

Chap 3. Orthogonality

3.1 Orthogonal Vectors and Subspaces

For $\mathbf{x} = (x_1, \dots, x_n)$, $\|\mathbf{x}\| = \sqrt{\mathbf{x}^T \cdot \mathbf{x}}$: the length of \mathbf{x} ,

$$\|\mathbf{x}\|^2 = x_1^2 + x_2^2 + \dots + x_n^2$$

$\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$: orthogonal $\Leftrightarrow \mathbf{x}^T \cdot \mathbf{y} = x_1 y_1 + \dots + x_n y_n = 0$
the inner product of \mathbf{x} and \mathbf{y} .

3A

The inner product $\mathbf{x}^T \cdot \mathbf{y} = 0$ iff \mathbf{x} and \mathbf{y} are orthogonal.

If $\mathbf{x}^T \cdot \mathbf{y} > 0$, their angle is less than 90° .

If $\mathbf{x}^T \cdot \mathbf{y} < 0$, their angle is greater than 90° .

Useful Fact: If nonzero vectors $\mathbf{v}_1, \dots, \mathbf{v}_k$ are mutually orthogonal, then those vectors are linearly independent.

Proof Suppose $c_1 \mathbf{v}_1 + \dots + c_k \mathbf{v}_k = \mathbf{0}$.

$$\Rightarrow \mathbf{v}_i (c_1 \mathbf{v}_1 + \dots + c_k \mathbf{v}_k) = c_i \mathbf{v}_i^T \cdot \mathbf{v}_i = 0$$

Since \mathbf{v}_i are nonzero, $\mathbf{v}_i^T \cdot \mathbf{v}_i \neq 0$ and $c_i = 0$ \square

The coordinate vectors $\mathbf{e}_1, \dots, \mathbf{e}_n \in \mathbb{R}^n$ are mutually orthogonal unit vectors. When they are rotated, the result is a new orthonormal basis: a new system of mutually orthogonal unit vectors in \mathbb{R}^n . In \mathbb{R}^2 , $\mathbf{v}_1 = (\cos \theta, \sin \theta)$, $\mathbf{v}_2 = (-\sin \theta, \cos \theta)$ form an orthonormal basis.

Orthogonal Subspaces

3B Two subspaces V and W of \mathbb{R}^n are orthogonal if every vector $v \in V$ is orthogonal to every $w \in W$:
 $v^T \cdot w = 0$ for all v and w .

3C] Fundamental Theorem of orthogonality

The row space is orthogonal to the nullspace (in \mathbb{R}^n).

The column space is orthogonal to the left nullspace (in \mathbb{R}^m).

$$\text{For } \mathbf{x} \in N(A) \Leftrightarrow A\mathbf{x} = \mathbf{0} \Rightarrow \mathbf{v}^T \mathbf{x} = \mathbf{z}^T A \mathbf{x} = 0$$

$$\forall \mathbf{v} \in C(A^T) \Leftrightarrow \mathbf{v} = A^T \mathbf{z} \text{ for some } \mathbf{z}$$

$$\therefore N(A) \perp C(A^T)$$

Definition (orthogonal Complement)

V : a subspace of \mathbb{R}^n .

V^\perp : the orthogonal complement of V is the space of
 all vectors orthogonal to V .
 $\hookrightarrow "V_{\text{perp}}"$.

3D] Fundamental Theorem of Linear Algebra, Part II

The nullspace is the orthogonal complement of the row space

The left nullspace is the orthogonal complement
of the column space in \mathbb{R}^m . in \mathbb{R}^n

3E] $A\mathbf{x} = \mathbf{b}$ is solvable iff $\mathbf{y}^T \mathbf{b} = 0$ whenever $\mathbf{y}^T A = 0$

$$\text{For } \mathbf{b} \in C(A) \Leftrightarrow \mathbf{b} \in (N(A^T))^{\perp}$$

Ex:

$$A = \begin{bmatrix} 1 & -1 & 0 \\ 0 & 1 & -1 \\ -1 & 0 & 1 \end{bmatrix} \Rightarrow \mathbf{y} = (1, 1, 1)^T \in N(A^T)$$

since $\mathbf{y}^T A = 0$.

Thus, $A\mathbf{x} = \mathbf{b}$ is solvable $\Leftrightarrow b_1 + b_2 + b_3 = 0$.
 $\hookrightarrow \mathbf{y}^T \mathbf{b}$.

The Matrix and the Subspaces

If $V = V^\perp$, then $V = V^\perp$ and $\dim V + \dim V^\perp = n$.

In other words $V^\perp = V$. The whole space \mathbb{R}^n can be decomposed into two perpendicular parts.

3F From the row space to the column space, A is actually invertible. Every vector b in the column space comes from exactly one vector x_r in the row space.

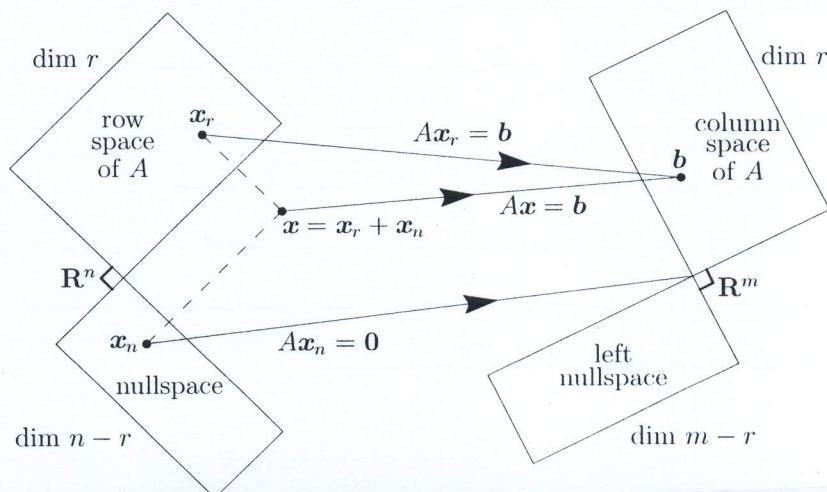


Figure 3.4 The true action $Ax = A(x_{\text{row}} + x_{\text{null}})$ of any m by n matrix.

For $\mathbf{b} \in C(A)$, $\mathbf{b} = A\mathbf{x} = A\mathbf{x}_r$, where $\mathbf{x} = \mathbf{x}_r + \mathbf{x}_n$ and $A\mathbf{x}_n = 0$.

Let \mathbf{x}' in the row space and $A\mathbf{x}' = \mathbf{b}$,

then $A(\mathbf{x}_r - \mathbf{x}') = \mathbf{b} - \mathbf{b} = \mathbf{0}$, and $\mathbf{x}_r - \mathbf{x}' \in N(A)$

and $\mathbf{x}_r - \mathbf{x}'$ is also in the row space $\Rightarrow \mathbf{x}_r = \mathbf{x}'$.

Exactly one vector in the row space is carried to \mathbf{b} . \square

o. A^T goes from \mathbb{R}^m to \mathbb{R}^n and from $C(A)$ to $C(A^T)$.

But $A^T \neq A^{-1}$ in general. A^T moves the spaces correctly, but not the individual vectors. When A^T fails to exist, the best substitute is the pseudo inverse A^+ :

$A^T A \mathbf{x} = \mathbf{x}$ for $\mathbf{x} \in C(A^T)$ and $A^+ \mathbf{y} = \mathbf{0}$ for $\mathbf{y} \in N(A^T)$.

3.2 Cosines and Projections onto Lines

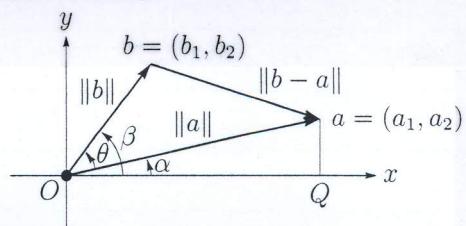
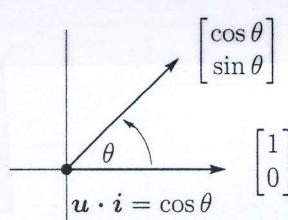


Figure 3.6 The cosine of the angle $\theta = \beta - \alpha$ using inner products.

Cosine formula: $\cos \theta = \cos \beta \cos \alpha + \sin \beta \sin \alpha = \frac{a_1 b_1 + a_2 b_2}{\|a\| \|b\|}$

$$\frac{b_1}{\|b\|} \quad \frac{a_1}{\|a\|} \quad \frac{b_2}{\|b\|} \quad \frac{a_2}{\|a\|}$$

[3G] $\cos \theta = \frac{a^T b}{\|a\| \|b\|}$

Law of Cosines: $\|b - a\|^2 = \|b\|^2 + \|a\|^2 - 2\|b\|\|a\|\cos \theta$.

$$\Rightarrow \|b^T b - 2a^T b + a^T a\| = \|b^T b + a^T a - 2\|b\|\|a\|\cos \theta$$

$$\therefore a^T b = \|a\| \|b\| \cos \theta$$

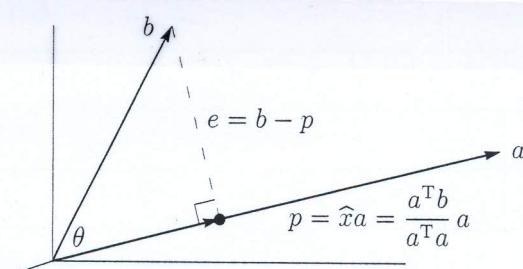


Figure 3.7 The projection p of b onto a , with $\cos \theta = \frac{a^T b}{\|a\| \|b\|}$.

[3H] The projection of b onto the line in the direction of a is

$$p = \hat{x}a = \frac{a^T b}{a^T a} a \quad \begin{bmatrix} \text{eg } (b - \hat{x}a) \perp a \\ \Leftrightarrow a^T(b - \hat{x}a) = 0 \end{bmatrix}$$

Ex1: All vectors a and b satisfy the Schwartz Inequality:

$$|\alpha^T b| \leq \|\alpha\| \|\beta\|,$$

which is $|\cos \theta| \leq 1$ in \mathbb{R}^n .

$$\begin{aligned} \|\beta - \frac{\alpha^T \beta}{\alpha^T \alpha} \alpha\|^2 &= \|\beta\|^2 - 2 \frac{(\alpha^T \beta)^2}{\alpha^T \alpha} + \frac{(\alpha^T \beta)^2}{\alpha^T \alpha} \alpha^T \alpha \\ &= \frac{(\beta^T \beta)(\alpha^T \alpha) - (\alpha^T \beta)^2}{\alpha^T \alpha} \geq 0 \end{aligned}$$

The equality holds iff b is a multiple of a .

Ex1: Project $b = (1, 2, 3)$ onto the line through $a = (1, 1, 1)$.

$$\hat{x} = \frac{\alpha^T b}{\alpha^T \alpha} = \frac{6}{3} = 2 \text{ and } p = \hat{x}a = (2, 2, 2).$$

$$\cos \theta = \frac{\|p\|}{\|b\|} = \frac{\sqrt{12}}{\sqrt{14}} \text{ and } \cos \theta = \frac{\alpha^T b}{\|\alpha\| \|b\|} = \frac{6}{\sqrt{3} \sqrt{14}}.$$
$$\Rightarrow 6 \leq \sqrt{3} \sqrt{14} \text{ and } \sqrt{36} \leq \sqrt{42}.$$

The cosine is less than 1 since b is not parallel to a .

Projection Matrix of Rank 1

$$P = a \cdot \frac{a^T b}{a^T a} \Rightarrow P = \frac{a \cdot a^T}{a^T a} : \text{projection matrix}$$

① P is a symmetric matrix

② Its square is itself: $P^2 = P$.

Ex2: $a = (1, 1, 1)$

$$P = \frac{a \cdot a^T}{a^T a} = \frac{1}{3} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{bmatrix}$$

Remark on scaling:

The projection matrix is the same if α is scaled.

$$\alpha = \begin{bmatrix} \alpha \\ \alpha \\ \alpha \end{bmatrix} \text{ gives } P = \frac{1}{3\alpha^2} \begin{bmatrix} \alpha \\ \alpha \\ \alpha \end{bmatrix} \begin{bmatrix} \alpha & \alpha & \alpha \end{bmatrix} = \begin{bmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{bmatrix}$$

If α has a unit length, $P = \alpha \cdot \alpha^T$ ($\alpha^T \cdot \alpha = 1$).

Ex3: Project onto the θ -direction in the xy-plane.

The line goes through $\alpha = (\cos\theta, \sin\theta)$.

$$P = \frac{\alpha \cdot \alpha^T}{\alpha^T \alpha} = \frac{1}{\underbrace{\cos^2\theta + \sin^2\theta}_1} \begin{bmatrix} \cos\theta \\ \sin\theta \end{bmatrix} \begin{bmatrix} \cos\theta & \sin\theta \end{bmatrix} = \begin{bmatrix} \cos^2\theta & \cos\theta \sin\theta \\ \cos\theta \sin\theta & \sin^2\theta \end{bmatrix}$$

- o This matrix P was discovered in Section 2.6.
Now we know P in any dimensions.
- o To project \mathbf{b} onto α , multiply by the projection matrix
 $P = P\mathbf{b}$.

3.3 Projections and Least Squares

Least squares solution for $\alpha x = lb$ by minimizing

$$E^2 = \| \alpha x - lb \| ^2 = (a_1 x - b_1)^2 + \dots + (a_m x - b_m)^2.$$

$$\frac{1}{2} \frac{dE^2}{dx} = (a_1 x - b_1) a_1 + \dots + (a_m x - b_m) a_m = 0$$

3K The least squares solution to a problem $\alpha x = lb$ in one unknown is $\hat{x} = \alpha^T lb / \alpha^T \alpha$.

Orthogonality : $\alpha^T (lb - \hat{x} \alpha) = \alpha^T lb - \frac{\alpha^T lb}{\alpha^T \alpha} \cdot \alpha^T \alpha = 0$

of α and ϵ

ϵ : the error vector

Least-Squares Problem with Several Variables

$Ax = lb$, where A is an $m \times n$ matrix.

$\Rightarrow E = \| Ax - lb \|$: error, which is the distance from lb to the point Ax in the column space.

- o. Locate the point $p = A\hat{x}$ that is closer to lb than any other point in the column space.
- o. The error $e = lb - A\hat{x}$ must be perpendicular to the column space.

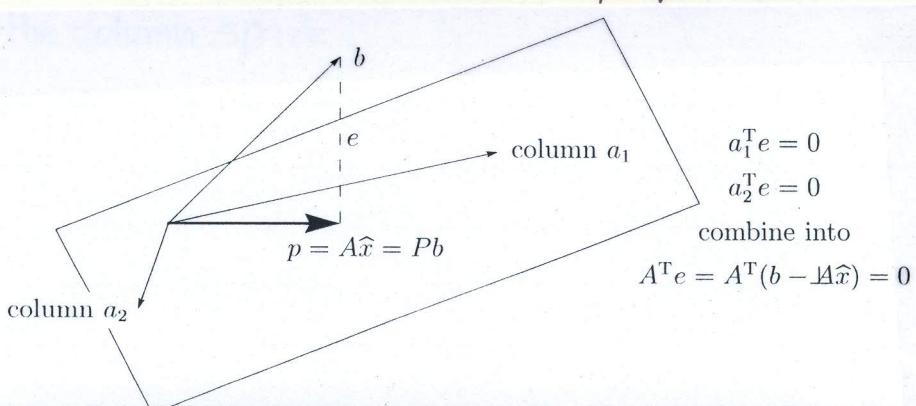


Figure 3.8 Projection onto the column space of a 3 by 2 matrix.

① All vectors perpendicular to the column space lie in the left nullspace. Thus the error $e = \mathbf{b} - A\hat{x}$ must be in the nullspace of A^T :

$$A^T(\mathbf{b} - A\hat{x}) = \mathbf{0} \text{ or } A^T A \hat{x} = A^T \mathbf{b}.$$

② The error vector must be perpendicular to each column a_1, \dots, a_n of A :

$$\begin{cases} a_1^T(\mathbf{b} - A\hat{x}) = 0 \\ \vdots \\ a_n^T(\mathbf{b} - A\hat{x}) = 0 \end{cases} \quad \text{or} \quad \begin{bmatrix} a_1^T \\ \vdots \\ a_n^T \end{bmatrix} \begin{bmatrix} \mathbf{b} - A\hat{x} \end{bmatrix} = \mathbf{0}.$$

This is again $A^T(\mathbf{b} - A\hat{x}) = \mathbf{0}$ and $A^T A \hat{x} = A^T \mathbf{b}$.

o. Taking partial derivatives of $E^2 = (A\hat{x} - \mathbf{b})^T(A\hat{x} - \mathbf{b})$ gives the same $\underline{\partial E^2 / \partial \hat{x}} = \underline{\partial A^T A \hat{x} - \partial A^T \mathbf{b} / \partial \hat{x}} = \mathbf{0}$.

↳ symmetric square matrix

BL When $A\hat{x} = \mathbf{b}$ is inconsistent, its least-squares solution minimizes $\|A\hat{x} - \mathbf{b}\|^2$:

Normal equations: $A^T A \hat{x} = A^T \mathbf{b}$.

$A^T A$ is invertible exactly when the columns of A are linearly independent! Then,

Best estimate \hat{x} : $\hat{x} = (A^T A)^{-1} A^T \mathbf{b}$.

The projection of \mathbf{b} onto the column space is the nearest point $A\hat{x}$:

Projection: $\hat{\mathbf{p}} = A\hat{x} = A(A^T A)^{-1} A^T \mathbf{b}$.

Ex: $A = \begin{bmatrix} 1 & 2 \\ 1 & 3 \\ 0 & 0 \end{bmatrix}, \mathbf{b} = \begin{bmatrix} 4 \\ 5 \\ 6 \end{bmatrix}$ $A\hat{x} = \mathbf{b}$ has no solution
 $A^T A \hat{x} = A^T \mathbf{b}$ gives the best \hat{x} .

The projection of $\mathbf{b} = (4, 5, 6)$ is $\hat{\mathbf{p}} = (4, 5, 0)$.

Solving the normal equations:

$$ATA = \begin{bmatrix} 1 & 1 & 0 \\ 2 & 3 & 0 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 1 & 3 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 2 & 5 \\ 5 & 13 \end{bmatrix}$$

$$\hat{x} = (ATA)^{-1} A^T b = \begin{bmatrix} 13 & -5 \\ -5 & 2 \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 \\ 2 & 3 & 0 \end{bmatrix} \begin{bmatrix} 4 \\ 5 \\ 6 \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}.$$

Projection $P = A\hat{x} = \begin{bmatrix} 1 & 2 \\ 1 & 3 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 4 \\ 5 \\ 0 \end{bmatrix}.$

Remarks: $\rightarrow lb = Ax$ (a combination of the columns)

① $lb \in C(A) \Rightarrow P = A(ATA)^{-1} A^T Ax = A\hat{x} = lb.$

\therefore The projection of lb is still lb .

② $lb \in N(A^T) \Rightarrow P = A(ATA)^{-1} A^T lb = A(ATA)^{-1} "0" = 0$
left nullspace \rightarrow

$\therefore lb$ projects to the zero vector

③ When A is square and invertible, $C(A) = \mathbb{R}^n$. \rightarrow the whole space

$$\begin{aligned} A \text{ invertible} \Rightarrow P &= A(ATA)^{-1} A^T lb \\ &= A \cdot A^{-1}(A^T)^{-1} A^T lb = lb. \end{aligned}$$

\therefore Every vector projects to itself, $P = lb$, $\hat{x} = x$.

⊕ A has only one column $\Rightarrow A^T A = \alpha^T \alpha$ and $\hat{x} = \alpha^T b / \alpha^T \alpha$.

The Cross-Product Matrix $A^T A$

o. $A^T A$ has the same nullspace as A .

From ① $Ax = 0 \Rightarrow A^T A x = 0$

② $A^T A x = 0 \Rightarrow x^T A^T A x = 0$, $\|Ax\|^2 = 0$, $Ax = 0$.

\rightarrow the same is true for $A^T A$

3M If A has independent columns, then

$A^T A$ is square, symmetric, and invertible.

Projection Matrices (We assume the independence of the columns of A in what follows)

$$P = A(ATA)^{-1}AT : \text{projection matrix}$$

- o. Pb is the component of b in the column space, and the error $e = b - Pb = (I-P)b$ is in $N(AT)$, the orthogonal complement of $C(A)$.
 $(I-P)$ is also a projection matrix!
 $\therefore Pb$ and $(I-P)b$ are two perpendicular components of b .

[3N] The projection matrix $P = A(ATA)^{-1}AT$ has two basic properties: (i) $P^2 = P$ and (ii) $P^T = P$.

Conversely, any symmetric matrix with $P^2 = P$ represents a projection.

<pf>

$$\begin{aligned} o. P^2 &= A(ATA)^{-1}AT A(ATA)^{-1}AT = A(ATA)^{-1}AT = P \\ P^T &= (AT)^T (ATA)^{-1}T AT = A((ATA)^{-1})^T AT \\ &= A(ATA)^{-1}AT = P. \end{aligned}$$

o. For the converse, from $P^2 = P$ and $P^T = P$,

We can show the error vector $b - Pb$ is orthogonal to the column space. For any vector $Pc \in C(P)$,

$$(b - Pb)^T Pc = b^T (I - P)^T P c = b^T (P - P^2) c = 0.$$

Thus, $b - Pb$ is orthogonal to the space, and Pb is the projection onto the column space. \square

Ex1: A : an invertible 4×4 matrix. $\Rightarrow C(A) = \mathbb{R}^4$.

$$P = A(ATA)^{-1}AT = A A^{-1}(A^T)^{-1}AT = I.$$

$$I: \text{symmetric}, \quad I^2 = I, \quad b - Ib = 0.$$

The projection onto the whole space is the Identity matrix.

Least-Squares Fitting of Data

In a series of experiments, we expect the output b to be a linear function of the input t . We look for a straight line: $b = C + Dt$, with two unknowns C, D .

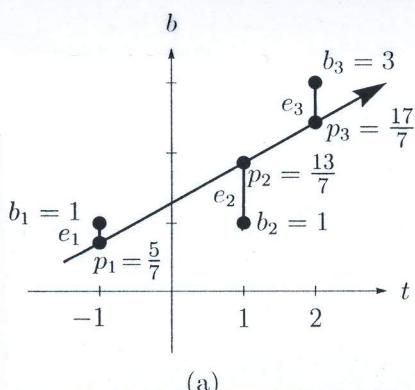
$$\begin{cases} C + Dt_1 = b_1 \\ C + Dt_2 = b_2 \\ \vdots \\ C + Dt_m = b_m \end{cases} \text{ or } \begin{bmatrix} 1 & t_1 \\ 1 & t_2 \\ \vdots & \vdots \\ 1 & t_m \end{bmatrix} \begin{bmatrix} C \\ D \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix}$$

\Downarrow

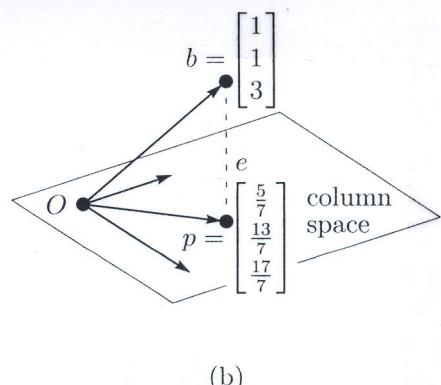
$$A \times = \mathbf{b}$$

The best solution (\hat{C}, \hat{D}) is the \hat{x} that minimizes $E^2 = \| \mathbf{b} - A \hat{x} \|^2 = (b_1 - C - Dt_1)^2 + \dots + (b_m - C - Dt_m)^2$.

The vector $p = A \hat{x}$ is as close as possible to \mathbf{b} . Of all straight lines $b = C + Dt$, we choose the one that best fits the data. On the graph, the errors are the vertical distances $b - C - Dt$ to the straight line (not perpendicular distances!).



(a)



(b)

Figure 3.9 Straight-line approximation matches the projection p of b .

Ex 2: Three measurements:

$$b_1=1 \text{ at } t_1=-1; b_2=1 \text{ at } t_2=1; b_3=3 \text{ at } t_3=2$$

$$Ax = Ib \text{ is } \begin{cases} C - D = 1 \\ C + D = 1 \\ C + 2D = 3 \end{cases} \text{ or } \begin{bmatrix} 1 & -1 \\ 1 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} C \\ D \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 3 \end{bmatrix}$$

$$A^T A x = A^T b \text{ is } \begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix} \begin{bmatrix} C \\ D \end{bmatrix} = \begin{bmatrix} 5 \\ 6 \end{bmatrix}$$

$$\therefore \hat{C} = \frac{9}{7} \text{ and } \hat{D} = \frac{4}{7} \text{ and the best line is } \frac{9}{7}t + \frac{4}{7} \neq$$

[30] The measurements b_1, \dots, b_m are given at distinct points t_1, \dots, t_m . Then the straight line $\hat{C} + \hat{D}t$ which minimizes E^2 comes from least squares:

$$A^T A \begin{bmatrix} \hat{C} \\ \hat{D} \end{bmatrix} = A^T b \text{ or } \begin{bmatrix} m \sum t_i \\ \sum t_i \sum t_i^2 \end{bmatrix} \begin{bmatrix} \hat{C} \\ \hat{D} \end{bmatrix} = \begin{bmatrix} \sum b_i \\ \sum t_i b_i \end{bmatrix}$$

Remark: Given a mixture of two radioactive chemicals with known half-lives λ and μ , we want to know their unknown amounts C and D : $b = C e^{-\lambda t} + D e^{-\mu t}$.

$$Ax = Ib \text{ is } \begin{cases} C e^{-\lambda t_1} + D e^{-\mu t_1} \approx b_1 \\ \vdots \\ C e^{-\lambda t_m} + D e^{-\mu t_m} \approx b_m. \end{cases}$$

The least-squares principle will give optimal \hat{C} and \hat{D} .

- o. But, if we knew the amounts C and D , and were trying to discover the decay rates λ and μ , this is a problem in nonlinear least squares. Setting the derivatives of E^2 to zero will give nonlinear equations for the optimal λ and μ .

Weighted Least Squares

The estimate \hat{x} of weight from two observations

$x = b_1$ and $x = b_2$:

$$\begin{bmatrix} 1 \\ 1 \end{bmatrix} [x] = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \Rightarrow A^T A \hat{x} = A^T b, \quad x \hat{x} = b_1 + b_2$$

Weighted error: $E^2 = w_1^2(x - b_1)^2 + w_2^2(x - b_2)^2$

$$\frac{dE^2}{dx} = 2 \left[w_1^2(x - b_1) + w_2^2(x - b_2) \right] = 0 \text{ at } \hat{x}_w = \frac{w_1^2 b_1 + w_2^2 b_2}{w_1^2 + w_2^2}$$

The least squares solution to $WAx = Wb$ is \hat{x}_w

Weighted Normal Equation: $(A^T W^T W A) \hat{x}_w = A^T W^T W b$.

- o. The projection $A \hat{x}_w$ is still the point in the column space that is closest to b under a new meaning of closeness. The perpendicularity test involves $(W\mathbf{y})^T (W\mathbf{x}) = 0$ instead of $\mathbf{y}^T \mathbf{x} = 0$. The matrix $W^T W$ appears in the middle. In this new sense, the projection $A \hat{x}_w$ and the error $b - A \hat{x}_w$ are again perpendicular.
- o. The inner product of \mathbf{x} and \mathbf{y} is generalized to $\mathbf{y}^T C \mathbf{x}$, where $C = W^T W$ is a symmetric matrix. For an orthogonal matrix $W = Q$, $C = Q^T Q = I$ and the inner product is not new.

For any invertible matrix W , these rules define a new inner product and length:

Weighted by W : $(\mathbf{x}, \mathbf{y})_W = (W\mathbf{y})^T (W\mathbf{x})$ and $\|\mathbf{x}\|_W = \|W\mathbf{x}\|$.

- o. If the errors in the b_i are independent of each other, and their variances are σ_i^2 , then the right weights are $w_i = 1/\sigma_i^2$. A more accurate measurement, with a smaller variance, gets a heavier weight.
- o. The observations may not be independent. Then W has off-diagonal terms. The best unbiased matrix $C = W^T W$ is the inverse of the covariance matrix - whose i,j entry is the expected value of (error in b_i) times (error in b_j). The main diagonal of C^{-1} contains the variances $\sigma_i^2 \rightarrow$ the average of (error in b_i)².

Ex 3:

Two bridge partners both guess (after the bidding) the total number of spades they hold. For each guess, the errors $-1, 0, 1$ might have equal probability $\frac{1}{3}$.

Then the expected error is zero and the variance is $\frac{2}{3}$.

$$E(e) = \frac{1}{3}(-1) + \frac{1}{3}(0) + \frac{1}{3}(1) = 0$$

$$E(e^2) = \frac{1}{3}(-1)^2 + \frac{1}{3}(0)^2 + \frac{1}{3}(1)^2 = \frac{2}{3}.$$

The two guesses are dependent but not identical. Say the chance that they are both too high or both too low is zero, but the chance of opposite errors is $\frac{1}{3}$. Then $E(e_1 e_2) = \frac{1}{3}(-1)$, and the inverse of the covariance matrix is $W^T W$:

$$\begin{bmatrix} E(e_1^2) & E(e_1 e_2) \\ E(e_1 e_2) & E(e_2^2) \end{bmatrix}^{-1} = \begin{bmatrix} \frac{2}{3} & -\frac{1}{3} \\ -\frac{1}{3} & \frac{2}{3} \end{bmatrix}^{-1} = C = W^T W.$$

This matrix goes into the middle of the weighted normal equations.

3.4 Orthogonal Bases and Gram-Schmidt

[3P] The vectors q_1, \dots, q_n are orthonormal if

$$q_i^T q_j = \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i = j \end{cases}$$

A matrix with orthonormal columns will be called Q .

The standard basis e_1, \dots, e_n consists of the columns of I . We can rotate these vectors to form other orthonormal bases. A subspace of \mathbb{R}^n may not contain the standard vectors e_i , but an orthonormal basis can be constructed for the space. This construction is known as Gram-Schmidt orthogonalization.

Orthogonal Matrices

[3Q] If Q has orthonormal columns, then $Q^T Q = I$.

An orthogonal matrix is a square matrix with orthonormal columns. Then $Q^T = Q^{-1}$. (Note that $Q^T Q = I$ even if Q is rectangular. But then Q^T is only a left-inverse.)

$$\text{Ex1: } Q = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \Rightarrow Q^T = Q^{-1} = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix}$$

$$\text{Ex2: } P = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \Rightarrow P^{-1} = P^T = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$$

- o. A reflection also forms an orthogonal matrix. Geometrically, an orthogonal Q is either a rotation or the product of a rotation and a reflection.

BR Multiplication by any Q preserves lengths:

$$\|Q\mathbf{x}\| = \|\mathbf{x}\| \text{ for every vector } \mathbf{x}.$$

It also preserves inner products and angles, since

$$(\mathbf{Q}\mathbf{x})^T(\mathbf{Q}\mathbf{y}) = \mathbf{x}^T Q^T Q \mathbf{y} = \mathbf{x}^T \mathbf{y}.$$

o. Write $\mathbf{l}\mathbf{b}$ as a combination $\mathbf{l}\mathbf{b} = x_1 q_1 + x_2 q_2 + \dots + x_n q_n$.

$$q_i^T \mathbf{l}\mathbf{b} = x_i q_i^T q_i \Rightarrow x_i = q_i^T \mathbf{l}\mathbf{b}$$

$$\therefore \mathbf{l}\mathbf{b} = (q_1^T \mathbf{l}\mathbf{b}) q_1 + (q_2^T \mathbf{l}\mathbf{b}) q_2 + \dots + (q_n^T \mathbf{l}\mathbf{b}) q_n$$

o. $x_1 q_1 + \dots + x_n q_n = \mathbf{l}\mathbf{b} \Rightarrow Q\mathbf{x} = \mathbf{l}\mathbf{b} \Rightarrow \mathbf{x} = Q^{-1}\mathbf{l}\mathbf{b} = Q^T \mathbf{l}\mathbf{b}$.

o. $\|\mathbf{l}\mathbf{b}\|^2 = (q_1^T \mathbf{l}\mathbf{b})^2 + (q_2^T \mathbf{l}\mathbf{b})^2 + \dots + (q_n^T \mathbf{l}\mathbf{b})^2 = \|Q^T \mathbf{l}\mathbf{b}\|^2$.

Remark: Since $Q^T = Q^{-1}$, we have $Q \cdot Q^T = I$. The rows of a square matrix are orthonormal whenever the columns are.

Rectangular Matrices with Orthonormal Columns

BR If Q has orthonormal columns, the least-square problem is easy.

$Q\mathbf{x} = \mathbf{l}\mathbf{b}$: rectangular system with no solution for most $\mathbf{l}\mathbf{b}$.

$Q^T Q \hat{\mathbf{x}} = Q^T \mathbf{l}\mathbf{b}$: normal equation for the best $\hat{\mathbf{x}}$ - in which $Q^T Q = I$.

$$\hat{\mathbf{x}} = Q^T \mathbf{l}\mathbf{b} : \hat{x}_i = q_i^T \mathbf{l}\mathbf{b}$$

$\mathbf{P} = Q\hat{\mathbf{x}}$: the projection of $\mathbf{l}\mathbf{b}$ is $(q_1^T \mathbf{l}\mathbf{b}) q_1 + \dots + (q_n^T \mathbf{l}\mathbf{b}) q_n$

$\mathbf{P} = Q Q^T \mathbf{l}\mathbf{b}$: the projection matrix is $\mathbf{P} = Q Q^T$

Ex 3: Projection of $\mathbf{l}\mathbf{b} = (x_1, y, z)$ onto the xy -plane is $\mathbf{P} = (x, y, 0)$

$$q_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \text{ and } (q_1^T \mathbf{l}\mathbf{b}) q_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}; q_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \text{ and } (q_2^T \mathbf{l}\mathbf{b}) q_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

$$\Rightarrow \mathbf{P} = q_1 q_1^T + q_2 q_2^T = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \text{ and } \mathbf{P} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} x \\ y \\ 0 \end{bmatrix}.$$

Ex4: When the measurement times average to zero, fitting a straight line leads to orthogonal columns.

Take $t_1 = -3, t_2 = 0, t_3 = 3$. The attempt to fit $y = C + Dt$ leads to

$$\begin{cases} C + Dt_1 = y_1 \\ C + Dt_2 = y_2 \\ C + Dt_3 = y_3 \end{cases} \quad \text{or} \quad \begin{bmatrix} 1 & -3 \\ 1 & 0 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} C \\ D \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

We can project \mathbf{y} separately onto each column:

$$\hat{C} = \frac{[1 \ 1 \ 1][y_1 \ y_2 \ y_3]^T}{1^2 + 1^2 + 1^2}, \quad \hat{D} = \frac{[-3 \ 0 \ 3][y_1 \ y_2 \ y_3]^T}{(-3)^2 + 0^2 + 3^2}$$

$\hat{C} = (y_1 + y_2 + y_3)/3$ is the best fit by a horizontal line, whereas $\hat{D}t$ is the best fit by a straight line through the origin. The sum $\hat{C} + \hat{D}t$ is the best fit by any straight line.

If the average $\bar{t} = (t_1 + \dots + t_m)/m$ is not zero, then the time can be shifted by \bar{t} . Let $y = C + Dt = c + d(t - \bar{t})$.

$$\hat{C} = \frac{[1 \ \dots \ 1][y_1 \ \dots \ y_m]^T}{1^2 + 1^2 + \dots + 1^2} = \frac{y_1 + \dots + y_m}{m}$$

$$\hat{d} = \frac{[(t_1 - \bar{t}) \ \dots \ (t_m - \bar{t})][y_1 \ \dots \ y_m]^T}{(t_1 - \bar{t})^2 + \dots + (t_m - \bar{t})^2} = \frac{\sum (t_i - \bar{t}) y_i}{\sum (t_i - \bar{t})^2}$$

The best \hat{C} is the mean and we also get a convenient formula for \hat{d} . The earlier $A^T A$ had the off-diagonal entries $\sum t_i$, and shifting the time by \bar{t} made these entries zero. The shift is an example of the Gram-Schmidt process.

The Gram-Schmidt Process

Given three independent vectors a, b, c , let $q_1 = a/\|a\|$.

Second vector $B = b - (q_1^T b) q_1$ and $q_2 = B/\|B\|$.

Third vector $C = c - (q_1^T c) q_1 - (q_2^T c) q_2$ and $q_3 = C/\|C\|$.

Ex5: Gram-Schmidt

$$a = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}, b = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, c = \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix} \Rightarrow q_1 = a/\sqrt{2} = \begin{bmatrix} 1/\sqrt{2} \\ 0 \\ 1/\sqrt{2} \end{bmatrix}$$

$$B = b - (q_1^T b) q_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} - \frac{1}{\sqrt{2}} \begin{bmatrix} 1/\sqrt{2} \\ 0 \\ 1/\sqrt{2} \end{bmatrix} = \begin{bmatrix} 1/2 \\ 0 \\ -1/2 \end{bmatrix} \Rightarrow q_2 = \begin{bmatrix} 1/\sqrt{2} \\ 0 \\ -1/\sqrt{2} \end{bmatrix}$$

$$C = c - (q_1^T c) q_1 - (q_2^T c) q_2 = \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix} - \sqrt{2} \begin{bmatrix} 1/\sqrt{2} \\ 0 \\ 1/\sqrt{2} \end{bmatrix} - \sqrt{2} \begin{bmatrix} 1/\sqrt{2} \\ 0 \\ -1/\sqrt{2} \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

Orthonormal basis $Q = [q_1 \ q_2 \ q_3] = \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} & 0 \\ 0 & 0 & 1 \\ 1/\sqrt{2} & -1/\sqrt{2} & 0 \end{bmatrix} \quad q_3$

[3T] The Gram-Schmidt process starts with independent vectors a_1, \dots, a_n , and ends with orthonormal vectors q_1, \dots, q_n . At step j , it subtracts from a_j its components in the directions q_1, \dots, q_{j-1} :

$$A_j = a_j - (q_1^T a_j) q_1 - \dots - (q_{j-1}^T a_j) q_{j-1}.$$

Then q_j is the unit vector $A_j/\|A_j\|$.

Remark on the calculation

It is easier to compute the orthogonal a, B, C without forcing their lengths to equal one. Then square roots enter only at the end, when dividing by those lengths.

The Factorization $A = QR$

$$Q = (q_1^T A) \cdot q_1$$

$$lb = (q_1^T lb) \cdot q_1 + (q_2^T lb) q_2$$

$$C = (q_1^T C) \cdot q_1 + (q_2^T C) q_2 + (q_3^T C) \cdot q_3$$

QR factors : $A = [A \mid b \mid C] = [q_1 \ q_2 \ q_3] \begin{bmatrix} q_1^T A & q_1^T b & q_1^T C \\ q_2^T A & q_2^T b & q_2^T C \\ q_3^T A & q_3^T b & q_3^T C \end{bmatrix} = QR$

Ex: $A = \begin{bmatrix} 1 & 1 & 2 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} & 0 \\ 0 & 0 & 1 \\ 1/\sqrt{2} & -1/\sqrt{2} & 0 \end{bmatrix} \begin{bmatrix} \sqrt{2} & 1/\sqrt{2} & \sqrt{2} \\ 1/\sqrt{2} & \sqrt{2} & 1 \end{bmatrix} = QR$

The length of a, b, c

3U Every $m \times n$ matrix with independent columns can be factored into $A = QR$. The columns of Q are orthonormal, and R is upper triangular and invertible. When $m=n$ and all matrices are square, Q becomes an orthogonal matrix.

The orthogonalization simplifies the least-squares problem $A\hat{x} = lb$. The normal equations become easier since $A^T A = R^T Q^T Q R = R^T R$.

The fundamental equation $A^T A \hat{x} = A^T lb$ simplifies to a triangular system

$$R^T R \hat{x} = R^T Q^T lb \quad \text{or} \quad R \hat{x} = Q^T lb.$$

Instead of solving $QR\hat{x} = lb$, we solve $R\hat{x} = Q^T lb$ by back-substitution. The real cost is the mn^2 operations of Gram-Schmidt, which are needed to find Q and R .

Function Spaces

① Lengths and Inner Products.

$$\|f\|^2 = \int_0^{2\pi} (f(x))^2 dx, \quad (f, g) = \int_0^{2\pi} f(x)g(x) dx$$

The Schwartz inequality is still satisfied:

$$|(f, g)| \leq \|f\| \cdot \|g\|$$

② Gram-Schmidt for Functions.

There is no interval $[a, b]$ on which $(1, x^2) = \int_a^b x^2 dx = 0$.

Therefore the closest parabola to $f(x)$ is not the sum of its projections onto $1, x, x^2$. On the interval $[0, 1]$,

$$A^T A = \begin{bmatrix} (1, 1) & (1, x) & (1, x^2) \\ (x, 1) & (x, x) & (x, x^2) \\ (x^2, 1) & (x^2, x) & (x^2, x^2) \end{bmatrix} = \begin{bmatrix} \int_0^1 1^2 dx & \int_0^1 x dx & \int_0^1 x^2 dx \\ \int_0^1 x dx & \int_0^1 x^2 dx & \int_0^1 x^3 dx \\ \int_0^1 x^2 dx & \int_0^1 x^3 dx & \int_0^1 x^4 dx \end{bmatrix} = \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{3} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} \end{bmatrix}$$

This is the ill-conditioned Hilbert matrix with a large inverse. The situation becomes impossible if we add a few more axes. It is virtually hopeless to solve $A^T A \mathbf{x} = A^T \mathbf{b}$ for the closest polynomial of degree 10. More precisely, it is hopeless to solve this by Gauss elimination. Every roundoff error would be amplified by more than 10^{13} .

The right idea is to switch to orthogonal axes (by Gram-Schmidt). On the interval $[-1, 1]$, we have

$$(1, x) = \int_{-1}^1 x dx = 0, \quad (x, x^2) = \int_{-1}^1 x^3 dx = 0.$$

Starting with $v_1 = 1, v_2 = x$,

$$\text{Orthogonalize } v_3 = x^2 - \frac{(1, x^2)}{(1, 1)} 1 - \frac{(x, x^2)}{(x, x)} x = x^2 - \frac{1}{3} x$$

$1, x, x^2 - \frac{1}{3} x$: Legendre polynomials

The closest polynomial is now computable by projecting onto each of the Legendre polynomials.