



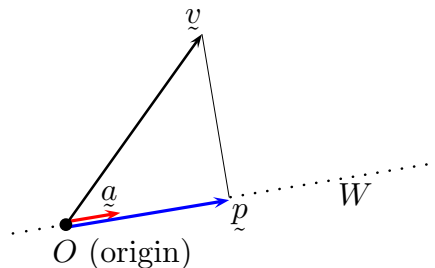
Projection onto a line through the origin

Suppose that $\underline{a} \in \mathbb{R}^n$ is a nonzero vector, and let W be the one-dimensional subspace spanned by \underline{a} . Thus W consists of all scalar multiples of \underline{a} . Geometrically, it is a straight line through the origin in n -dimensional space. If $\underline{v} \in \mathbb{R}^n$ is arbitrary then, as we saw in the first week, the projection of \underline{v} onto W is the vector $\underline{p} = A\underline{x}$, where \underline{x} is the unique solution of the linear system $(A^T A)\underline{x} = A^T \underline{v}$. Here A is any matrix whose columns form a basis for W . In the case we are currently considering, A will have only one column, since $\dim W = 1$; indeed, we may take $A = \underline{a}$ (an $n \times 1$ matrix). Now, substituting in to the formula, we find that

$$\underline{p} = \underline{a}(\underline{a}^T \underline{a})^{-1}(\underline{a}^T \underline{v}) = \frac{\underline{a} \cdot \underline{v}}{\underline{a} \cdot \underline{a}} \underline{a},$$

since the product of a row vector by a column vector can be rewritten in terms of the dot product.

Geometrically, \underline{p} is that scalar multiple of \underline{a} such that $\underline{v} - \underline{p}$ is perpendicular to W :



It is easy to check that \underline{p} depends on the direction of \underline{a} but not on its length: indeed, $\|\underline{p}\| = \|\underline{v}\| \cos \theta$, where θ is the angle between \underline{a} and \underline{p} .

For example, let $\underline{a} = (1, -1, 1)^T$, $\underline{v} = (0, 1, 2)^T$. Then $\underline{a} \cdot \underline{v} = 0 - 1 + 2 = 1$ and $\underline{a} \cdot \underline{a} = 1^2 + (-1)^2 + 1^2 = 3$. So the projection of \underline{v} onto $\text{Span}(V)$ is

$$\underline{p} = \frac{1}{3} \underline{a} = \frac{1}{3} \begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix} = \begin{pmatrix} 1/3 \\ -1/3 \\ 1/3 \end{pmatrix}.$$

“Least Squares” approximations

The theory we have been discussing has an important application in the analysis of experimental data, which we now describe.

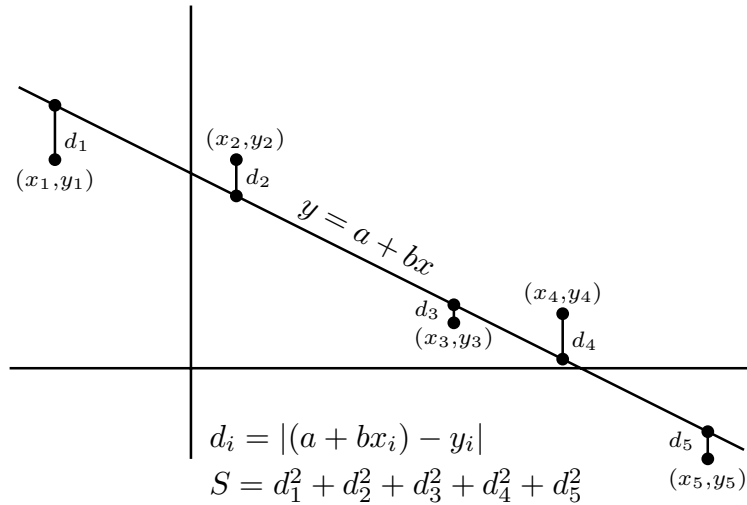
Suppose that we have set of points $\{(x_1, y_1), (x_2, y_2), \dots, (x_k, y_k)\}$ in the (x, y) -plane, and we wish to find the *straight line of best fit*. That is, we want to find a and b so that the line $y = a + bx$ goes as close as is possible to the points $(x_1, y_1), (x_2, y_2), \dots, (x_k, y_k)$.

The solution to this problem, of course, depends on how one defines closeness. For this purpose we consider the sum of the squares of the vertical distances from the data points to the line. If L is the line $y = a + bx$, then the point on L with x -coordinate x_i is $(x_i, a + bx_i)$, and this is the point on L that lies vertically above or below the data point (x_i, y_i) . The distance between these two points is $|(a + bx_i) - y_i|$, and the quantity

$$S = ((a + bx_1) - y_1)^2 + ((a + bx_2) - y_2)^2 + \dots + ((a + bx_k) - y_k)^2$$

is the sum of the squares of the vertical distances from the data points to L . Clearly $S \geq 0$, and S can only be zero if $a + bx_i = y_i$ for all i , which would mean that all the

data points lie on L . Our aim is to vary L so as to minimize S ; the line for which S is minimal is called the line of best fit.



The observation that relates this problem to vector spaces is that S can be interpreted as the square of the length of a certain vector \underline{z} in \mathbb{R}^k . Specifically $S = \|\underline{z}\|^2$, where

$$\begin{aligned} \underline{z} &= \begin{pmatrix} (a + bx_1) - y_1 \\ (a + bx_2) - y_2 \\ \vdots \\ (a + bx_k) - y_k \end{pmatrix} \\ &= a \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} + b \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \end{pmatrix} - \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_k \end{pmatrix}, \end{aligned}$$

and our task is to find a and b that minimize $\|\underline{z}\|$. (Minimizing $\|\underline{z}\|$ is of course equivalent to minimizing $\|\underline{z}\|^2$.) In other words, we must find a and b such that

$$\underline{p} = a \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} + b \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \end{pmatrix}$$

is as close as possible to $\underline{y} = (y_1, y_2, \dots, y_k)^T$. If we define W to be the subspace of \mathbb{R}^k spanned by the two column vectors $(1, 1, \dots, 1)^T$ and $(x_1, x_2, \dots, x_k)^T$ then \underline{p} is the element of W that is closest to \underline{y} ; that is, \underline{p} is the projection of \underline{y} onto W .[†] So, by the formulas we have already derived,

$$\begin{pmatrix} a \\ b \end{pmatrix} = (A^T A)^{-1} A^T \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_k \end{pmatrix}$$

[†] If the x_i are all equal then $\dim W = 1$, and the formulas below do not apply. However, in this case all the data points lie on a vertical line, and this is the line of best fit.

where A is the matrix whose two columns are $(1, 1, \dots, 1)^T$ and $(x_1, x_2, \dots, x_k)^T$. (For this problem we do not need to calculate the vector \underline{p} , since the coefficients a and b are what we want.)

As an example, let us find the line of best fit for the four data points $(0, 0)$, $(1, 1)$, $(3, 2)$, $(4, 5)$. In the notation used above,

$$A^T A = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 3 & 4 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 1 & 1 \\ 1 & 3 \\ 1 & 4 \end{pmatrix} = \begin{pmatrix} 4 & 8 \\ 8 & 26 \end{pmatrix},$$

and we have to find a and b satisfying

$$\begin{pmatrix} 4 & 8 \\ 8 & 26 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = A^T \begin{pmatrix} 0 \\ 1 \\ 2 \\ 5 \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 3 & 4 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 2 \\ 5 \end{pmatrix} = \begin{pmatrix} 8 \\ 27 \end{pmatrix}.$$

Let us use the row operations technique to solve this system:

$$\begin{aligned} \left(\begin{array}{cc|c} 4 & 8 & 8 \\ 8 & 26 & 27 \end{array} \right) &\xrightarrow{R_1 := \frac{1}{4}R_1} \left(\begin{array}{cc|c} 1 & 2 & 2 \\ 8 & 26 & 27 \end{array} \right) \xrightarrow{R_2 := R_2 - 8R_1} \left(\begin{array}{cc|c} 1 & 2 & 2 \\ 0 & 10 & 11 \end{array} \right) \\ &\xrightarrow{R_2 := 0.1R_2} \left(\begin{array}{cc|c} 1 & 2 & 2 \\ 0 & 1 & 1.1 \end{array} \right) \xrightarrow{R_1 := R_1 - 2R_2} \left(\begin{array}{cc|c} 1 & 0 & -0.2 \\ 0 & 1 & 1.1 \end{array} \right). \end{aligned}$$

So the equation of the line of best fit is $y = -0.2 + 1.1x$.

Best fitting parabolas, cubic curves, etc.

Given points (x_1, y_1) , (x_2, y_2) , \dots , (x_k, y_k) , we may want to find the best fitting polynomial equation of some specified degree. Suppose that the degree is n , so that the equation we want to find has the form $y = a_0 + a_1x + a_2x^2 + \dots + a_nx^n$. As in the straight line case, we use the sum of the squares of the vertical distances from the data points to the curve as a measure of how well the curve fits the data. The vertical distance from the point (x_i, y_i) to the curve is the distance between the points (x_i, y_i) and $(x_i, a_0 + a_1x_i + a_2x_i^2 + \dots + a_nx_i^n)$, and this is given by $|(a_0 + a_1x_i + a_2x_i^2 + \dots + a_nx_i^n) - y_i|$. So the quantity to be minimized is

$$S = \sum_{i=1}^k ((a_0 + a_1x_i + a_2x_i^2 + \dots + a_nx_i^n) - y_i)^2,$$

which is the square of the length of the vector

$$\begin{aligned} \underline{z} &= \begin{pmatrix} (a_0 + a_1x_1 + a_2x_1^2 + \dots + a_nx_1^n) - y_1 \\ (a_0 + a_1x_2 + a_2x_2^2 + \dots + a_nx_2^n) - y_2 \\ \vdots \\ (a_0 + a_1x_k + a_2x_k^2 + \dots + a_nx_k^n) - y_k \end{pmatrix} \\ &= a_0 \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} + a_1 \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \end{pmatrix} + a_2 \begin{pmatrix} x_1^2 \\ x_2^2 \\ \vdots \\ x_k^2 \end{pmatrix} + \dots + a_n \begin{pmatrix} x_1^n \\ x_2^n \\ \vdots \\ x_k^n \end{pmatrix} - \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_k \end{pmatrix}. \end{aligned}$$

We see immediately that the length of z is minimal when

$$\tilde{p} = a_0 \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} + a_1 \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \end{pmatrix} + a_2 \begin{pmatrix} x_1^2 \\ x_2^2 \\ \vdots \\ x_k^2 \end{pmatrix} + \cdots + a_n \begin{pmatrix} x_1^n \\ x_2^n \\ \vdots \\ x_k^n \end{pmatrix}$$

is the projection of $(y_1, y_2, \dots, y_k)^T$ onto the space W spanned by the vectors

$$\begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}, \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \end{pmatrix}, \begin{pmatrix} x_1^2 \\ x_2^2 \\ \vdots \\ x_k^2 \end{pmatrix}, \dots, \begin{pmatrix} x_1^n \\ x_2^n \\ \vdots \\ x_k^n \end{pmatrix}.$$

It can be shown that, provided there are at least $n + 1$ distinct x_i 's, these vectors are linearly independent, and therefore constitute a basis for W . We shall assume that this condition is satisfied.

By the same reasoning that we used in our discussion of the straight line of best fit, we conclude that

$$(A^T A) \begin{pmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{pmatrix} = A^T \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ a_k \end{pmatrix},$$

where A is the matrix whose columns constitute the above basis for W . That is,

$$A = \begin{pmatrix} 1 & x_1 & x_1^2 & \dots & x_1^n \\ 1 & x_2 & x_2^2 & \dots & x_2^n \\ 1 & x_3 & x_3^2 & \dots & x_3^n \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_k & x_k^2 & \dots & x_k^n \end{pmatrix}.$$

To illustrate this, let us find the parabola of best fit through the same four points that we used in our example of the line of best fit: $(0, 0)$, $(1, 1)$, $(3, 2)$, $(4, 5)$. This time we have

$$A^T A = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 3 & 4 \\ 0 & 1 & 9 & 16 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 3 & 9 \\ 1 & 4 & 16 \end{pmatrix} = \begin{pmatrix} 4 & 8 & 26 \\ 8 & 26 & 92 \\ 26 & 92 & 338 \end{pmatrix}$$

and so we must solve

$$\begin{pmatrix} 4 & 8 & 26 \\ 8 & 26 & 92 \\ 26 & 92 & 338 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 3 & 4 \\ 0 & 1 & 9 & 16 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \\ 2 \\ 5 \end{pmatrix} = \begin{pmatrix} 8 \\ 27 \\ 99 \end{pmatrix}.$$

To save ourselves some calculation, we can get MAGMA to do this. Of course, we first have to change it into a row vector problem by taking transposes. The equation becomes

$$(a \ b \ c) \begin{pmatrix} 4 & 8 & 26 \\ 8 & 26 & 92 \\ 26 & 92 & 338 \end{pmatrix} = (8 \ 27 \ 99).$$

Now we can use MAGMA's `Solution` function:

```
> R := RealField();
> V := VectorSpace(R,3);
> M := KMatrixSpace(R,3,3);
> v := V![8,27,99];
> ATA := M![4,8,26,8,26,92,26,92,338];
> Solution(ATA,v);
> (3/10 -7/30 1/3)
```

We conclude that the equation of the parabola of best fit is $y = \frac{3}{10} - \frac{7}{30}x + \frac{1}{3}x^2$.

Orthonormal bases

Definition. A set $\{v_1, v_2, \dots, v_k\}$ is called an *orthogonal* set if $v_i \cdot v_j = 0$ whenever $i \neq j$.

In other words, a set of vectors is said to be orthogonal if distinct elements of the set are perpendicular to one another. For example,

$$\underline{i} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \quad \underline{j} = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \quad \underline{k} = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

form an orthogonal set of vectors in \mathbb{R}^3 .

Definition. A set $\{v_1, v_2, \dots, v_k\}$ is called an *orthonormal* set if it is orthogonal and $v_i \cdot v_i = 1$ for all i .

The set $\{\underline{i}, \underline{j}, \underline{k}\}$ is also an orthonormal set in \mathbb{R}^3 . The following three vectors form an orthogonal set that is not orthonormal:

$$\begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}, \quad \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix}, \quad \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}.$$

Whenever we have an orthogonal set of vectors that are all nonzero, we can produce an orthonormal set of vectors by a process known as *normalization*, which consists of replacing each v in the set by $(1/\|v\|)v$. In other words, divide each vector by its own length: this produces a scalar multiple of the vector having length 1. In the example above, the orthonormal set obtained in this way is

$$\begin{pmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \\ 0 \end{pmatrix}, \quad \begin{pmatrix} 1/\sqrt{2} \\ -1/\sqrt{2} \\ 0 \end{pmatrix}, \quad \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}.$$

The following theorem states one of the key facts about orthogonal sets.

Theorem. *If $\{v_1, v_2, \dots, v_k\}$ is an orthogonal set of nonzero vectors in \mathbb{R}^n then it is also a linearly independent set.*

Proof. Suppose that $\lambda_1, \lambda_2, \dots, \lambda_k$ are scalars such that $\lambda_1 v_1 + \lambda_2 v_2 + \dots + \lambda_k v_k = \mathbf{0}$. Then, for any i ,

$$\begin{aligned} 0 &= \mathbf{0} \cdot v_i = (\lambda_1 v_1 + \lambda_2 v_2 + \dots + \lambda_k v_k) \cdot v_i \\ &= \lambda_1 (v_1 \cdot v_i) + \lambda_2 (v_2 \cdot v_i) + \dots + \lambda_k (v_k \cdot v_i) \\ &= \lambda_i (v_i \cdot v_i), \end{aligned}$$

all the other terms being zero since $v_j \cdot v_i = 0$ when $j \neq i$. Now since we have assumed that $v_i \neq \mathbf{0}$ it follows that $v_i \cdot v_i = \|v_i\|^2 \neq 0$, and so the above equation shows that $\lambda_i = 0$. Furthermore, i was arbitrary, and so we have shown that the only solution of $\sum_{i=1}^n \lambda_i v_i = \mathbf{0}$ is given by $\lambda_i = 0$ for all i . That is, v_1, v_2, \dots, v_k are linearly independent, as claimed. \square

By this theorem, an orthogonal set of nonzero vectors necessarily constitutes a basis for the subspace it spans. Normalization will then yield an orthonormal basis for this subspace.

Projections using orthogonal bases

Orthogonal bases are important because many formulas become much simpler to work with when expressed in terms of orthogonal bases. The formula for the projection onto a subspace is a case in point.

Let A be an $n \times k$ matrix. The columns of A form an orthonormal set of vectors if and only if $A^T A = I$, the identity matrix. This fact is quite straightforward to prove, as follows. Let v_1, v_2, \dots, v_k be the columns of A . Then $v_1^T, v_2^T, \dots, v_k^T$ are the rows of A^T , and, by the way matrix multiplication is defined, the (i, j) entry of $A^T A$ is $v_i^T v_j$, the product of the i -th row of A^T and the j -th column of A . But the product of a row vector by a column vector can alternatively be expressed as the dot product of two column vectors; so in fact we can say that the (i, j) entry of $A^T A$ is the dot product of the i -th and j -th columns of A . The columns of A comprise an orthonormal set if and only if $v_i \cdot v_j$ is zero for $i \neq j$ and 1 when $i = j$. Since $v_i \cdot v_j$ is the (i, j) -entry of $A^T A$, this condition is equivalent to the main-diagonal entries of $A^T A$ being 1 and the other entries being 0, and of course this is the same as saying that $A^T A = I$. So the columns of A form an orthonormal set if and only if $A^T A = I$, as claimed.

In the course of the above discussion, we showed that if A is a matrix with k columns v_1, v_2, \dots, v_k , then

$$A^T A = \begin{pmatrix} v_1 \cdot v_1 & v_1 \cdot v_2 & \dots & v_1 \cdot v_k \\ v_2 \cdot v_1 & v_2 \cdot v_2 & \dots & v_2 \cdot v_k \\ \vdots & \vdots & & \vdots \\ v_k \cdot v_1 & v_k \cdot v_2 & \dots & v_k \cdot v_k \end{pmatrix}. \quad (1)$$

This result is true for all matrices A , without the assumption that the columns are orthonormal. When the columns are orthonormal, the right hand side of Eq. (1) is just the identity matrix.

Definition. The matrix on the right hand side of Eq. (1) above is called the *Gram matrix* of the set of vectors $\{v_1, v_2, \dots, v_k\}$.

We turn now to consideration of projections onto a subspace for which we have an orthogonal basis. First, let us suppose that $\{a_1, a_2, \dots, a_k\}$ is an orthonormal set of

vectors in \mathbb{R}^n , and let $W = \text{Span}(a_1, a_2, \dots, a_k)$. If $v \in \mathbb{R}^n$ then the projection of v onto W is given by

$$p = A(A^T A)^{-1} A^T v = AA^T v \quad (2)$$

(since $A^T A = I$). Note that although we know that $A^T A = I$, we have no simple formula for AA^T . To evaluate the right hand side of Eq. (2) we bracket it as $A(A^T v)$, and observe that

$$A^T v = \begin{pmatrix} a_1^T \\ a_2^T \\ \vdots \\ a_k^T \end{pmatrix} v = \begin{pmatrix} a_1^T v \\ a_2^T v \\ \vdots \\ a_k^T v \end{pmatrix} = \begin{pmatrix} a_1 \cdot v \\ a_2 \cdot v \\ \vdots \\ a_k \cdot v \end{pmatrix}.$$

So we find that

$$\begin{aligned} p &= A \begin{pmatrix} a_1 \cdot v \\ a_2 \cdot v \\ \vdots \\ a_k \cdot v \end{pmatrix} = (a_1 \quad a_2 \quad \dots \quad a_k) \begin{pmatrix} a_1 \cdot v \\ a_2 \cdot v \\ \vdots \\ a_k \cdot v \end{pmatrix} \\ &= (a_1 \cdot v)a_1 + (a_2 \cdot v)a_2 + \dots + (a_k \cdot v)a_k. \end{aligned}$$

Comparing this with the formula that we obtained earlier for the projection of v onto a 1-dimensional subspace, we see that in fact p is the sum of the projections of v on to the one-dimensional subspaces spanned by the a_i , for each i from 1 to k . (The projection of v onto the one-dimensional space $\text{Span}(a_i)$ is $\frac{a_i \cdot v}{a_i \cdot a_i} a_i = (a_i \cdot v)a_i$, since $a_i \cdot a_i = 1$.)