



# BLIND DECONVOLUTION: A MATTER OF NORM

By Dianne P. O'Leary

**W**E CONTINUE THE SPECTROSCOPY PROBLEM FROM THE LAST ISSUE, TRYING TO RECONSTRUCT A TRUE SPECTRUM FROM AN OBSERVED ONE. AGAIN, WE'LL USE

blind deconvolution, but this time we'll impose some constraints on the error matrix  $E$ , leading to a more difficult problem to solve but often a more useful reconstruction.

## Counts

Recall that we have the counts of Figure 1, measured by a spectrometer (to view the figure, see the solution to last issue's homework on page 63). Suppose we have particles whose energy ranges from  $e_{lo}$  to  $e_{high}$ , and we define some intermediate energy levels  $e_{l_0} = e_0 < e_1 < \dots < e_{n_b-1} < e_{n_b} = e_{high}$ . This creates  $n_b$  bins, where the count for the  $j$ th bin is the number of particles determined to have energies between  $e_{j-1}$  and  $e_j$ . Our spectrometer records  $n_b$  counts, one for each bin.

One way to model this system is to try to determine the correct counts  $f_j$ , the correct blurring given the measured counts  $g_j, j = 1, \dots, n_b$ .

This gives us a matrix equation  $(K + E)f \approx g + r$ , where  $E$  accounts for errors in modeling the spectrometer's blur, and  $r$  accounts for errors in counts. The matrix entry  $k_{j\ell}$  is the probability that a particle whose energy is in the interval  $[e_{\ell-1}, e_\ell]$  is assigned to bin  $j$  ( $j, \ell = 1, \dots, n_b$ ).

We assume that there is significant error in both  $g$  and  $K$ , but we note that in our data, each bin's properties are the same, so the rows of  $K$  have a pattern: for example, if the  $m \times n$  matrix  $K$  were  $5 \times 5$ , we would notice that

$$K = \begin{bmatrix} k_5 & k_4 & k_3 & k_2 & k_1 \\ k_6 & k_5 & k_4 & k_3 & k_2 \\ k_7 & k_6 & k_5 & k_4 & k_3 \\ k_8 & k_7 & k_6 & k_5 & k_4 \\ k_9 & k_8 & k_7 & k_6 & k_5 \end{bmatrix},$$

so there would be only  $m + n - 1 = 9$  distinct entries in  $K$ . A matrix of this form is called a *Toeplitz matrix*, and it is determined by the entries in its first row and column. Under this assumption, it might make sense to assume that the error matrix  $E$  also has this same structure, and therefore also depends on  $m + n - 1$  parameters instead of  $mn$ . We will gather these parameters in a vector called  $\hat{e}$ .

## Using the Euclidean Norm

In the previous homework, we solved the problem

$$\min_{E, r, f} \frac{1}{2} \|E\|_F^2 + \frac{1}{2} \|r\|_2^2, \tag{1}$$

where

$$r = g - (K + E)f. \tag{2}$$

With our new constraint that  $E$  be Toeplitz, our old solution isn't feasible, so we minimize the function, subject to the constraint, over all choices of  $f$  and  $\hat{e}$ .

Our goal is to find an effective algorithm to solve this problem, and we'll do it in several steps. The first is to derive a useful alternative expression for the matrix-vector product  $Ef$ .

### PROBLEM 1.

Show that  $Ef$  can be written as  $F\hat{e}$ , where  $\hat{e}$  is the vector that has entries  $\hat{e}_i$ , and  $F$  is a matrix whose entries depend on the entries in the vector  $f$ . In other words, find a matrix  $F$  so that  $Ef = F\hat{e}$ .

Let's use Newton's method to solve our minimization problem. Recall that if we minimize some function  $s(\mathbf{x})$ , then the Newton direction is the solution  $\mathbf{p}$  to the linear system

$$H(\mathbf{x})\mathbf{p} = -\nabla s(\mathbf{x}),$$

where  $\nabla s(\mathbf{x})$  is the gradient of  $s$  with respect to  $\mathbf{x}$ , and  $H(\mathbf{x})$

is the *Hessian* matrix, containing the second derivatives:  $b_{ij}(x) = \partial^2 s(\mathbf{x}) / \partial x_i \partial x_j$ . Let's derive a formula for  $\mathbf{p}$ .

### PROBLEM 2.

Derive the Newton direction for Equation 1. To do this, use the definitions of  $\mathbf{E}$  (in terms of  $\hat{\mathbf{e}}$ ) and  $\mathbf{r}$  (Equation 2), and then differentiate the function in Equation 1 with respect to  $\hat{\mathbf{e}}$  and  $\mathbf{f}$ .

Although the formula from Problem 2 is mathematically correct, it isn't the best computationally; Problem 3 provides a better alternative.

### PROBLEM 3.

Show that this Newton direction is approximately the same as the solution to the least squares problem

$$\min_{\Delta \hat{\mathbf{e}}, \Delta \mathbf{f}} \left\| \begin{bmatrix} \mathbf{F} & \mathbf{K} + \mathbf{E} \\ \mathbf{D} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \Delta \hat{\mathbf{e}} \\ \Delta \mathbf{f} \end{bmatrix} + \begin{bmatrix} -\mathbf{r} \\ \mathbf{D}\hat{\mathbf{e}} \end{bmatrix} \right\|_2,$$

where  $\mathbf{D}$  is a diagonal matrix of size  $(m+n-1) \times (m+n-1)$  with entries equal to the square roots of  $1, 2, \dots, n, \dots, n, n-1, \dots, 1$ . (In particular, the least squares solution is very close to the Newton direction if the model is good, so that  $\|\mathbf{r}\|$  is small.)

If we were to use this model on our spectroscopy data, the solution would be quite contaminated by error. Therefore, we make one further modification: we solve the problem

$$\min_{\mathbf{E}, \mathbf{f}} \frac{1}{2} \|\mathbf{E}\|_F^2 + \frac{1}{2} \|\mathbf{r}\|_2^2 + \frac{1}{2} \lambda^2 \|\mathbf{f}\|_2^2,$$

where the last term is a *Tikhonov regularization* term (as in "Image Deblurring: I Can See Clearly Now," vol. 5, no. 3, May/June 2003, pp. 82–84), with a fixed parameter  $\lambda$ , added to control the size of  $\mathbf{f}$ . In this case, the approximate Newton direction is computed from the solution to the least squares problem

$$\min_{\Delta \hat{\mathbf{e}}, \Delta \mathbf{f}} \left\| \begin{bmatrix} \mathbf{F} & \mathbf{K} + \mathbf{E} \\ \mathbf{D} & \mathbf{0} \\ \mathbf{0} & \lambda \mathbf{I} \end{bmatrix} \begin{bmatrix} \Delta \hat{\mathbf{e}} \\ \Delta \mathbf{f} \end{bmatrix} + \begin{bmatrix} -\mathbf{r} \\ \mathbf{D}\hat{\mathbf{e}} \\ \lambda \mathbf{f} \end{bmatrix} \right\|_2,$$

## Tools

You can use Matlab's `linprog` to solve the linear programming problems.

Choices of vector norms are discussed in many elementary textbooks, but James Ortega gives a particularly nice discussion.<sup>1</sup>

Armin Pruessner and I discuss the use of regularization plus norm choice.<sup>2</sup> References to earlier work using regularization or norm choice can be found in that article, too.

### References

1. J.M. Ortega, *Numerical Analysis: A Second Course*, Academic Press, 1972 (reprinted by SIAM Press).
2. A. Pruessner and D.P. O'Leary, "Blind Deconvolution Using a Regularized Structured Total Least Norm Approach," *SIAM J. Matrix Analysis and Applications*, vol. 24, 2003, pp. 1018–1037.

Now we put these pieces together to solve our problem.

### PROBLEM 4.

Write a function `[f, ehat, r, itn] = stls(K, g, lambda, tol)` that uses a variant of Newton's method to solve our Toeplitz-constrained problem in a stable and efficient way. Use the least squares problem from earlier to compute the approximate Newton direction. Start the iteration with  $\hat{\mathbf{e}} = \mathbf{0}$  and  $\mathbf{f} = \mathbf{0}$ . Stop the iteration when the norm of the approximate Newton step is smaller than `tol`, and set `itn` to the number of iterations. Provide documentation for your function. Use it on the data from the Web site ([www.computer.org/cise/homework](http://www.computer.org/cise/homework)), setting  $\lambda = 0.06$  and `tol` =  $10^{-3}$ . Plot the solution, and print the residual norm, the solution norm, and the number of iterations.

### Using Other Norms

If the errors in our data aren't distributed normally, we have several reasonable alternatives to the choice of the Euclidean norm for the minimization function. For example, instead of minimizing

$$\frac{1}{2} \|\mathbf{E}\|_F^2 + \frac{1}{2} \|\mathbf{r}\|_2^2 = \frac{1}{2} \left\| \begin{bmatrix} \mathbf{r} \\ \mathbf{D}\hat{\mathbf{e}} \end{bmatrix} \right\|_2^2,$$

we might instead minimize

$$\left\| \begin{bmatrix} \mathbf{r} \\ \mathbf{D}\hat{\mathbf{e}} \end{bmatrix} \right\|_p, \quad (3)$$

where if  $p = 1$ , the norm is defined as the sum of the absolute

values of the components, and if  $p = \infty$ , the norm is the maximum of the absolute values of the components. Either of these choices has the effect of reducing the effects of outliers in our measurements.

To derive an algorithm to solve this problem for  $p = 1$  or  $\infty$ , and to match our previous algorithm when  $p = 2$ , we reason this way. We need to satisfy the constraint

$$F\hat{\mathbf{e}} = E\mathbf{f},$$

even after we replace  $\mathbf{f}$  by  $\mathbf{f} + \Delta\mathbf{f}$  and  $\hat{\mathbf{e}}$  by  $\hat{\mathbf{e}} + \Delta\hat{\mathbf{e}}$ , so we require

$$(F + \Delta F)(\hat{\mathbf{e}} + \Delta\hat{\mathbf{e}}) = (E + \Delta E)(\mathbf{f} + \Delta\mathbf{f}),$$

where  $\Delta E$  is formed from  $\Delta\hat{\mathbf{e}}$ , and  $\Delta F$  is formed from  $\Delta\mathbf{f}$  so that  $\Delta F\hat{\mathbf{e}} = E\Delta\mathbf{f}$  and  $F\Delta\hat{\mathbf{e}} = \Delta E\mathbf{f}$ . This means that

$$\Delta F\Delta\hat{\mathbf{e}} = \Delta E\Delta\mathbf{f}.$$

Now let's examine the residual after we replace  $\mathbf{f}$  by  $\mathbf{f} + \Delta\mathbf{f}$  and  $\hat{\mathbf{e}}$  by  $\hat{\mathbf{e}} + \Delta\hat{\mathbf{e}}$ :

$$\begin{aligned} \mathbf{r}_{new} &= \mathbf{g} - (\mathbf{K} + \mathbf{E} + \Delta\mathbf{E})(\mathbf{f} + \Delta\mathbf{f}) \\ &= \mathbf{g} - (\mathbf{K} + \mathbf{E})\mathbf{f} - \Delta\mathbf{E}\mathbf{f} - (\mathbf{K} + \mathbf{E})\Delta\mathbf{f} - \Delta\mathbf{E}\Delta\mathbf{f}. \end{aligned}$$

If both  $\Delta\mathbf{f}$  and  $\Delta\hat{\mathbf{e}}$  are small, then the last term is negligible, and we can approximate

$$\mathbf{r}_{new} \approx \mathbf{r} - F\Delta\hat{\mathbf{e}} - (\mathbf{K} + \mathbf{E})\Delta\mathbf{f},$$

so that our minimization function (Equation 3) is approximated by

$$\left\| \begin{bmatrix} F & \mathbf{K} + \mathbf{E} \\ D & \mathbf{0} \end{bmatrix} \begin{bmatrix} \Delta\hat{\mathbf{e}} \\ \Delta\mathbf{f} \end{bmatrix} + \begin{bmatrix} -\mathbf{r} \\ D\hat{\mathbf{e}} \end{bmatrix} \right\|_p.$$

So, to compute our step, we need to minimize a function of this form; our next task is to develop an algorithm that does this.

### PROBLEM 5.

a. Show that when  $p = 1$ , minimizing

$$\left\| \begin{bmatrix} F & \mathbf{K} + \mathbf{E} \\ D & \mathbf{0} \end{bmatrix} \begin{bmatrix} \Delta\hat{\mathbf{e}} \\ \Delta\mathbf{f} \end{bmatrix} + \begin{bmatrix} -\mathbf{r} \\ D\hat{\mathbf{e}} \end{bmatrix} \right\|_1$$

over all choices of  $\Delta\mathbf{f}$  and  $\Delta\hat{\mathbf{e}}$  is equivalent to solving the linear programming problem

$$\min_{\Delta\hat{\mathbf{e}}, \Delta\mathbf{f}, \bar{\sigma}} \bar{\sigma} = \sum_{i=1}^m \bar{\sigma}_{1,i} + \sum_{i=1}^q \bar{\sigma}_{2,i} + \sum_{i=1}^n \bar{\sigma}_{3,i}$$

subject to

$$\begin{aligned} -\bar{\sigma}_1 &\leq F\Delta\hat{\mathbf{e}} + (\mathbf{K} + \mathbf{E})\Delta\mathbf{f} - \mathbf{r} &&\leq \bar{\sigma}_1 \\ -\bar{\sigma}_2 &\leq D\Delta\hat{\mathbf{e}} + D\hat{\mathbf{e}} &&\leq \bar{\sigma}_2 \\ -\bar{\sigma}_3 &\leq \lambda\Delta\mathbf{f} + \lambda\mathbf{f} &&\leq \bar{\sigma}_3 \end{aligned}$$

where  $\bar{\sigma}_1 \in \mathbf{R}^{m \times 1}$ ,  $\bar{\sigma}_2 \in \mathbf{R}^{q \times 1}$ , and  $\bar{\sigma}_3 \in \mathbf{R}^{n \times 1}$ .

b. Derive a similar linear program to compute  $\hat{\mathbf{e}} + \Delta\hat{\mathbf{e}}$ ,  $\Delta\mathbf{f}$  when  $p = \infty$ .

Let's see how the choice of norm affects our computed spectrum.

### PROBLEM 6.

Write a function `[f, ehat, r, itn] = stln1(K, g, lambda, tol)` that uses a variant of Newton's method to solve the problem when  $p = 1$ . Use the linear program to compute an approximate Newton direction. Start the iteration with  $\hat{\mathbf{e}} = \mathbf{0}$  and  $\mathbf{f} = \mathbf{0}$ . Stop the iteration when the norm of the approximate Newton step is smaller than `tol`, and set `itn` to the number of iterations. Use it on the data from the Web site ([www.computer.org/cise/homework](http://www.computer.org/cise/homework)), setting  $\lambda = 0.06$  and `tol` =  $10^{-3}$ . Plot the solution, and print the residual norm, the solution norm, and the number of iterations.

Repeat for the case  $p = \infty$ .

### Comparing Our Results

Recall that our goal is to reconstruct the spectrum of the particles fed into the spectrometer. Take some time now to compare the results we obtained using various problem formulations.

### PROBLEM 7.

Compare the results from Problems 4 and 6 with those of the last issue by answering these two questions: How does the quality of results compare? How does the amount of work compare?

## BLIND DECONVOLUTION: ERRORS, ERRORS EVERYWHERE

By Dianne P. O'Leary

CONSIDER FIGURE 1'S DATA, REPRESENTING COUNTS MEASURED BY A SPECTROMETER, MODELED BY  $\mathbf{Kf} \approx \mathbf{g}$ .

### PROBLEM 1.

Program the least squares algorithm and try it on the data of Figure 1 for various values of  $\tilde{n}$ , the number of singular values retained. The matrix  $\mathbf{K}$  is  $27 \times 22$ , and we assume that the true counts for the first two and the last three bins are zero. Note how ill-conditioned the original matrix  $\mathbf{K}$  is (by recording  $\sigma_1/\sigma_n$ ).

#### Answer:

See the posted program `problem1_and_3.m` at [www.computer.org/cise/homework](http://www.computer.org/cise/homework). The program isn't difficult, but it's important to make sure that you do the singular value decomposition (SVD) only once (at a cost of  $O(mn^3)$ ) and then form each of the trial solutions at a cost of  $O(n^2)$ . This requires using the associative law of multiplication.

In fact, it's possible to form each solution by updating a previous one (by adding the newly nonneglected term) at a cost of  $O(n)$ , which would be an even better algorithm, left as an exercise.

### PROBLEM 2.

Suppose we have the singular value decomposition (SVD) of  $[\mathbf{K}, \mathbf{g}] = \tilde{\mathbf{U}}\tilde{\Sigma}\tilde{\mathbf{V}}^T$ . Assume that  $\mathbf{K}$  has rank  $n$  and that  $\tilde{v}_{nm} > \tilde{v}_{n+1,n+1} \neq 0$ . Show that the solution to

$$\min_{\mathbf{E}, \mathbf{r}} \|\mathbf{E} \mathbf{r}\|_F,$$

subject to the constraint

$$[\mathbf{K} + \mathbf{E} \mathbf{g} + \mathbf{r}] \begin{bmatrix} \mathbf{f} \\ -1 \end{bmatrix} = \mathbf{0}$$

is

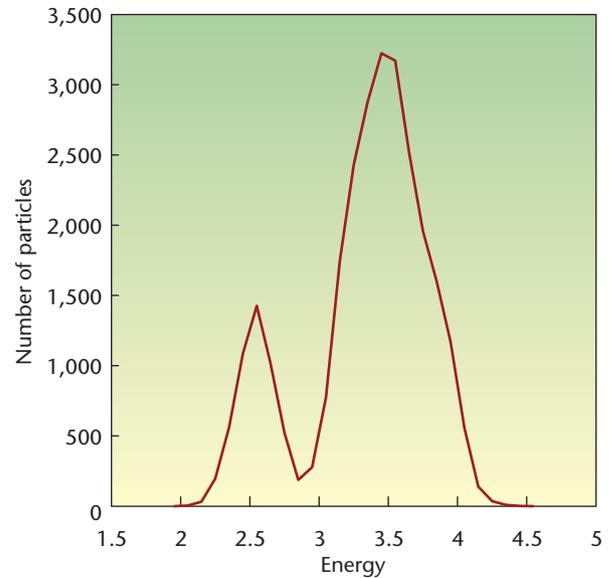


Figure 1. (Simulated) data from a spectrometer. Given that some particles have at most five different energy levels, determine these energies and the relative abundance of the particles.

$$[\mathbf{E} \mathbf{r}] = -\tilde{\sigma}_{n+1}\tilde{\mathbf{u}}_{n+1}\tilde{\mathbf{v}}_{n+1}^T,$$

with

$$\begin{bmatrix} \mathbf{f} \\ -1 \end{bmatrix} = -\frac{1}{\tilde{v}_{n+1,n+1}}\tilde{\mathbf{v}}_{n+1},$$

where  $\tilde{\mathbf{u}}_{n+1}$  is the  $(n+1)$ st column of  $\tilde{\mathbf{U}}$ , and  $\tilde{\mathbf{v}}_{n+1}$  is the  $(n+1)$ st column of  $\tilde{\mathbf{V}}$ .

Hint:

- First show that this solution satisfies the constraint and that the resulting  $\|\mathbf{E} \mathbf{r}\|_F = \tilde{\sigma}_{n+1}$ .
- Show that  $\|\tilde{\mathbf{U}}^T \mathbf{A} \tilde{\mathbf{V}}\|_F = \|\mathbf{A}\|_F$  for any matrix  $\mathbf{A}$  of size  $m \times (n+1)$ .
- Transform the problem to minimizing  $\|\tilde{\mathbf{E}} \tilde{\mathbf{r}}\|_F$  subject to  $(\tilde{\Sigma} + \tilde{\mathbf{E}})\tilde{\mathbf{f}} = \mathbf{0}$  for some vectors  $\tilde{\mathbf{f}}$  and  $\tilde{\mathbf{r}}$  and matrix  $\tilde{\mathbf{E}}$ , solve the problem in this coordinate system, and show that no solution gives a value of the minimization function smaller than  $\tilde{\sigma}_{n+1}$ .

Answer:

- We know that

$$[\mathbf{K} \mathbf{g}] = \sum_{i=1}^{n+1} \tilde{\sigma}_i \tilde{\mathbf{u}}_i \tilde{\mathbf{v}}_i^T,$$

so using the formula for  $[\mathbf{E} \mathbf{r}]$ , we see that

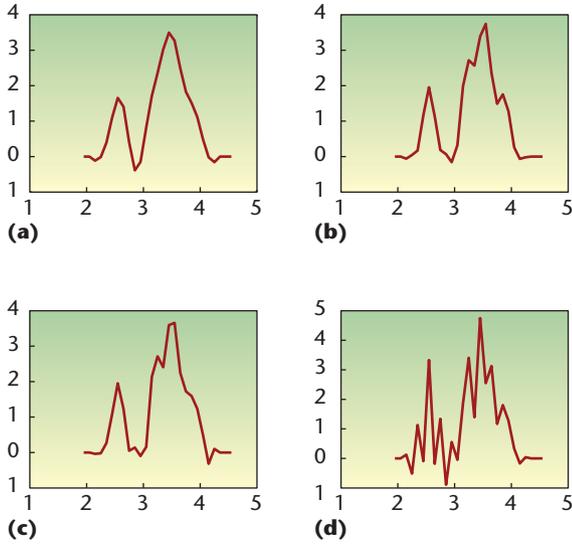


Figure A. Computed least squares solutions (counts versus energies) for various values of the cutoff parameter  $\tilde{n}$ . The solutions retain (a) 12 singular values, (b) 15 singular values, (c) 17 singular values, and (d) 21 singular values.

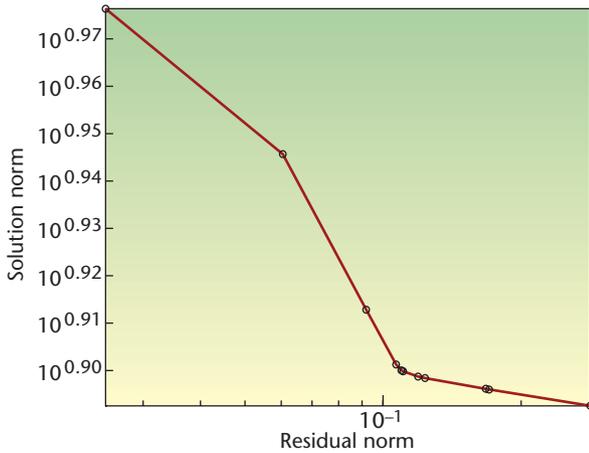


Figure B. The L-curve for least squares solutions.

$$[K \ g] + [E \ r] = \sum_{i=1}^{n+1} \tilde{\sigma}_i \tilde{\mathbf{u}}_i \tilde{\mathbf{v}}_i^T.$$

Now, since  $\tilde{\mathbf{v}}_{n+1}$  is orthogonal to  $\tilde{\mathbf{v}}_i$  for  $k = 1, \dots, n$ , it follows that

$$([K \ g] + [E \ r]) \begin{bmatrix} \mathbf{f} \\ -1 \end{bmatrix} = -\sum_{i=1}^n \tilde{\sigma}_i \tilde{\mathbf{u}}_i \tilde{\mathbf{v}}_i^T \frac{1}{\tilde{v}_{n+1, n+1}} \tilde{\mathbf{v}}_{n+1} = \mathbf{0}.$$

Note that  $\|[E, r]\|_F = \tilde{\sigma}_{n+1}$ .

b. This can be proven using the fact that  $\|A\|_F^2 = \text{tr}(A^T A)$ , where  $\text{tr}(B)$  is the trace of the matrix  $B$ , equal to the sum of its diagonal elements (or the sum of its eigenvalues). We can use the fact that  $\text{tr}(AB) = \text{tr}(BA)$ .

We can also prove it just from the definition of the Frobenius norm and the fact that  $\|Ux\|_2 = \|x\|_2$  for all vectors  $x$  and orthogonal matrices  $U$ . Using this fact, and letting  $\mathbf{a}_i$  be the  $i$ th column of  $A$ , we see that

$$\|UA\|_F^2 = \sum_{i=1}^n \|\mathbf{U}\mathbf{a}_i\|_2^2 = \sum_{i=1}^n \|\mathbf{a}_i\|_2^2 = \|A\|_F^2.$$

Similarly, letting  $\hat{\mathbf{a}}_i^T$  be the  $i$ th row of  $A$ ,

$$\|AV\|_F^2 = \sum_{i=1}^m \|\hat{\mathbf{a}}_i^T V\|_2^2 = \sum_{i=1}^m \|\hat{\mathbf{a}}_i\|_2^2 = \|A\|_F^2,$$

and the conclusion follows.

c. From the constraint

$$[K + E \ g + r] \begin{bmatrix} \mathbf{f} \\ -1 \end{bmatrix} = \mathbf{0},$$

we see that

$$\tilde{U}^T [K + E \ g + r] \tilde{V} \tilde{V}^T \begin{bmatrix} \mathbf{f} \\ -1 \end{bmatrix} = \mathbf{0},$$

so

$$(\tilde{\Sigma} + \tilde{E})\tilde{\mathbf{f}} = \mathbf{0}$$

where

$$\tilde{E} = \tilde{U}^T [E, r] \tilde{V} \text{ and } \tilde{\mathbf{f}} = \tilde{V}^T \begin{bmatrix} \mathbf{f} \\ -1 \end{bmatrix}.$$

From part b, we know that minimizing  $\|[E, r]\|_F$  is the same as minimizing  $\|\tilde{E}\|_F$ .

Therefore, to solve our problem, we want to make the smallest change to  $\tilde{\Sigma}$  that makes the matrix  $\tilde{\Sigma} + \tilde{E}$  rank deficient, so that the constraint can be satisfied by a nonzero  $\tilde{\mathbf{f}}$ . Changing the  $(n + 1, n + 1)$  element of  $\tilde{\Sigma}$  from  $\tilde{\sigma}_{n+1}$  to 0 certainly makes the constraint feasible (by setting the last component of  $\tilde{\mathbf{f}}$  nonzero and the other components zero). Any other change gives a bigger  $\|\tilde{E}\|_F$ . Thus, the smallest value of the minimization function is  $\tilde{\sigma}_{n+1}$ , and since we verified in part a that our solution has this value, we're finished.

If you don't find this argument convincing, we can be more precise by using a fact found in standard texts. For any matrix  $\mathbf{B}$  and vector  $\mathbf{z}$  for which  $\mathbf{Bz}$  is defined:  $\|\mathbf{Bz}\|_2 \leq \|\mathbf{B}\|_2 \|\mathbf{z}\|_2$ , where  $\|\mathbf{B}\|_2$  is defined to be the largest singular value of  $\mathbf{B}$ . Therefore,

- $\|\mathbf{B}\|_2 \leq \|\mathbf{B}\|_F$ , because we can see from part b and the SVD of  $\mathbf{B}$  that the Frobenius norm of  $\mathbf{B}$  is just the square root of the sum of the squares of its singular values.
- If  $(\tilde{\Sigma} + \tilde{\mathbf{E}})\tilde{\mathbf{f}} = \mathbf{0}$ , then  $\tilde{\Sigma}\tilde{\mathbf{f}} = -\tilde{\mathbf{E}}\tilde{\mathbf{f}}$ .
- $\tilde{\sigma}_{n+1}^2 \|\tilde{\mathbf{f}}\|_2^2 = \sum_{i=1}^{n+1} \tilde{\sigma}_{n+1}^2 \tilde{f}_i^2 \leq \sum_{i=1}^{n+1} \tilde{\sigma}_i^2 \tilde{f}_i^2 = \|\tilde{\Sigma}\tilde{\mathbf{f}}\|_2^2$ .
- Therefore,  $\tilde{\sigma}_{n+1} \|\tilde{\mathbf{f}}\|_2 \leq \|\tilde{\Sigma}\tilde{\mathbf{f}}\|_2 = \|\tilde{\mathbf{E}}\tilde{\mathbf{f}}\|_2 \leq \|\tilde{\mathbf{E}}\|_2 \|\tilde{\mathbf{f}}\|_2$ , so we conclude that  $\|\tilde{\mathbf{E}}\|_F \geq \tilde{\sigma}_{n+1}$ , and we have found a minimizing solution.

### PROBLEM 3.

Write a Matlab function to solve Model 2 using this truncated technique for various values of  $\tilde{n}$ . The input values should be  $\mathbf{K}$ ,  $\mathbf{g}$ , and a range of  $\tilde{n}$  values. Include appropriate documentation, and use your function to solve our problem.

#### Answer:

See the posted program `problem1_and_3.m` at [www.computer.org/cise/homework](http://www.computer.org/cise/homework).

### PROBLEM 4.

Write a brief summary of the results you obtained using Model 1 and Model 2 to solve the problem of Figure 1. Give your best estimate of the number of different peaks (energy levels) in the original data  $\mathbf{f}_{true}$ , the relative heights of the peaks, and the centers of the peaks. Make a convincing argument to justify your estimate and your choice of parameters ( $\delta$  and  $\tilde{n}$ ) for each method.

#### Answer:

##### Model 1: Least Squares

To estimate the variance of the error, note that in the least squares model, the last five components of the right-hand side  $\mathbf{U}^T \mathbf{g}$  can't be zeroed by any choice of  $\mathbf{f}$ , so if we believe the model, we believe that these are entirely due to error. All other components should have at least some data in addition to noise. Therefore, estimate the variance using the last five

## Tools

You can find more information about the L-curve (originally due to Chuck Lawson and Richard Hanson) in Per Christian Hansen's work;<sup>1</sup> the total least squares (TLS) truncation we used in this problem is discussed in a 1997 article.<sup>2</sup>

The standard reference on TLS is the book by Sabine Van Huffel and Joos Vanderwalle.<sup>3</sup>

An alternate way to "regularize" TLS is given in a 1999 work.

## References

1. P.C. Hansen, *Rank-Deficient and Discrete Ill-Posed Problems*, SIAM Press, 1998.
2. R. Fierro et al., "Regularization by Truncated Least Squares," *SIAM J. Scientific Computing*, vol. 18, 1997, pp. 1223–1241.
3. S. Van Huffel and J. Vanderwalle, *The Total Least Squares Problem*, SIAM Press, 1991.
4. G.H. Golub, P.C. Hansen, and D.P. O'Leary, "Tikhonov Regularization and Total Least Squares," *SIAM J. Matrix Analysis and Applications*, vol. 21, 1999, pp. 185–194.

to get  $\delta^2 = 1.2349 \times 10^{-4}$ .

The condition number of the matrix, the ratio of largest to smallest singular value, is 61.8455. This is a well-conditioned matrix! Most spectroscopy problems have a very ill-conditioned one (having a condition number of  $10^3$  or more). This is a clue that an error probably exists in the matrix, moving the small singular values away from zero.

We try various choices of  $\tilde{n}$ , the number of singular values retained, and show the results in Figure A. The discrepancy principle predicts the residual norm to be  $\delta\sqrt{m} = 0.0577$ . This is most closely matched by retaining 21 singular values, which gives seven peaks, contradicting the given information that there are at most five peaks. It also produces some rather large magnitude negative peaks, and we know that counts need to be nonnegative. Thus, the least squares model doesn't seem to fit the data well.

An alternate way to pick  $\tilde{n}$  is to use the L-curve. This is a plot of the log of the solution norm versus the log of the residual norm. It's called an L-curve because its shape often resembles that of the letter L. What we really want is a small residual and a solution norm that isn't unreasonably big. We could take the value of  $\tilde{n}$  at the corner of the L-curve, because if we take a smaller  $\tilde{n}$ , the residual norm increases quickly, and if we take a larger one, the solution norm increases quickly. This plot, shown in Figure B, advises that we should retain 15 to 17 singular values, and referring to

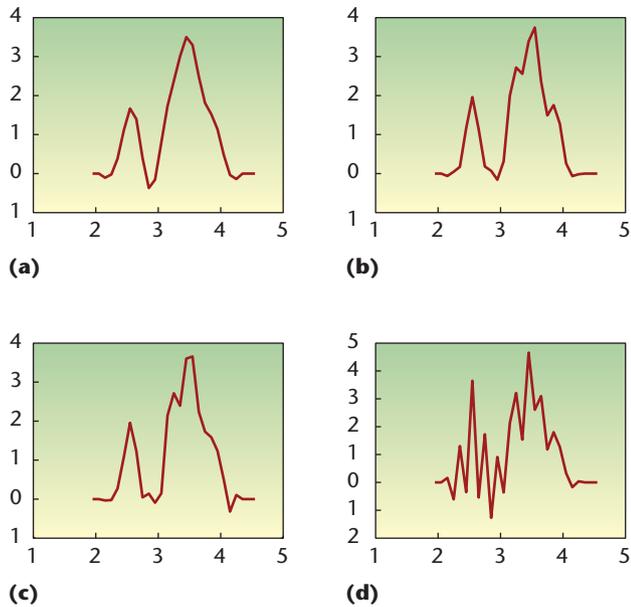


Figure C. Computed total least squares solutions (counts versus energies) for various values of the cutoff parameter  $\tilde{n}$ . The solutions retain (a) 12 singular values, (b) 15 singular values, (c) 17 singular values, and (d) 21 singular values.

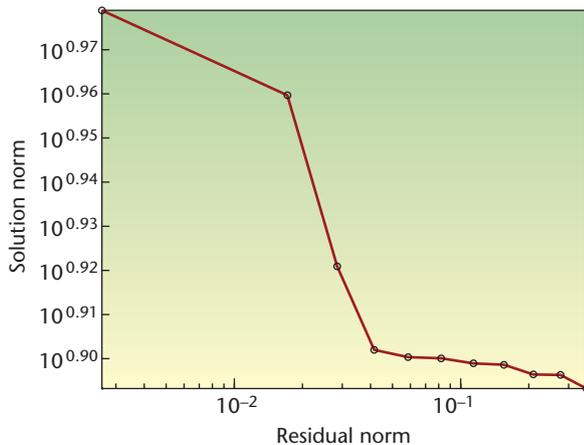


Figure D. The L-curve for total least squares solutions.

Figure A, this yields four peaks, consistent with our given information.

(The theoretical properties of the L-curve, as well as any other method that doesn't require the variance to be given in advance, aren't good, but this method is often more robust to errors in assumptions about the model, such as underestimating the variance or not accounting for errors in the matrix.)

An alternate heuristic is to look for a value of  $\tilde{n}$  that makes the residual look most like white noise, but because our error isn't normally distributed, this heuristic doesn't have

much meaning for our problem.

An excellent way to approach this problem is to generate your own test data, for which you know the true solution, and use it to gain insight into the choice of  $\tilde{n}$ .

**Model 2: Total Least Squares (TLS)**

Figure C shows sample solutions. The discrepancy principle doesn't give much insight for TLS, so we use more heuristic methods, giving us even less confidence in our choice.

For example, from Figure D we see that the L-curve corner occurs when 15 singular values are retained, giving a solution that looks very much like the L-curve least squares solution. Because the number of peaks is reasonable, and because there are only a small number of negative values in the solution (and these have small magnitude), we might accept this solution.

Now we need to extract the energies and estimated counts for the four types of particles. I have normalized them so that the count for the lowest energy peak is one. For the computed estimate to energy levels and counts,

Bin centers	2.55	3.25	3.55	3.85
Relative counts	1.00	1.39	1.91	0.90

A spectroscopist would actually estimate the counts by taking the integral under each of the four peaks, and then estimate the energy by the centroid of the peak, but this is difficult since three of the peaks aren't well separated.

The program used to generate the data is posted at [www.computer.org/cise/homework](http://www.computer.org/cise/homework). The variance of the error is  $10^{-4}$ . The true energy levels and counts are

Energy	2.54	3.25	3.53	3.85
Relative counts	1	1.5	2	1

So, despite all the errors, our computed solution estimates the energy levels to two digits and the relative counts to within 10 percent.

**Dianne P. O'Leary** is a professor of computer science and a faculty member in the Institute for Advanced Computer Studies and the Applied Mathematics Program at the University of Maryland. She received a BS in mathematics from Purdue University and a PhD in computer science from Stanford. She is a member of SIAM, the ACM, and the AWM. Contact her at [oleary@cs.umd.edu](mailto:oleary@cs.umd.edu); [www.cs.umd.edu/users/oleary/](http://www.cs.umd.edu/users/oleary/).