

Exercise 5.5.17 Let $A = \begin{bmatrix} 0 & a & b \\ a & 0 & c \\ b & c & 0 \end{bmatrix}$ and

$$B = \begin{bmatrix} c & a & b \\ a & b & c \\ b & c & a \end{bmatrix}.$$

- Show that $x^3 - (a^2 + b^2 + c^2)x - 2abc$ has real roots by considering A .
- Show that $a^2 + b^2 + c^2 \geq ab + ac + bc$ by considering B .

Exercise 5.5.18 Assume the 2×2 matrix A is similar to an upper triangular matrix. If $\operatorname{tr} A = 0 = \operatorname{tr} A^2$, show that $A^2 = 0$.

Exercise 5.5.19 Show that A is similar to A^T for all 2×2 matrices A . [Hint: Let $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$. If $c = 0$ treat the cases $b = 0$ and $b \neq 0$ separately. If $c \neq 0$, reduce to the case $c = 1$ using Exercise 5.5.12(d).]

Exercise 5.5.20 Refer to Section 3.4 on linear recurrences. Assume that the sequence x_0, x_1, x_2, \dots satisfies

$$x_{n+k} = r_0x_n + r_1x_{n+1} + \cdots + r_{k-1}x_{n+k-1}$$

for all $n \geq 0$. Define

$$A = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \\ r_0 & r_1 & r_2 & \cdots & r_{k-1} \end{bmatrix}, V_n = \begin{bmatrix} x_n \\ x_{n+1} \\ \vdots \\ x_{n+k-1} \end{bmatrix}.$$

Then show that:

- $V_n = A^n V_0$ for all n .
- $c_A(x) = x^k - r_{k-1}x^{k-1} - \cdots - r_1x - r_0$
- If λ is an eigenvalue of A , the eigenspace E_λ has dimension 1, and $\mathbf{x} = (1, \lambda, \lambda^2, \dots, \lambda^{k-1})^T$ is an eigenvector. [Hint: Use $c_A(\lambda) = 0$ to show that $E_\lambda = \mathbb{R}\mathbf{x}$.]
- A is diagonalizable if and only if the eigenvalues of A are distinct. [Hint: See part (c) and Theorem 5.5.4.]
- If $\lambda_1, \lambda_2, \dots, \lambda_k$ are distinct real eigenvalues, there exist constants t_1, t_2, \dots, t_k such that $x_n = t_1\lambda_1^n + \cdots + t_k\lambda_k^n$ holds for all n . [Hint: If D is diagonal with $\lambda_1, \lambda_2, \dots, \lambda_k$ as the main diagonal entries, show that $A^n = PD^nP^{-1}$ has entries that are linear combinations of $\lambda_1^n, \lambda_2^n, \dots, \lambda_k^n$.]

Exercise 5.5.21 Suppose A is 2×2 and $A^2 = 0$. If $\operatorname{tr} A \neq 0$ show that $A = 0$.

5.6 Best Approximation and Least Squares

Often an exact solution to a problem in applied mathematics is difficult to obtain. However, it is usually just as useful to find arbitrarily close approximations to a solution. In particular, finding “linear approximations” is a potent technique in applied mathematics. One basic case is the situation where a system of linear equations has no solution, and it is desirable to find a “best approximation” to a solution to the system. In this section best approximations are defined and a method for finding them is described. The result is then applied to “least squares” approximation of data.

Suppose A is an $m \times n$ matrix and \mathbf{b} is a column in \mathbb{R}^m , and consider the system

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

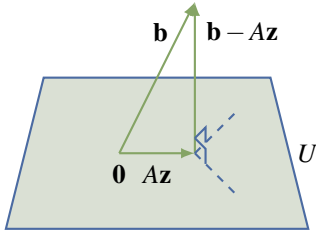
of m linear equations in n variables. This need not have a solution. However, given any column $\mathbf{z} \in \mathbb{R}^n$, the distance $\|\mathbf{b} - A\mathbf{z}\|$ is a measure of how far $A\mathbf{z}$ is from \mathbf{b} . Hence it is natural to ask whether there is a column \mathbf{z} in \mathbb{R}^n that is as close as possible to a solution in the sense that

$$\|\mathbf{b} - A\mathbf{z}\|$$

is the minimum value of $\|\mathbf{b} - \mathbf{Ax}\|$ as \mathbf{x} ranges over all columns in \mathbb{R}^n .

The answer is “yes”, and to describe it define

$$U = \{\mathbf{Ax} \mid \mathbf{x} \text{ lies in } \mathbb{R}^n\}$$



This is a subspace of \mathbb{R}^n (verify) and we want a vector \mathbf{Az} in U as close as possible to \mathbf{b} . That there is such a vector is clear geometrically if $n = 3$ by the diagram. In general such a vector \mathbf{Az} exists by a general result called the *projection theorem* that will be proved in Chapter 8 (Theorem 8.1.3). Moreover, the projection theorem gives a simple way to compute \mathbf{z} because it also shows that the vector $\mathbf{b} - \mathbf{Az}$ is *orthogonal* to every vector \mathbf{Ax} in U . Thus, for all \mathbf{x} in \mathbb{R}^n ,

$$\begin{aligned} 0 &= (\mathbf{Ax}) \cdot (\mathbf{b} - \mathbf{Az}) = (\mathbf{Ax})^T (\mathbf{b} - \mathbf{Az}) = \mathbf{x}^T A^T (\mathbf{b} - \mathbf{Az}) \\ &= \mathbf{x} \cdot [A^T (\mathbf{b} - \mathbf{Az})] \end{aligned}$$

In other words, the vector $A^T (\mathbf{b} - \mathbf{Az})$ in \mathbb{R}^n is orthogonal to *every* vector in \mathbb{R}^n and so must be zero (being orthogonal to itself). Hence \mathbf{z} satisfies

$$(A^T A)\mathbf{z} = A^T \mathbf{b}$$

Definition 5.14 Normal Equations

This is a system of linear equations called the **normal equations** for \mathbf{z} .

Note that this system can have more than one solution (see Exercise 5.6.5). However, the $n \times n$ matrix $A^T A$ is invertible if (and only if) the columns of A are linearly independent (Theorem 5.4.3); so, in this case, \mathbf{z} is uniquely determined and is given explicitly by $\mathbf{z} = (A^T A)^{-1} A^T \mathbf{b}$. However, the most efficient way to find \mathbf{z} is to apply gaussian elimination to the normal equations.

This discussion is summarized in the following theorem.

Theorem 5.6.1: Best Approximation Theorem

Let A be an $m \times n$ matrix, let \mathbf{b} be any column in \mathbb{R}^m , and consider the system

$$\mathbf{Ax} = \mathbf{b}$$

of m equations in n variables.

1. Any solution \mathbf{z} to the normal equations

$$(A^T A)\mathbf{z} = A^T \mathbf{b}$$

is a best approximation to a solution to $\mathbf{Ax} = \mathbf{b}$ in the sense that $\|\mathbf{b} - \mathbf{Az}\|$ is the minimum value of $\|\mathbf{b} - \mathbf{Ax}\|$ as \mathbf{x} ranges over all columns in \mathbb{R}^n .

2. If the columns of A are linearly independent, then $A^T A$ is invertible and \mathbf{z} is given uniquely by $\mathbf{z} = (A^T A)^{-1} A^T \mathbf{b}$.

We note in passing that if A is $n \times n$ and invertible, then

$$\mathbf{z} = (A^T A)^{-1} A^T \mathbf{b} = A^{-1} \mathbf{b}$$

is the solution to the system of equations, and $\|\mathbf{b} - A\mathbf{z}\| = 0$. Hence if A has independent columns, then $(A^T A)^{-1} A^T$ is playing the role of the inverse of the nonsquare matrix A . The matrix $A^T (A A^T)^{-1}$ plays a similar role when the rows of A are linearly independent. These are both special cases of the **generalized inverse** of a matrix A (see Exercise 5.6.14). However, we shall not pursue this topic here.

Example 5.6.1

The system of linear equations

$$3x - y = 4$$

$$x + 2y = 0$$

$$2x + y = 1$$

has no solution. Find the vector $\mathbf{z} = \begin{bmatrix} x_0 \\ y_0 \end{bmatrix}$ that best approximates a solution.

Solution. In this case,

$$A = \begin{bmatrix} 3 & -1 \\ 1 & 2 \\ 2 & 1 \end{bmatrix}, \text{ so } A^T A = \begin{bmatrix} 3 & 1 & 2 \\ -1 & 2 & 1 \end{bmatrix} \begin{bmatrix} 3 & -1 \\ 1 & 2 \\ 2 & 1 \end{bmatrix} = \begin{bmatrix} 14 & 1 \\ 1 & 6 \end{bmatrix}$$

is invertible. The normal equations $(A^T A)\mathbf{z} = A^T \mathbf{b}$ are

$$\begin{bmatrix} 14 & 1 \\ 1 & 6 \end{bmatrix} \mathbf{z} = \begin{bmatrix} 14 \\ -3 \end{bmatrix}, \text{ so } \mathbf{z} = \frac{1}{83} \begin{bmatrix} 87 \\ -56 \end{bmatrix}$$

Thus $x_0 = \frac{87}{83}$ and $y_0 = \frac{-56}{83}$. With these values of x and y , the left sides of the equations are, approximately,

$$3x_0 - y_0 = \frac{317}{83} = 3.82$$

$$x_0 + 2y_0 = \frac{-25}{83} = -0.30$$

$$2x_0 + y_0 = \frac{118}{83} = 1.42$$

This is as close as possible to a solution.

Example 5.6.2

The average number g of goals per game scored by a hockey player seems to be related linearly to two factors: the number x_1 of years of experience and the number x_2 of goals in the preceding 10 games. The data on the following page were collected on four players. Find the linear function

$g = a_0 + a_1x_1 + a_2x_2$ that best fits these data.

g	x_1	x_2
0.8	5	3
0.8	3	4
0.6	1	5
0.4	2	1

Solution. If the relationship is given by $g = r_0 + r_1x_1 + r_2x_2$, then the data can be described as follows:

$$\begin{bmatrix} 1 & 5 & 3 \\ 1 & 3 & 4 \\ 1 & 1 & 5 \\ 1 & 2 & 1 \end{bmatrix} \begin{bmatrix} r_0 \\ r_1 \\ r_2 \end{bmatrix} = \begin{bmatrix} 0.8 \\ 0.8 \\ 0.6 \\ 0.4 \end{bmatrix}$$

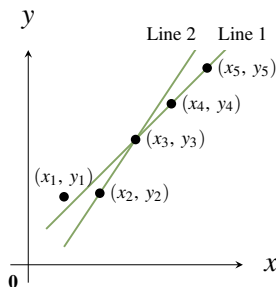
Using the notation in Theorem 5.6.1, we get

$$\begin{aligned} \mathbf{z} &= (A^T A)^{-1} A^T \mathbf{b} \\ &= \frac{1}{42} \begin{bmatrix} 119 & -17 & -19 \\ -17 & 5 & 1 \\ -19 & 1 & 5 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 5 & 3 & 1 & 2 \\ 3 & 4 & 5 & 1 \end{bmatrix} \begin{bmatrix} 0.8 \\ 0.8 \\ 0.6 \\ 0.4 \end{bmatrix} = \begin{bmatrix} 0.14 \\ 0.09 \\ 0.08 \end{bmatrix} \end{aligned}$$

Hence the best-fitting function is $g = 0.14 + 0.09x_1 + 0.08x_2$. The amount of computation would have been reduced if the normal equations had been constructed and then solved by gaussian elimination.

Least Squares Approximation

In many scientific investigations, data are collected that relate two variables. For example, if x is the number of dollars spent on advertising by a manufacturer and y is the value of sales in the region in question, the manufacturer could generate data by spending x_1, x_2, \dots, x_n dollars at different times and measuring the corresponding sales values y_1, y_2, \dots, y_n .



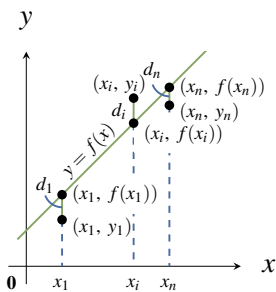
Suppose it is known that a linear relationship exists between the variables x and y —in other words, that $y = a + bx$ for some constants a and b . If the data are plotted, the points $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ may appear to lie on a straight line and estimating a and b requires finding the “best-fitting” line through these data points. For example, if five data points occur as shown in the diagram, line 1 is clearly a better fit than line 2. In general, the problem is to find the values of the constants a and b such that the line $y = a + bx$ best approximates the data in question. Note that an exact fit would be obtained if a and b were such that $y_i = a + bx_i$ were true for each data point (x_i, y_i) . But this is too much to expect. Ex-

perimental errors in measurement are bound to occur, so the choice of a and b should be made in such a way that the errors between the observed values y_i and the corresponding fitted values $a + bx_i$ are in some sense minimized. Least squares approximation is a way to do this.

The first thing we must do is explain exactly what we mean by the *best fit* of a line $y = a + bx$ to an observed set of data points $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$. For convenience, write the linear function $r_0 + r_1x$ as

$$f(x) = r_0 + r_1x$$

so that the fitted points (on the line) have coordinates $(x_1, f(x_1)), \dots, (x_n, f(x_n))$.



The second diagram is a sketch of what the line $y = f(x)$ might look like. For each i the observed data point (x_i, y_i) and the fitted point $(x_i, f(x_i))$ need not be the same, and the distance d_i between them measures how far the line misses the observed point. For this reason d_i is often called the **error** at x_i , and a natural measure of how close the line $y = f(x)$ is to the observed data points is the sum $d_1 + d_2 + \dots + d_n$ of all these errors. However, it turns out to be better to use the sum of squares

$$S = d_1^2 + d_2^2 + \dots + d_n^2$$

as the measure of error, and the line $y = f(x)$ is to be chosen so as to make this sum as small as possible. This line is said to be the **least squares approximating line** for the data points $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$.

The square of the error d_i is given by $d_i^2 = [y_i - f(x_i)]^2$ for each i , so the quantity S to be minimized is the sum:

$$S = [y_1 - f(x_1)]^2 + [y_2 - f(x_2)]^2 + \dots + [y_n - f(x_n)]^2$$

Note that all the numbers x_i and y_i are *given* here; what is required is that the *function* f be chosen in such a way as to minimize S . Because $f(x) = r_0 + r_1x$, this amounts to choosing r_0 and r_1 to minimize S . This problem can be solved using Theorem 5.6.1. The following notation is convenient.

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \text{and} \quad f(\mathbf{x}) = \begin{bmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_n) \end{bmatrix} = \begin{bmatrix} r_0 + r_1x_1 \\ r_0 + r_1x_2 \\ \vdots \\ r_0 + r_1x_n \end{bmatrix}$$

Then the problem takes the following form: Choose r_0 and r_1 such that

$$S = [y_1 - f(x_1)]^2 + [y_2 - f(x_2)]^2 + \dots + [y_n - f(x_n)]^2 = \|\mathbf{y} - f(\mathbf{x})\|^2$$

is as small as possible. Now write

$$M = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \quad \text{and} \quad \mathbf{r} = \begin{bmatrix} r_0 \\ r_1 \end{bmatrix}$$

Then $M\mathbf{r} = f(\mathbf{x})$, so we are looking for a column $\mathbf{r} = \begin{bmatrix} r_0 \\ r_1 \end{bmatrix}$ such that $\|\mathbf{y} - M\mathbf{r}\|^2$ is as small as possible. In other words, we are looking for a best approximation \mathbf{z} to the system $M\mathbf{r} = \mathbf{y}$. Hence Theorem 5.6.1 applies directly, and we have

Theorem 5.6.2

Suppose that n data points $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ are given, where at least two of x_1, x_2, \dots, x_n are distinct. Put

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad M = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$$

Then the least squares approximating line for these data points has equation

$$y = z_0 + z_1x$$

where $\mathbf{z} = \begin{bmatrix} z_0 \\ z_1 \end{bmatrix}$ is found by gaussian elimination from the normal equations

$$(M^T M)\mathbf{z} = M^T \mathbf{y}$$

The condition that at least two of x_1, x_2, \dots, x_n are distinct ensures that $M^T M$ is an invertible matrix, so \mathbf{z} is unique:

$$\mathbf{z} = (M^T M)^{-1} M^T \mathbf{y}$$

Example 5.6.3

Let data points $(x_1, y_1), (x_2, y_2), \dots, (x_5, y_5)$ be given as in the accompanying table. Find the least squares approximating line for these data.

x	y
1	1
3	2
4	3
6	4
7	5

Solution. In this case we have

$$\begin{aligned} M^T M &= \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_5 \end{bmatrix} \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_5 \end{bmatrix} \\ &= \begin{bmatrix} 5 & x_1 + \cdots + x_5 \\ x_1 + \cdots + x_5 & x_1^2 + \cdots + x_5^2 \end{bmatrix} = \begin{bmatrix} 5 & 21 \\ 21 & 111 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} \text{and } M^T \mathbf{y} &= \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_5 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_5 \end{bmatrix} \\ &= \begin{bmatrix} y_1 + y_2 + \cdots + y_5 \\ x_1 y_1 + x_2 y_2 + \cdots + x_5 y_5 \end{bmatrix} = \begin{bmatrix} 15 \\ 78 \end{bmatrix} \end{aligned}$$

so the normal equations $(M^T M)\mathbf{z} = M^T \mathbf{y}$ for $\mathbf{z} = \begin{bmatrix} z_0 \\ z_1 \end{bmatrix}$ become

$$\begin{bmatrix} 5 & 21 \\ 21 & 111 \end{bmatrix} \begin{bmatrix} z_0 \\ z_1 \end{bmatrix} = \begin{bmatrix} 15 \\ 78 \end{bmatrix}$$

The solution (using gaussian elimination) is $\mathbf{z} = \begin{bmatrix} z_0 \\ z_1 \end{bmatrix} = \begin{bmatrix} 0.24 \\ 0.66 \end{bmatrix}$ to two decimal places, so the least squares approximating line for these data is $y = 0.24 + 0.66x$. Note that $M^T M$ is indeed invertible here (the determinant is 114), and the exact solution is

$$\mathbf{z} = (M^T M)^{-1} M^T \mathbf{y} = \frac{1}{114} \begin{bmatrix} 111 & -21 \\ -21 & 5 \end{bmatrix} \begin{bmatrix} 15 \\ 78 \end{bmatrix} = \frac{1}{114} \begin{bmatrix} 27 \\ 75 \end{bmatrix} = \frac{1}{38} \begin{bmatrix} 9 \\ 25 \end{bmatrix}$$

Least Squares Approximating Polynomials

Suppose now that, rather than a straight line, we want to find a polynomial

$$y = f(x) = r_0 + r_1 x + r_2 x^2 + \cdots + r_m x^m$$

of degree m that best approximates the data pairs $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$. As before, write

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \text{and} \quad f(\mathbf{x}) = \begin{bmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_n) \end{bmatrix}$$

For each x_i we have two values of the variable y , the observed value y_i , and the computed value $f(x_i)$. The problem is to choose $f(x)$ —that is, choose r_0, r_1, \dots, r_m —such that the $f(x_i)$ are as close as possible to the y_i . Again we define “as close as possible” by the least squares condition: We choose the r_i such that

$$\|\mathbf{y} - f(\mathbf{x})\|^2 = [y_1 - f(x_1)]^2 + [y_2 - f(x_2)]^2 + \cdots + [y_n - f(x_n)]^2$$

is as small as possible.

Definition 5.15 Least Squares Approximation

A polynomial $f(x)$ satisfying this condition is called a **least squares approximating polynomial** of degree m for the given data pairs.

If we write

$$M = \begin{bmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^m \\ 1 & x_2 & x_2^2 & \cdots & x_2^m \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n & x_n^2 & \cdots & x_n^m \end{bmatrix} \quad \text{and} \quad \mathbf{r} = \begin{bmatrix} r_0 \\ r_1 \\ \vdots \\ r_m \end{bmatrix}$$

we see that $f(\mathbf{x}) = M\mathbf{r}$. Hence we want to find \mathbf{r} such that $\|\mathbf{y} - M\mathbf{r}\|^2$ is as small as possible; that is, we want a best approximation \mathbf{z} to the system $M\mathbf{r} = \mathbf{y}$. Theorem 5.6.1 gives the first part of Theorem 5.6.3.

Theorem 5.6.3

Let n data pairs $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ be given, and write

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad M = \begin{bmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^m \\ 1 & x_2 & x_2^2 & \cdots & x_2^m \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n & x_n^2 & \cdots & x_n^m \end{bmatrix} \quad \mathbf{z} = \begin{bmatrix} z_0 \\ z_1 \\ \vdots \\ z_m \end{bmatrix}$$

1. If \mathbf{z} is any solution to the normal equations

$$(M^T M)\mathbf{z} = M^T \mathbf{y}$$

then the polynomial

$$z_0 + z_1x + z_2x^2 + \cdots + z_mx^m$$

is a least squares approximating polynomial of degree m for the given data pairs.

2. If at least $m + 1$ of the numbers x_1, x_2, \dots, x_n are distinct (so $n \geq m + 1$), the matrix $M^T M$ is invertible and \mathbf{z} is uniquely determined by

$$\mathbf{z} = (M^T M)^{-1} M^T \mathbf{y}$$

Proof. It remains to prove (2), and for that we show that the columns of M are linearly independent (Theorem 5.4.3). Suppose a linear combination of the columns vanishes:

$$r_0 \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} + r_1 \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + \cdots + r_m \begin{bmatrix} x_1^m \\ x_2^m \\ \vdots \\ x_n^m \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

If we write $q(x) = r_0 + r_1x + \cdots + r_mx^m$, equating coefficients shows that

$$q(x_1) = q(x_2) = \cdots = q(x_n) = 0$$

Hence $q(x)$ is a polynomial of degree m with at least $m + 1$ distinct roots, so $q(x)$ must be the zero polynomial (see Appendix D or Theorem 6.5.4). Thus $r_0 = r_1 = \cdots = r_m = 0$ as required. \square

Example 5.6.4

Find the least squares approximating quadratic $y = z_0 + z_1x + z_2x^2$ for the following data points.

$$(-3, 3), (-1, 1), (0, 1), (1, 2), (3, 4)$$

Solution. This is an instance of Theorem 5.6.3 with $m = 2$. Here

$$\mathbf{y} = \begin{bmatrix} 3 \\ 1 \\ 1 \\ 2 \\ 4 \end{bmatrix} \quad M = \begin{bmatrix} 1 & -3 & 9 \\ 1 & -1 & 1 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 3 & 9 \end{bmatrix}$$

Hence,

$$M^T M = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ -3 & -1 & 0 & 1 & 3 \\ 9 & 1 & 0 & 1 & 9 \end{bmatrix} \begin{bmatrix} 1 & -3 & 9 \\ 1 & -1 & 1 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 3 & 9 \end{bmatrix} = \begin{bmatrix} 5 & 0 & 20 \\ 0 & 20 & 0 \\ 20 & 0 & 164 \end{bmatrix}$$

$$M^T \mathbf{y} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ -3 & -1 & 0 & 1 & 3 \\ 9 & 1 & 0 & 1 & 9 \end{bmatrix} \begin{bmatrix} 3 \\ 1 \\ 1 \\ 2 \\ 4 \end{bmatrix} = \begin{bmatrix} 11 \\ 4 \\ 66 \end{bmatrix}$$

The normal equations for \mathbf{z} are

$$\begin{bmatrix} 5 & 0 & 20 \\ 0 & 20 & 0 \\ 20 & 0 & 164 \end{bmatrix} \mathbf{z} = \begin{bmatrix} 11 \\ 4 \\ 66 \end{bmatrix} \quad \text{whence } \mathbf{z} = \begin{bmatrix} 1.15 \\ 0.20 \\ 0.26 \end{bmatrix}$$

This means that the least squares approximating quadratic for these data is $y = 1.15 + 0.20x + 0.26x^2$.

Other Functions

There is an extension of Theorem 5.6.3 that should be mentioned. Given data pairs $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, that theorem shows how to find a polynomial

$$f(x) = r_0 + r_1x + \cdots + r_mx^m$$

such that $\|\mathbf{y} - f(\mathbf{x})\|^2$ is as small as possible, where \mathbf{x} and $f(\mathbf{x})$ are as before. Choosing the appropriate polynomial $f(x)$ amounts to choosing the coefficients r_0, r_1, \dots, r_m , and Theorem 5.6.3 gives a formula for the optimal choices. Here $f(x)$ is a linear combination of the functions $1, x, x^2, \dots, x^m$ where the r_i are the coefficients, and this suggests applying the method to other functions. If $f_0(x), f_1(x), \dots, f_m(x)$ are given functions, write

$$f(x) = r_0f_0(x) + r_1f_1(x) + \cdots + r_mf_m(x)$$

where the r_i are real numbers. Then the more general question is whether r_0, r_1, \dots, r_m can be found such that $\|\mathbf{y} - f(\mathbf{x})\|^2$ is as small as possible where

$$f(\mathbf{x}) = \begin{bmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_m) \end{bmatrix}$$

Such a function $f(\mathbf{x})$ is called a **least squares best approximation** for these data pairs of the form $r_0f_0(x) + r_1f_1(x) + \cdots + r_mf_m(x)$, r_i in \mathbb{R} . The proof of Theorem 5.6.3 goes through to prove

Theorem 5.6.4

Let n data pairs $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ be given, and suppose that $m + 1$ functions $f_0(x), f_1(x), \dots, f_m(x)$ are specified. Write

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad M = \begin{bmatrix} f_0(x_1) & f_1(x_1) & \cdots & f_m(x_1) \\ f_0(x_2) & f_1(x_2) & \cdots & f_m(x_2) \\ \vdots & \vdots & \cdots & \vdots \\ f_0(x_n) & f_1(x_n) & \cdots & f_m(x_n) \end{bmatrix} \quad \mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{bmatrix}$$

1. If \mathbf{z} is any solution to the normal equations

$$(M^T M)\mathbf{z} = M^T \mathbf{y}$$

then the function

$$z_0f_0(x) + z_1f_1(x) + \cdots + z_mf_m(x)$$

is the best approximation for these data among all functions of the form $r_0f_0(x) + r_1f_1(x) + \cdots + r_mf_m(x)$ where the r_i are in \mathbb{R} .

2. If $M^T M$ is invertible (that is, if $\text{rank}(M) = m + 1$), then \mathbf{z} is uniquely determined; in fact, $\mathbf{z} = (M^T M)^{-1}(M^T \mathbf{y})$.

Clearly Theorem 5.6.4 contains Theorem 5.6.3 as a special case, but there is no simple test in general for whether $M^T M$ is invertible. Conditions for this to hold depend on the choice of the functions $f_0(x), f_1(x), \dots, f_m(x)$.

Example 5.6.5

Given the data pairs $(-1, 0)$, $(0, 1)$, and $(1, 4)$, find the least squares approximating function of the form $r_0x + r_12^x$.

Solution. The functions are $f_0(x) = x$ and $f_1(x) = 2^x$, so the matrix M is

$$M = \begin{bmatrix} f_0(x_1) & f_1(x_1) \\ f_0(x_2) & f_1(x_2) \\ f_0(x_3) & f_1(x_3) \end{bmatrix} = \begin{bmatrix} -1 & 2^{-1} \\ 0 & 2^0 \\ 1 & 2^1 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} -2 & 1 \\ 0 & 2 \\ 2 & 4 \end{bmatrix}$$

In this case $M^T M = \frac{1}{4} \begin{bmatrix} 8 & 6 \\ 6 & 21 \end{bmatrix}$ is invertible, so the normal equations

$$\frac{1}{4} \begin{bmatrix} 8 & 6 \\ 6 & 21 \end{bmatrix} \mathbf{z} = \begin{bmatrix} 4 \\ 9 \end{bmatrix}$$

have a unique solution $\mathbf{z} = \frac{1}{11} \begin{bmatrix} 10 \\ 16 \end{bmatrix}$. Hence the best-fitting function of the form $r_0x + r_12^x$ is

$$\bar{f}(x) = \frac{10}{11}x + \frac{16}{11}2^x. \text{ Note that } \bar{f}(\mathbf{x}) = \begin{bmatrix} \bar{f}(-1) \\ \bar{f}(0) \\ \bar{f}(1) \end{bmatrix} = \begin{bmatrix} \frac{-2}{11} \\ \frac{16}{11} \\ \frac{42}{11} \end{bmatrix}, \text{ compared with } \mathbf{y} = \begin{bmatrix} 0 \\ 1 \\ 4 \end{bmatrix}$$

Exercises for 5.6

Exercise 5.6.1 Find the best approximation to a solution of each of the following systems of equations.

<p>a. $x + y - z = 5$ $2x - y + 6z = 1$ $3x + 2y - z = 6$ $-x + 4y + z = 0$</p>	<p>b. $3x + y + z = 6$ $2x + 3y - z = 1$ $2x - y + z = 0$ $3x - 3y + 3z = 8$</p>
--	---

Exercise 5.6.2 Find the least squares approximating line $y = z_0 + z_1x$ for each of the following sets of data points.

a. $(1, 1), (3, 2), (4, 3), (6, 4)$

b. $(2, 4), (4, 3), (7, 2), (8, 1)$

c. $(-1, -1), (0, 1), (1, 2), (2, 4), (3, 6)$

d. $(-2, 3), (-1, 1), (0, 0), (1, -2), (2, -4)$

Exercise 5.6.3 Find the least squares approximating quadratic $y = z_0 + z_1x + z_2x^2$ for each of the following sets of data points.

a. $(0, 1), (2, 2), (3, 3), (4, 5)$

b. $(-2, 1), (0, 0), (3, 2), (4, 3)$

Exercise 5.6.4 Find a least squares approximating function of the form $r_0x + r_1x^2 + r_22^x$ for each of the following sets of data pairs.

- $(-1, 1), (0, 3), (1, 1), (2, 0)$
- $(0, 1), (1, 1), (2, 5), (3, 10)$

y	x_1	x_2	x_3
28	50	18	10
30	40	20	16
21	35	14	10
23	40	12	12
23	30	16	14

Exercise 5.6.5 Find the least squares approximating function of the form $r_0 + r_1x^2 + r_2\sin\frac{\pi x}{2}$ for each of the following sets of data pairs.

- $(0, 3), (1, 0), (1, -1), (-1, 2)$
- $(-1, \frac{1}{2}), (0, 1), (2, 5), (3, 9)$

Exercise 5.6.6 If M is a square invertible matrix, show that $\mathbf{z} = M^{-1}\mathbf{y}$ (in the notation of Theorem 5.6.3).

Exercise 5.6.7 Newton's laws of motion imply that an object dropped from rest at a height of 100 metres will be at a height $s = 100 - \frac{1}{2}gt^2$ metres t seconds later, where g is a constant called the acceleration due to gravity. The values of s and t given in the table are observed. Write $x = t^2$, find the least squares approximating line $s = a + bx$ for these data, and use b to estimate g .

Then find the least squares approximating quadratic $s = a_0 + a_1t + a_2t^2$ and use the value of a_2 to estimate g .

t	1	2	3
s	95	80	56

Exercise 5.6.8 A naturalist measured the heights y_i (in metres) of several spruce trees with trunk diameters x_i (in centimetres). The data are as given in the table. Find the least squares approximating line for these data and use it to estimate the height of a spruce tree with a trunk of diameter 10 cm.

x_i	5	7	8	12	13	16
y_i	2	3.3	4	7.3	7.9	10.1

Exercise 5.6.9 The yield y of wheat in bushels per acre appears to be a linear function of the number of days x_1 of sunshine, the number of inches x_2 of rain, and the number of pounds x_3 of fertilizer applied per acre. Find the best fit to the data in the table by an equation of the form $y = r_0 + r_1x_1 + r_2x_2 + r_3x_3$. [Hint: If a calculator for inverting $A^T A$ is not available, the inverse is given in the answer.]

Exercise 5.6.10

- Use $m = 0$ in Theorem 5.6.3 to show that the best-fitting horizontal line $y = a_0$ through the data points $(x_1, y_1), \dots, (x_n, y_n)$ is

$$y = \frac{1}{n}(y_1 + y_2 + \dots + y_n)$$

the average of the y coordinates.

- Deduce the conclusion in (a) without using Theorem 5.6.3.

Exercise 5.6.11 Assume $n = m + 1$ in Theorem 5.6.3 (so M is square). If the x_i are distinct, use Theorem 3.2.6 to show that M is invertible. Deduce that $\mathbf{z} = M^{-1}\mathbf{y}$ and that the least squares polynomial is the interpolating polynomial (Theorem 3.2.6) and actually passes through all the data points.

Exercise 5.6.12 Let A be any $m \times n$ matrix and write $K = \{\mathbf{x} \mid A^T A \mathbf{x} = \mathbf{0}\}$. Let \mathbf{b} be an m -column. Show that, if \mathbf{z} is an n -column such that $\|\mathbf{b} - A\mathbf{z}\|$ is minimal, then all such vectors have the form $\mathbf{z} + \mathbf{x}$ for some $\mathbf{x} \in K$. [Hint: $\|\mathbf{b} - A\mathbf{y}\|$ is minimal if and only if $A^T A \mathbf{y} = A^T \mathbf{b}$.]

Exercise 5.6.13 Given the situation in Theorem 5.6.4, write

$$f(x) = r_0p_0(x) + r_1p_1(x) + \dots + r_m p_m(x)$$

Suppose that $f(x)$ has at most k roots for any choice of the coefficients r_0, r_1, \dots, r_m , not all zero.

- Show that $M^T M$ is invertible if at least $k + 1$ of the x_i are distinct.
- If at least two of the x_i are distinct, show that there is always a best approximation of the form $r_0 + r_1 e^x$.
- If at least three of the x_i are distinct, show that there is always a best approximation of the form $r_0 + r_1 x + r_2 e^x$. [Calculus is needed.]

Exercise 5.6.14 If A is an $m \times n$ matrix, it can be proved that there exists a unique $n \times m$ matrix $A^\#$ satisfying the following four conditions: $AA^\#A = A$; $A^\#AA^\# = A^\#$; $AA^\#$ and $A^\#A$ are symmetric. The matrix $A^\#$ is called the **generalized inverse** of A , or the **Moore-Penrose inverse**.

- a. If A is square and invertible, show that $A^\# = A^{-1}$.
- b. If $\text{rank } A = m$, show that $A^\# = A^T(AA^T)^{-1}$.
- c. If $\text{rank } A = n$, show that $A^\# = (A^T A)^{-1}A^T$.

5.7 An Application to Correlation and Variance

Suppose the heights h_1, h_2, \dots, h_n of n men are measured. Such a data set is called a **sample** of the heights of all the men in the population under study, and various questions are often asked about such a sample: What is the average height in the sample? How much variation is there in the sample heights, and how can it be measured? What can be inferred from the sample about the heights of all men in the population? How do these heights compare to heights of men in neighbouring countries? Does the prevalence of smoking affect the height of a man?

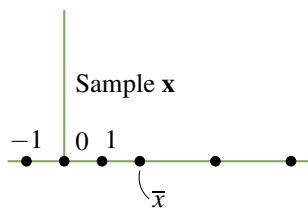
The analysis of samples, and of inferences that can be drawn from them, is a subject called *mathematical statistics*, and an extensive body of information has been developed to answer many such questions. In this section we will describe a few ways that linear algebra can be used.

It is convenient to represent a sample $\{x_1, x_2, \dots, x_n\}$ as a **sample vector**¹⁵ $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_n]$ in \mathbb{R}^n . This being done, the dot product in \mathbb{R}^n provides a convenient tool to study the sample and describe some of the statistical concepts related to it. The most widely known statistic for describing a data set is the **sample mean** \bar{x} defined by¹⁶

$$\bar{x} = \frac{1}{n}(x_1 + x_2 + \dots + x_n) = \frac{1}{n} \sum_{i=1}^n x_i$$

The mean \bar{x} is “typical” of the sample values x_i , but may not itself be one of them. The number $x_i - \bar{x}$ is called the **deviation** of x_i from the mean \bar{x} . The deviation is positive if $x_i > \bar{x}$ and it is negative if $x_i < \bar{x}$. Moreover, the sum of these deviations is zero:

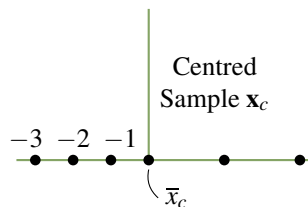
$$\sum_{i=1}^n (x_i - \bar{x}) = \left(\sum_{i=1}^n x_i \right) - n\bar{x} = n\bar{x} - n\bar{x} = 0 \tag{5.6}$$



This is described by saying that the sample mean \bar{x} is *central* to the sample values x_i .

If the mean \bar{x} is subtracted from each data value x_i , the resulting data $x_i - \bar{x}$ are said to be **centred**. The corresponding data vector is

$$\mathbf{x}_c = [x_1 - \bar{x} \ x_2 - \bar{x} \ \dots \ x_n - \bar{x}]$$



and (5.6) shows that the mean $\bar{x}_c = 0$. For example, we have plotted the sample $\mathbf{x} = [-1 \ 0 \ 1 \ 4 \ 6]$ in the first diagram. The mean is $\bar{x} = 2$,

¹⁵We write vectors in \mathbb{R}^n as row matrices, for convenience.

¹⁶The mean is often called the “average” of the sample values x_i , but statisticians use the term “mean”.