

Example 3.4 Orthogonality and Projections. We illustrate these concepts by continuing with the least squares problem in Examples 3.1 and 3.3. At the solution $\mathbf{x}^T = [1236, 1943, 2416]$, the residual vector,

$$\mathbf{r} = \mathbf{b} - \mathbf{A}\mathbf{x} = \mathbf{b} - \mathbf{y} = \begin{bmatrix} 1237 \\ 1941 \\ 2417 \\ 711 \\ 1177 \\ 475 \end{bmatrix} - \begin{bmatrix} 1236 \\ 1943 \\ 2416 \\ 707 \\ 1180 \\ 473 \end{bmatrix} = \begin{bmatrix} 1 \\ -2 \\ 1 \\ 4 \\ -3 \\ 2 \end{bmatrix},$$

is orthogonal to each column of \mathbf{A} , i.e., $\mathbf{A}^T \mathbf{r} = \mathbf{o}$. The orthogonal projector onto $\text{span}(\mathbf{A})$ is given by

$$\mathbf{P} = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T = \frac{1}{4} \begin{bmatrix} 2 & 1 & 1 & -1 & -1 & 0 \\ 1 & 2 & 1 & 1 & 0 & -1 \\ 1 & 1 & 2 & 0 & 1 & 1 \\ -1 & 1 & 0 & 2 & 1 & -1 \\ -1 & 0 & 1 & 1 & 2 & 1 \\ 0 & -1 & 1 & -1 & 1 & 2 \end{bmatrix},$$

and the orthogonal projector onto $\text{span}(\mathbf{A})^\perp$ is given by

$$\mathbf{P}_\perp = \mathbf{I} - \mathbf{P} = \frac{1}{4} \begin{bmatrix} 2 & -1 & -1 & 1 & 1 & 0 \\ -1 & 2 & -1 & -1 & 0 & 1 \\ -1 & -1 & 2 & 0 & -1 & -1 \\ 1 & -1 & 0 & 2 & -1 & 1 \\ 1 & 0 & -1 & -1 & 2 & -1 \\ 0 & 1 & -1 & 1 & -1 & 2 \end{bmatrix},$$

so that $\mathbf{b} = \mathbf{P}\mathbf{b} + \mathbf{P}_\perp \mathbf{b} = \mathbf{y} + \mathbf{r}$.

3.3 Sensitivity and Conditioning

We turn now to the sensitivity and conditioning of linear least squares problems. First, we must extend the notion of matrix condition number to include rectangular matrices. The definition of condition number for a square matrix given in Section 2.3.3 makes use of the matrix inverse. A nonsquare matrix \mathbf{A} does not have an inverse in the conventional sense, but it is possible to define a *pseudoinverse*, denoted by \mathbf{A}^+ , that behaves like an inverse in many respects (see Exercise 3.32). We will later see a more general definition that applies to any matrix, but for now we consider only matrices \mathbf{A} with full column rank, in which case $\mathbf{A}^T \mathbf{A}$ is nonsingular and we define the pseudoinverse of \mathbf{A} to be

$$\mathbf{A}^+ = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T.$$

Trivially, we see that $\mathbf{A}^+\mathbf{A} = \mathbf{I}$, and from Section 3.2.2 we see that $\mathbf{P} = \mathbf{A}\mathbf{A}^+$ is an orthogonal projector onto $\text{span}(\mathbf{A})$, so that the solution to the least squares problem $\mathbf{A}\mathbf{x} \cong \mathbf{b}$ is given by

$$\mathbf{x} = \mathbf{A}^+\mathbf{b}.$$

We now define the condition number of an $m \times n$ matrix with $\text{rank}(\mathbf{A}) = n$ to be

$$\text{cond}(\mathbf{A}) = \|\mathbf{A}\|_2 \cdot \|\mathbf{A}^+\|_2.$$

By convention, $\text{cond}(\mathbf{A}) = \infty$ if $\text{rank}(\mathbf{A}) < n$. Just as the condition number of a square matrix measures closeness to singularity, the condition number of a rectangular matrix measures closeness to rank deficiency.

Whereas the conditioning of a square linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$ depends only on the matrix \mathbf{A} , the conditioning of a least squares problem $\mathbf{A}\mathbf{x} \cong \mathbf{b}$ depends on the right-hand-side vector \mathbf{b} as well as the matrix \mathbf{A} , and thus $\text{cond}(\mathbf{A})$ alone does not suffice to characterize sensitivity. In particular, if \mathbf{b} lies near $\text{span}(\mathbf{A})$, then a small perturbation in \mathbf{b} changes $\mathbf{y} = \mathbf{P}\mathbf{b}$ relatively little. But if \mathbf{b} is nearly orthogonal to $\text{span}(\mathbf{A})$, on the other hand, then $\mathbf{y} = \mathbf{P}\mathbf{b}$ itself will be relatively small, so that a small change in \mathbf{b} can cause a relatively large change in \mathbf{y} , and hence in the least squares solution \mathbf{x} . Thus, for a given \mathbf{A} , we would expect a least squares problem with a \mathbf{b} that yields a large residual (i.e., a poor fit to the data) to be more sensitive than one with a small residual (i.e., a good fit to the data). An appropriate measure of the closeness of \mathbf{b} to $\text{span}(\mathbf{A})$ is the ratio

$$\frac{\|\mathbf{A}\mathbf{x}\|_2}{\|\mathbf{b}\|_2} = \frac{\|\mathbf{y}\|_2}{\|\mathbf{b}\|_2} = \cos(\theta),$$

where θ is the angle between \mathbf{b} and \mathbf{y} (see Fig. 3.2). Thus, we expect greater sensitivity when this ratio is small, so that θ is near $\pi/2$.

We now make a more quantitative assessment of the sensitivity of the solution \mathbf{x} of a least squares problem $\mathbf{A}\mathbf{x} \cong \mathbf{b}$, where \mathbf{A} has full column rank. For simplicity, we will consider perturbations in \mathbf{b} and \mathbf{A} separately. For a perturbed right-hand-side vector $\mathbf{b} + \Delta\mathbf{b}$, the perturbed solution is given by the normal equations

$$\mathbf{A}^T\mathbf{A}(\mathbf{x} + \Delta\mathbf{x}) = \mathbf{A}^T(\mathbf{b} + \Delta\mathbf{b}).$$

Because $\mathbf{A}^T\mathbf{A}\mathbf{x} = \mathbf{A}^T\mathbf{b}$, we then have

$$\mathbf{A}^T\mathbf{A}\Delta\mathbf{x} = \mathbf{A}^T\Delta\mathbf{b},$$

so that

$$\Delta\mathbf{x} = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\Delta\mathbf{b} = \mathbf{A}^+\Delta\mathbf{b}.$$

Taking norms, we obtain

$$\|\Delta\mathbf{x}\|_2 \leq \|\mathbf{A}^+\|_2 \cdot \|\Delta\mathbf{b}\|_2.$$

Dividing both sides by $\|\mathbf{x}\|_2$, we obtain the bound

$$\begin{aligned} \frac{\|\Delta\mathbf{x}\|_2}{\|\mathbf{x}\|_2} &\leq \|\mathbf{A}^+\|_2 \frac{\|\Delta\mathbf{b}\|_2}{\|\mathbf{x}\|_2} \\ &= \text{cond}(\mathbf{A}) \frac{\|\mathbf{b}\|_2}{\|\mathbf{A}\|_2 \cdot \|\mathbf{x}\|_2} \frac{\|\Delta\mathbf{b}\|_2}{\|\mathbf{b}\|_2} \\ &\leq \text{cond}(\mathbf{A}) \frac{\|\mathbf{b}\|_2}{\|\mathbf{Ax}\|_2} \frac{\|\Delta\mathbf{b}\|_2}{\|\mathbf{b}\|_2} \\ &= \text{cond}(\mathbf{A}) \frac{1}{\cos(\theta)} \frac{\|\Delta\mathbf{b}\|_2}{\|\mathbf{b}\|_2}. \end{aligned}$$

Thus, the condition number for the least squares solution \mathbf{x} with respect to perturbations in \mathbf{b} depends on $\text{cond}(\mathbf{A})$ and also on the angle θ between \mathbf{b} and \mathbf{Ax} (see Fig. 3.2). In particular, the condition number is approximately $\text{cond}(\mathbf{A})$ when the residual is small, so that $\cos(\theta) \approx 1$, but the condition number can be arbitrarily worse than $\text{cond}(\mathbf{A})$ when the residual is large, so that $\cos(\theta) \approx 0$.

For a perturbed matrix $\mathbf{A} + \mathbf{E}$, the perturbed solution is given by the normal equations

$$(\mathbf{A} + \mathbf{E})^T (\mathbf{A} + \mathbf{E})(\mathbf{x} + \Delta\mathbf{x}) = (\mathbf{A} + \mathbf{E})^T \mathbf{b}.$$

Noting that $\mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{b}$, dropping second-order terms (i.e., products of small perturbations), and rearranging, we then have

$$\begin{aligned} \mathbf{A}^T \mathbf{A} \Delta\mathbf{x} &\approx \mathbf{E}^T \mathbf{b} - \mathbf{E}^T \mathbf{Ax} - \mathbf{A}^T \mathbf{Ex} \\ &= \mathbf{E}^T (\mathbf{b} - \mathbf{Ax}) - \mathbf{A}^T \mathbf{Ex} \\ &= \mathbf{E}^T \mathbf{r} - \mathbf{A}^T \mathbf{Ex}, \end{aligned}$$

so that

$$\Delta\mathbf{x} \approx (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{E}^T \mathbf{r} - (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Ex} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{E}^T \mathbf{r} - \mathbf{A}^+ \mathbf{Ex}.$$

Taking norms, we obtain

$$\|\Delta\mathbf{x}\|_2 \lesssim \|(\mathbf{A}^T \mathbf{A})^{-1}\|_2 \cdot \|\mathbf{E}\|_2 \cdot \|\mathbf{r}\|_2 + \|\mathbf{A}^+\|_2 \cdot \|\mathbf{E}\|_2 \cdot \|\mathbf{x}\|_2.$$

Dividing both sides by $\|\mathbf{x}\|_2$ and using the fact that $\|\mathbf{A}\|_2^2 \cdot \|(\mathbf{A}^T \mathbf{A})^{-1}\|_2 = [\text{cond}(\mathbf{A})]^2$, we obtain the bound

$$\begin{aligned} \frac{\|\Delta\mathbf{x}\|_2}{\|\mathbf{x}\|_2} &\lesssim \|(\mathbf{A}^T \mathbf{A})^{-1}\|_2 \cdot \|\mathbf{E}\|_2 \frac{\|\mathbf{r}\|_2}{\|\mathbf{x}\|_2} + \|\mathbf{A}^+\|_2 \cdot \|\mathbf{E}\|_2 \\ &= [\text{cond}(\mathbf{A})]^2 \frac{\|\mathbf{E}\|_2}{\|\mathbf{A}\|_2} \frac{\|\mathbf{r}\|_2}{\|\mathbf{Ax}\|_2} + \text{cond}(\mathbf{A}) \frac{\|\mathbf{E}\|_2}{\|\mathbf{A}\|_2} \\ &\leq \left([\text{cond}(\mathbf{A})]^2 \frac{\|\mathbf{r}\|_2}{\|\mathbf{Ax}\|_2} + \text{cond}(\mathbf{A}) \right) \frac{\|\mathbf{E}\|_2}{\|\mathbf{A}\|_2} \\ &= ([\text{cond}(\mathbf{A})]^2 \tan(\theta) + \text{cond}(\mathbf{A})) \frac{\|\mathbf{E}\|_2}{\|\mathbf{A}\|_2}. \end{aligned}$$

Thus, the condition number for the least squares solution \mathbf{x} with respect to perturbations in \mathbf{A} depends on $\text{cond}(\mathbf{A})$ and also on the angle θ between \mathbf{b} and $\mathbf{A}\mathbf{x}$ (see Fig. 3.2). In particular, the condition number is approximately $\text{cond}(\mathbf{A})$ when the residual is small, so that $\tan(\theta) \approx 0$, but the condition number is effectively squared for a moderate residual, and becomes arbitrarily large when the residual is larger still. These sensitivity results will not only enable us to assess the quality of least squares solutions, but will also play an important role in understanding the relative merits of the various algorithms for computing such solutions numerically.

Example 3.5 Sensitivity and Conditioning. We again illustrate these concepts by continuing with Examples 3.1, 3.3, and 3.4. The pseudoinverse is given by

$$\mathbf{A}^+ = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T = \frac{1}{4} \begin{bmatrix} 2 & 1 & 1 & -1 & -1 & 0 \\ 1 & 2 & 1 & 1 & 0 & -1 \\ 1 & 1 & 2 & 0 & 1 & 1 \end{bmatrix}.$$

The matrix norms can be computed to obtain

$$\|\mathbf{A}\|_2 = 2, \quad \|\mathbf{A}^+\|_2 = 1,$$

so that

$$\text{cond}(\mathbf{A}) = \|\mathbf{A}\|_2 \cdot \|\mathbf{A}^+\|_2 = 2.$$

From the ratio

$$\cos(\theta) = \frac{\|\mathbf{A}\mathbf{x}\|_2}{\|\mathbf{b}\|_2} = \frac{\|\mathbf{y}\|_2}{\|\mathbf{b}\|_2} \approx \frac{3640.8761}{3640.8809} \approx 0.99999868,$$

we see that the angle θ between \mathbf{b} and \mathbf{y} is about 0.001625, which is very tiny, as expected for a problem with a very close fit to the data. From the small condition number and small angle θ , we conclude that this particular least squares problem is well-conditioned.

Example 3.6 Condition-Squaring Effect. Consider the matrix and perturbation

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ \epsilon & -\epsilon \\ 0 & 0 \end{bmatrix}, \quad \mathbf{E} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ -\epsilon & \epsilon \end{bmatrix},$$

where $\epsilon \ll 1$, say around $\sqrt{\epsilon_{\text{mach}}}$, for which we have

$$\text{cond}(\mathbf{A}) = 1/\epsilon, \quad \|\mathbf{E}\|_2 / \|\mathbf{A}\|_2 = \epsilon.$$

For the right-hand-side vector $\mathbf{b} = [1 \ 0 \ \epsilon]^T$, we have $\|\Delta\mathbf{x}\|_2 / \|\mathbf{x}\|_2 = 0.5$, so the relative perturbation in the solution is about equal to $\text{cond}(\mathbf{A})$ times the relative perturbation in \mathbf{A} . There is no condition-squaring effect for this right-hand side because the residual is small and $\tan(\theta) \approx \epsilon$, effectively suppressing the condition-squared term in the perturbation bound.