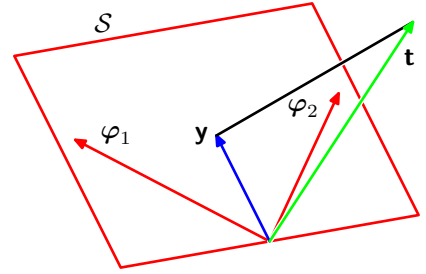


**Figure 4.3** Geometrical interpretation of the least-squares solution in an  $N$ -dimensional space whose axes are the values of  $t_1, \dots, t_N$ . The least-squares regression function is obtained by finding the orthogonal projection of the data vector  $\mathbf{t}$  onto the subspace spanned by the basis functions  $\phi_j(\mathbf{x})$  in which each basis function is viewed as a vector  $\varphi_j$  of length  $N$  with elements  $\phi_j(\mathbf{x}_n)$ .



#### 4.1.4 Geometry of least squares

At this point, it is instructive to consider the geometrical interpretation of the least-squares solution. To do this, we consider an  $N$ -dimensional space whose axes are given by the  $t_n$ , so that  $\mathbf{t} = (t_1, \dots, t_N)^T$  is a vector in this space. Each basis function  $\phi_j(\mathbf{x}_n)$ , evaluated at the  $N$  data points, can also be represented as a vector in the same space, denoted by  $\varphi_j$ , as illustrated in Figure 4.3. Note that  $\varphi_j$  corresponds to the  $j$ th column of  $\Phi$ , whereas  $\phi(\mathbf{x}_n)$  corresponds to the transpose of the  $n$ th row of  $\Phi$ . If the number  $M$  of basis functions is smaller than the number  $N$  of data points, then the  $M$  vectors  $\phi_j(\mathbf{x}_n)$  will span a linear subspace  $S$  of dimensionality  $M$ . We define  $\mathbf{y}$  to be an  $N$ -dimensional vector whose  $n$ th element is given by  $y(\mathbf{x}_n, \mathbf{w})$ , where  $n = 1, \dots, N$ . Because  $\mathbf{y}$  is an arbitrary linear combination of the vectors  $\varphi_j$ , it can live anywhere in the  $M$ -dimensional subspace. The sum-of-squares error (4.11) is then equal (up to a factor of  $1/2$ ) to the squared Euclidean distance between  $\mathbf{y}$  and  $\mathbf{t}$ . Thus, the least-squares solution for  $\mathbf{w}$  corresponds to that choice of  $\mathbf{y}$  that lies in subspace  $S$  and is closest to  $\mathbf{t}$ . Intuitively, from Figure 4.3, we anticipate that this solution corresponds to the orthogonal projection of  $\mathbf{t}$  onto the subspace  $S$ . This is indeed the case, as can easily be verified by noting that the solution for  $\mathbf{y}$  is given by  $\Phi \mathbf{w}_{\text{ML}}$  and then confirming that this takes the form of an orthogonal projection.

In practice, a direct solution of the normal equations can lead to numerical difficulties when  $\Phi^T \Phi$  is close to singular. In particular, when two or more of the basis vectors  $\varphi_j$  are co-linear, or nearly so, the resulting parameter values can have large magnitudes. Such near degeneracies will not be uncommon when dealing with real data sets. The resulting numerical difficulties can be addressed using the technique of *singular value decomposition*, or *SVD* (Deisenroth, Faisal, and Ong, 2020). Note that the addition of a regularization term ensures that the matrix is non-singular, even in the presence of degeneracies.

#### 4.1.5 Sequential learning

The maximum likelihood solution (4.14) involves processing the entire training set in one go and is known as a *batch* method. This can become computationally costly for large data sets. If the data set is sufficiently large, it may be worthwhile to use *sequential* algorithms, also known as *online* algorithms, in which the data points are considered one at a time and the model parameters updated after each such presentation. Sequential learning is also appropriate for real-time applications in which the data observations arrive in a continuous stream and predictions must be

#### Exercise 4.4