix

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

is given by

$$|A| = ad - bc$$

**Inverse** For square matrices where  $N = M$ , an important matrix is the *inverse* matrix  $A^{-1}$ , which is defined by

$$AA^{-1} = I$$

where  $I$  is the *identity matrix*

$$I = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \cdots & & 1 \end{bmatrix}$$

In general, like the determinant, inverses take some work to calculate, and you should find a numerical routine. For practice, however, it is useful to remember that the inverse of the  $2 \times 2$  matrix

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

is given by

$$A^{-1} = \frac{\begin{bmatrix} d & -b \\ -c & a \end{bmatrix}}{|A|}$$

**Trace** The *trace* of a matrix  $A$  is the sum of its diagonal elements; that is,

$$\text{Tr}(A) = \sum_{i=1}^N a_{ii}$$

**Positive Definite** A matrix  $A$  is *positive definite* if for every  $\mathbf{x}$ ,

$$\mathbf{x}^T A \mathbf{x} > 0$$

and *positive semidefinite* if

$$\mathbf{x}^T A \mathbf{x} \geq 0$$

**Orthonormal Transformation** A transformation matrix  $A$  is orthonormal when

$$A^{-1} = A^T$$

As a consequence

$$AA^T = I$$

## 1 Coordinate Systems

Let us start by considering how coordinate systems represent multi-dimensional data. Multiplying a matrix by a vector is a special case of matrix multiplication where

$$\mathbf{y} = A\mathbf{x}$$

This can be written as  $y_i = \sum_{k=1}^N a_{kj}x_j, i = 1, \dots, M$ . Alternatively we can see the transformation as a *linear combination* of the columns of  $A$ :

$$\mathbf{y} = \mathbf{a}_1x_1 + \mathbf{a}_2x_2 + \dots + \mathbf{a}_Nx_N$$

Often in manipulating vectors it is implicitly assumed that they are described with respect to an orthogonal coordinate system. Hence the actual coordinate vectors are not discussed. In the general case, however, the right coordinate system for data might not be orthogonal. To develop this point, consider first that the vectors  $\mathbf{a}_i$  have a special interpretation as a coordinate system or *basis* for a multidimensional space. For example, in the traditional basis in three dimensions,

$$\mathbf{a}_1 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \mathbf{a}_2 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \text{ and } \mathbf{a}_3 = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

allows  $\mathbf{y}$  to be written as

$$\mathbf{y} = \mathbf{a}_1y_1 + \mathbf{a}_2y_2 + \mathbf{a}_3y_3$$

A fundamentally important property of coordinate systems is that they are only describable *relative* to one another. For example,  $\mathbf{y}$  is described in terms of the basis vectors  $\mathbf{a}_i$ .

This basis is orthogonal, since

$$\mathbf{a}_i \cdot \mathbf{a}_j = 0$$

for all  $i$  and  $j$  such that  $i \neq j$ , but it turns out that a nonorthogonal basis would also work. For example,

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -1 & 1 & -1 \end{bmatrix}$$

would still allow  $\mathbf{y}$  to be represented (although the coefficients would of course be different). However, the matrix

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -1 & 1 & 0 \end{bmatrix}$$

would not work. In this case the reason is easy to see: there is no way of representing the third component. In general, to represent  $n$ -dimensional vectors, the basis must *span* the space. A general condition for this is that the columns of  $A$  must be *linearly independent*. Formally this means that the only way you could write

$$\mathbf{a}_1x_1 + \mathbf{a}_2x_2 + \cdots + \mathbf{a}_Nx_N = \mathbf{0}$$

would be if  $x_i = 0$  for all  $i$ .

What happens when the columns of an  $n$ -dimensional matrix do not span the space? The dimension of the matrix is equal to the number of linearly independent

vectors, which is known as the *rank* of the matrix. When the rank  $r$  is less than the dimension  $N$ , the vectors are said to span an  $r$ -dimensional *subspace*.

In some cases it is desirable for a matrix to have less than full rank. For example, for the equation

$$A\mathbf{x} = \mathbf{0}$$

to have a nontrivial solution, the columns of  $A$  must be linearly dependent. Why? This equation is just a rewritten version of the previous equation. If the columns were linearly independent, then the only way the equations could be satisfied would be to have  $x_i = 0$  for all  $i$ . But for a nontrivial solution the  $x_i$  should be nonzero. Hence for this to happen the columns must be linearly dependent. For example, in three dimensions this can happen when all three of the vectors are in a plane.

In contrast, for the equation

$$A\mathbf{x} = \mathbf{c}$$

to have a unique solution the converse is true, because in order to have a solution, now the vector  $\mathbf{c}$  must be expressed in terms of a linear combination of the columns  $\mathbf{a}_i$ . For this statement to be generally true the columns must span the space of  $\mathbf{c}$ ; hence, together with the vector  $\mathbf{c}$ , they must be linearly dependent. These two cases are illustrated in Figure 2.



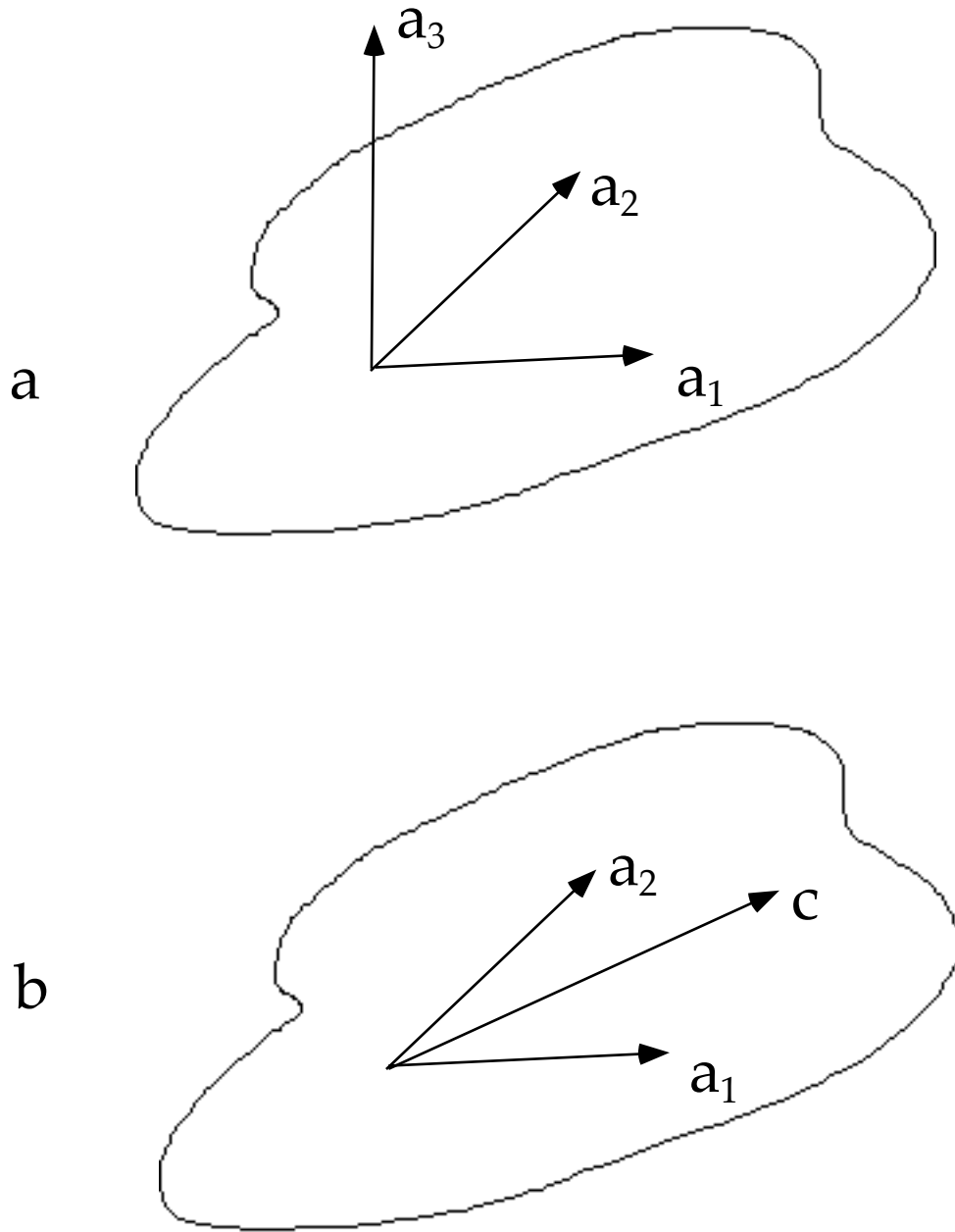


Figure 2: (a) The columns of  $A$  are linearly independent in three dimensions if the three column vectors are not coplanar. (b) For the case  $A\mathbf{x} = \mathbf{c}$  the vector

## Eigenvalues and Eigenvectors

At this point you should be used to the idea that any matrix can be thought of as representing a coordinate system. When a vector is multiplied by such a matrix, the general result is that the magnitude and direction of the resultant vector are different from the original. However, there is a very important special case. For any matrix, there are vector directions such that the matrix multiplication only changes the magnitude of the vector, leaving the direction unchanged. For these special directions, matrix multiplication reduces to scalar multiplication. The following example shows a case for the matrix

$$\begin{bmatrix} 3 & 1 \\ 2 & 2 \end{bmatrix}$$

where

$$\begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

is a special direction for the matrix, since multiplying it by the matrix just results in scaling the vector by a factor  $\lambda = 4$ ; that is,

$$\begin{bmatrix} 3 & 1 \\ 2 & 2 \end{bmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = 4 \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

In the general case, if a vector  $\mathbf{v}$  lies along one of these directions,

$$W\mathbf{v} = \lambda\mathbf{v}$$

where  $\lambda$  is a scalar. Vectors that lie along these special directions are known as *eigenvectors*, and the scalars associated with a transformation matrix are known as

*eigenvalues*. Finding the eigenvalues of an  $n \times n$  matrix for arbitrary  $n$  requires a trip to the recipe book, starting with the solution of an  $n^{\text{th}}$ -order polynomial to find the eigenvalues, but it is useful to work it out for an easy two-dimensional case, as follows:

$$\begin{bmatrix} 3 & 1 \\ 2 & 2 \end{bmatrix} \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} = \lambda \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$$

or, in other words,

$$\begin{bmatrix} 3 - \lambda & 1 \\ 2 & 2 - \lambda \end{bmatrix} \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

From the previous section we know that for this equation to have a solution, the columns of the matrix must be linearly dependent, and thus  $|W| = 0$ . Thus

$$(3 - \lambda)(2 - \lambda) - 2 = 0$$

which can be solved to find the two eigenvalues  $\lambda_1 = 4$  and  $\lambda_2 = 1$ . Now for the eigenvectors. Substituting  $\lambda_1 = 4$  into the equation results in

$$\begin{bmatrix} -1 & 1 \\ 2 & -2 \end{bmatrix} \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

Now this set of equations is degenerate, meaning that there is only one useful equation in two unknowns. As a consequence there is an infinity of solutions, and you must pick one arbitrarily. Pick  $v_1 = 1$ . Then  $v_2 = 1$ . Thus the eigenvector associated with  $\lambda_1 = 4$  is

$$\begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

As an easy exercise you can find the eigenvector associated with  $\lambda_2 = 1$ .

Now, for any particular matrix, why not pick the eigenvectors as the basis? It turns out that this is a good thing to do, since the effect is to transform the matrix into another matrix whose only nonzero elements are on the diagonal. Furthermore, these diagonal elements are the eigenvalues. The effect is to reduce matrix multiplication in the old basis to scalar multiplication in the new basis.

## Changing Coordinates

What happens to transformations when the coordinate basis is changed? Suppose that the coordinate transformation is given by

$$\mathbf{x}^* = A\mathbf{x}$$

$$\mathbf{y}^* = A\mathbf{y}$$

Given the transformation

$$\mathbf{y} = W\mathbf{x}$$

what happens to  $W$  when the coordinate system is changed to the starred system? That is, for some  $W^*$  it will be true that

$$\mathbf{y}^* = W^*\mathbf{x}^*$$

What is the relation between  $W$  and  $W^*$ ? One way to find out is to change back to the original system, transform by  $W$ , and then transform back to the starred system; that is,

$$\mathbf{x} = A^{-1}\mathbf{x}^*$$

$$\mathbf{y} = W\mathbf{x}$$

$$\mathbf{y}^* = A\mathbf{y}$$

Putting these transformations together:

$$\mathbf{y}^* = AWA^{-1}\mathbf{x}^*$$

Since the vector transformation taken by the two different routes should be the same, it must be true that

$$W^* = AWA^{-1}$$

Matrices related in this way are called *similar*.

## Eigenvalue Transformations

Now let's relate this discussion to eigenvectors. Suppose that the eigenvectors have been chosen as the basis set. Then for a given eigenvector  $\mathbf{y}_i$ ,

$$W\mathbf{y}_i = \lambda\mathbf{y}_i$$

and if  $Y$  is a matrix whose columns are the eigenvectors  $\mathbf{y}_i$ , then

$$WY = Y\Lambda$$

Here  $\Lambda$  is a matrix whose only nonzero components are the diagonal elements  $\lambda_i$ . Premultiplying both sides by  $Y^{-1}$ ,

$$Y^{-1}WY = \Lambda$$

What this equation means is that given a matrix  $W$ , the transformation it defines can always be simplified to that of a matrix whose only nonzero elements are diagonal by transforming to coordinates that use its eigenvectors as a basis. Furthermore, those elements are the eigenvalues.

### Example

To check this result, let us use the earlier example where

$$W = \begin{bmatrix} 3 & 1 \\ 2 & 2 \end{bmatrix}$$

and

$$Y = \begin{bmatrix} 1 & 1 \\ 1 & -2 \end{bmatrix}$$

and

$$\Lambda = \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix}$$

Let us pick a particular vector as  $\mathbf{x}$

$$\mathbf{x} = \begin{pmatrix} 3 \\ 4 \end{pmatrix}$$

so that  $\mathbf{y}$  is given by

$$\mathbf{y} = \begin{pmatrix} 13 \\ 14 \end{pmatrix}$$

First note that

$$Y^{-1} = \frac{1}{3} \begin{bmatrix} 2 & 1 \\ 1 & -1 \end{bmatrix}$$

so that  $\mathbf{x}^* = Y^{-1}\mathbf{x}$ , which is

$$\begin{pmatrix} \frac{10}{3} \\ -\frac{1}{3} \end{pmatrix}$$

and  $\mathbf{y}^* = Y^{-1}\mathbf{y}$ , which is

$$\begin{pmatrix} \frac{40}{3} \\ -\frac{1}{3} \end{pmatrix}$$

But from this you can see that  $\mathbf{y}^* = \Lambda\mathbf{x}^*$ .

## Properties of Eigenvalues

1. An eigenvalue matrix  $\Lambda$  is invariant under any orthogonal transformation.
2. If all its eigenvalues are positive, a matrix  $A$  is positive definite.
3. The trace of  $A$  is the sum of all its eigenvalues and is invariant under any orthogonal transformation.
4. The trace of  $A^m$  is the sum of all its eigenvalues and is invariant under any orthogonal transformation.
5. The determinant of  $A$  is equal to the product of all its eigenvalues and is invariant under any orthogonal transformation.



## Random Vectors

In this case vectors are drawn from some random distribution that captures the natural variations in the world. A random vector  $\mathbf{X}$  is specified by a probability density function  $p(\mathbf{X})$ , where formally

$$p(\mathbf{X}) = \lim_{\Delta \mathbf{x}_i \rightarrow 0} \frac{P(\mathbf{X} \in \mathbf{I})}{\prod_i \Delta x_i}$$

where

$$\mathbf{I} = \{\mathbf{X} : x_i < X_i \leq x_i + \Delta x_i, \forall i\}$$

Although a random vector is fully characterized by its density function, such functions are often difficult to determine or mathematically complex to use. These limitations motivate modeling distributions with functions that can be described by a low number of parameters. The most important of such parameters are the *mean vector* and *covariance matrix*, which are just generalizations of the mean and variance of scalar random variables to vector random variables.

The mean vector is defined by

$$\mathbf{M} = E\{\mathbf{X}\} = \int \mathbf{X} p(\mathbf{X}) d\mathbf{X}$$

and the covariance matrix by

$$\Sigma = E\{(\mathbf{X} - \mathbf{M})(\mathbf{X} - \mathbf{M})^T\}$$

In practice, with real data you will use the sample mean vector and sample covariance matrix. Where  $\mathbf{X}^k$ ,  $k = 1, N$  are the samples,

$$M = \frac{1}{N} \sum_{k=1}^N \mathbf{X}^k$$

$$\Sigma = \frac{1}{N} \sum_{k=1}^N (\mathbf{X}^k - M)(\mathbf{X}^k - M)^T$$

You will need more than three samples in practice, but to illustrate the calculations, suppose

$$\mathbf{X}^1 = \begin{pmatrix} -1 \\ 3 \\ 1 \end{pmatrix}, \mathbf{X}^2 = \begin{pmatrix} 2 \\ 1 \\ -1 \end{pmatrix}, \mathbf{X}^3 = \begin{pmatrix} 2 \\ 2 \\ 3 \end{pmatrix}$$

Then the mean value is

$$M = \frac{1}{3} \begin{pmatrix} 3 \\ 6 \\ 3 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \\ 1 \end{pmatrix}$$

So that

$$\mathbf{X}^1 - M = \begin{pmatrix} -2 \\ 1 \\ 0 \end{pmatrix}, \mathbf{X}^2 - M = \begin{pmatrix} 1 \\ -1 \\ -2 \end{pmatrix}, \mathbf{X}^3 - M = \begin{pmatrix} 1 \\ 0 \\ 2 \end{pmatrix}$$

and the covariance matrix is given by

$$\begin{aligned} \Sigma &= \frac{1}{3} \left\{ \begin{bmatrix} 4 & -2 & 0 \\ -2 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} + \begin{bmatrix} 1 & -1 & -2 \\ -1 & 1 & 2 \\ -2 & 2 & 4 \end{bmatrix} + \begin{bmatrix} 1 & 0 & 2 \\ 0 & 0 & 0 \\ 2 & 0 & 4 \end{bmatrix} \right\} \\ &= \frac{1}{3} \begin{bmatrix} 6 & -3 & 0 \\ -3 & 2 & 2 \\ -2 & 2 & 8 \end{bmatrix} \end{aligned}$$

## High-Dimensional Spaces

Suppose, as is often the case, that the dimension of the space is extremely large. Now the standard way to proceed would be to choose eigenvectors  $u_k$  and eigenvalues  $\lambda_k$  of the sample covariance matrix  $\Sigma$  where

$$\begin{aligned}\Sigma &= \frac{1}{M} \sum_{n=1}^M \mathbf{X}_n \mathbf{X}_n^T \\ &= A A^T\end{aligned}$$

where

$$A = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_M]$$

an  $M \times N$  matrix of  $M$  data samples. The problem with this tack is that it is infeasible owing to the high dimensionality of the matrix  $\Sigma$ . Since for an image the dimension of  $\mathbf{X}$  is  $n^2$ , then the dimension of  $\Sigma$  is  $n^2 \times n^2$ . For typical values of  $n$ , say 256, this is impossibly large. Salvation comes from the fact that the matrix  $\Sigma$  may be approximated by a matrix of lower rank. That is, most of the variation can be captured by projecting the data onto a subspace whose dimension is much less than the dimension of the space.

Rather than finding the eigenvectors of the larger system, consider finding the eigenvectors of the  $M \times M$  system

$$A^T A \mathbf{v} = \mu \mathbf{v} \tag{1}$$

Premultiplying both sides by  $A$ ,

$$A A^T A \mathbf{v} = \mu A \mathbf{v}$$

What this equation shows is that if  $\mathbf{v}$  is an eigenvector of  $A^T A$ , then  $A \mathbf{v}$  is an eigenvector of  $\Sigma$ . Furthermore, the eigenvalues of the smaller system are the same as those of the much larger system. It turns out also that these are the  $M$  largest eigenvalues. So to find the eigenvectors of the larger system, first find the eigenvalues and eigenvectors of the smaller system, and then multiply the eigenvectors by  $A$ .

**Example: Face Recognition** A lovely example of representing data is from face recognition. The task is to recognize the identity of an image of a face. The face image is described by an  $N \times N$  array of brightness values. Given  $M$  exemplars of images with known identities, the objective is to take a new image and identify it with the training image to which it is most similar. The key element is in the similarity metric. It should be chosen to score the essential variations in the data. From the last section, the way to discover the essential variations is with principal components analysis, which identifies the eigenvalues of the covariance matrix of all the data.

To begin, identify the training set of images as  $\mathbf{I}_1, \mathbf{I}_2, \dots, \mathbf{I}_M$ . These are shown in Figure 3. To work with these it is useful to subtract the bias introduced, as all the brightness levels are positive. Thus first identify the “average face” (shown in Figure ??)

$$\mathbf{I}_{ave} = \frac{1}{M} \sum_{n=1}^M \mathbf{I}_n$$

and then convert the training set by subtracting the average,

$$\mathbf{X}_i = \mathbf{I}_i - \mathbf{I}_{ave}, i = 1, \dots, M$$

Now use Equation 1 to find  $M$  eigenvectors  $v_k$  and eigenvalues  $\lambda_k$ .

From the “short” eigenvectors (of length  $M$ ), the larger eigenvectors  $\mathbf{u}_k$ , termed eigenfaces, can be constructed using  $\mathbf{v}_k$ , as follows:

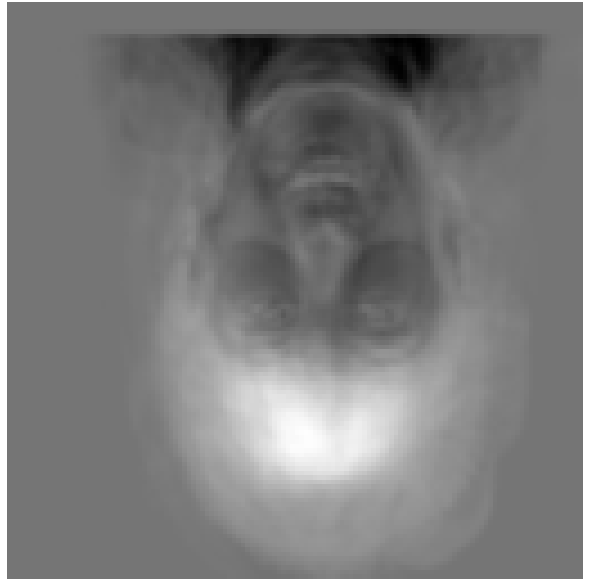
$$\mathbf{u}_i = \sum_{k=1}^M v_{ik} \mathbf{X}_k$$



Figure 3: A database of 12 figures used to calculate eigenvectors.

Figure ?? shows the first seven eigenvectors calculated from Equation 1.









Now that the direction principal variations have been calculated in the space of faces, this information can be used to classify a new face in terms of the faces in the data set. To do so, compute the coordinates of the new image in the eigenvector space  $\mathbf{\Omega} = (\omega_1, \omega_2, \dots, \omega_M)$  as follows:

$$\omega_k = u_k^T (\mathbf{I} - \mathbf{I}_{ave}), k = 1, \dots, M$$

Next compare  $\mathbf{\Omega}$  to the  $\mathbf{\Omega}$ s for each of the classes to pick the closest; that is, pick the class  $k$  that minimizes

$$\|\mathbf{\Omega} - \mathbf{\Omega}_k\|$$