

## Stability Analysis of the Pseudo-Inverse

We have seen that if we make indirect observations  $\mathbf{y} \in \mathbb{R}^M$  of an unknown vector  $\mathbf{x}_0 \in \mathbb{R}^N$  through a  $M \times N$  matrix  $\mathbf{A}$ ,  $\mathbf{y} = \mathbf{A}\mathbf{x}_0$ , then applying the pseudo-inverse of  $\mathbf{A}$  gives us the least squares estimate of  $\mathbf{x}_0$ :

$$\hat{\mathbf{x}}_{\text{ls}} = \mathbf{A}^\dagger \mathbf{y} = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T \mathbf{y},$$

where  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  is the singular value decomposition (SVD) of  $\mathbf{A}$ .

We will now discuss what happens if our measurements contain *noise* — the analysis here will be very similar to when we looked at the stability of solving square sym+def systems, and in fact this is one of the main reasons we introduced the SVD.

Suppose we observe

$$\mathbf{y} = \mathbf{A}\mathbf{x}_0 + \mathbf{e},$$

where  $\mathbf{e} \in \mathbb{R}^M$  is an unknown perturbation. Say that we again apply the pseudo-inverse to  $\mathbf{y}$  in an attempt to recover  $\mathbf{x}$ :

$$\hat{\mathbf{x}}_{\text{ls}} = \mathbf{A}^\dagger \mathbf{y} = \mathbf{A}^\dagger \mathbf{A}\mathbf{x}_0 + \mathbf{A}^\dagger \mathbf{e}$$

What effect does the presence of the noise vector  $\mathbf{e}$  had on our estimate of  $\mathbf{x}_0$ ? We answer this question by comparing  $\hat{\mathbf{x}}_{\text{ls}}$  to the reconstruction we would obtain if we used standard least-squares on perfectly noise-free observations  $\mathbf{y}_{\text{clean}} = \mathbf{A}\mathbf{x}_0$ . This noise-free recon-

struction can be written as

$$\begin{aligned}
 \mathbf{x}_{\text{pinv}} &= \mathbf{A}^\dagger \mathbf{y}_{\text{clean}} = \mathbf{A}^\dagger \mathbf{A} \mathbf{x}_0 \\
 &= \mathbf{V} \boldsymbol{\Sigma}^{-1} \mathbf{U}^T \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T \mathbf{x}_0 \\
 &= \mathbf{V} \mathbf{V}^T \mathbf{x}_0 \\
 &= \sum_{r=1}^R \langle \mathbf{x}_0, \mathbf{v}_r \rangle \mathbf{v}_r.
 \end{aligned}$$

The vector  $\mathbf{x}_{\text{pinv}}$  is the orthogonal projection of  $\mathbf{x}_0$  onto the row space (everything orthogonal to the null space) of  $\mathbf{A}$ . If  $\mathbf{A}$  has full column rank ( $R = N$ ), then  $\mathbf{x}_{\text{pinv}} = \mathbf{x}_0$ . If not, then the application of  $\mathbf{A}$  destroys the part of  $\mathbf{x}_0$  that is not in  $\mathbf{x}_{\text{pinv}}$ , and so we only attempt to recover the “visible” components. In some sense,  $\mathbf{x}_{\text{pinv}}$  contains all of the components of  $\mathbf{x}_0$  that  $\mathbf{A}$  does not completely remove, and has them preserved perfectly.

The reconstruction error (relative to  $\mathbf{x}_{\text{pinv}}$  is)

$$\|\hat{\mathbf{x}}_{\text{ls}} - \mathbf{x}_{\text{pinv}}\|_2^2 = \|\mathbf{A}^\dagger \mathbf{e}\|_2^2 = \|\mathbf{V} \boldsymbol{\Sigma}^{-1} \mathbf{U}^T \mathbf{e}\|_2^2. \quad (1)$$

Now suppose for a moment that the error has unit norm,  $\|\mathbf{e}\|_2^2 = 1$ . Then the worst case for (1) is given by

$$\underset{\mathbf{e} \in \mathbb{R}^M}{\text{maximize}} \quad \|\mathbf{V} \boldsymbol{\Sigma}^{-1} \mathbf{U}^T \mathbf{e}\|_2^2 \quad \text{subject to} \quad \|\mathbf{e}\|_2 = 1.$$

Since the columns of  $\mathbf{U}$  are orthonormal,  $\|\mathbf{U}^T \mathbf{e}\|_2^2 \leq \|\mathbf{e}\|_2^2$ , and the above is equivalent to

$$\max_{\boldsymbol{\beta} \in \mathbb{R}^R: \|\boldsymbol{\beta}\|_2=1} \|\mathbf{V} \boldsymbol{\Sigma}^{-1} \boldsymbol{\beta}\|_2^2. \quad (2)$$

Also, for any vector  $\mathbf{z} \in \mathbb{R}^R$ , we have

$$\|\mathbf{V}\mathbf{z}\|_2^2 = \langle \mathbf{V}\mathbf{z}, \mathbf{V}\mathbf{z} \rangle = \langle \mathbf{z}, \mathbf{V}^T \mathbf{V}\mathbf{z} \rangle = \langle \mathbf{z}, \mathbf{z} \rangle = \|\mathbf{z}\|_2^2,$$

since the columns of  $\mathbf{V}$  are orthonormal. So we can simplify (2) to

$$\underset{\boldsymbol{\beta} \in \mathbb{R}^R}{\text{maximize}} \|\boldsymbol{\Sigma}^{-1}\boldsymbol{\beta}\|_2^2 \quad \text{subject to} \quad \|\boldsymbol{\beta}\|_2 = 1.$$

The worst case  $\boldsymbol{\beta}$  (you should verify this at home) will have a 1 in the entry corresponding to the largest entry in  $\boldsymbol{\Sigma}^{-1}$ , and will be zero everywhere else. Thus

$$\max_{\boldsymbol{\beta} \in \mathbb{R}^R: \|\boldsymbol{\beta}\|_2=1} \|\boldsymbol{\Sigma}^{-1}\boldsymbol{\beta}\|_2^2 = \max_{r=1, \dots, R} \sigma_r^{-2} = \frac{1}{\sigma_R^2}.$$

(Recall that by convention, we order the singular values so that  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_R$ .)

Returning to the reconstruction error (1), we now see that

$$\|\hat{\mathbf{x}}_{\text{ls}} - \mathbf{x}_{\text{pinv}}\|_2^2 = \|\mathbf{V}\boldsymbol{\Sigma}^{-1}\mathbf{U}^T\mathbf{e}\|_2^2 \leq \frac{1}{\sigma_R^2}\|\mathbf{e}\|_2^2.$$

Since  $\mathbf{U}$  is an  $M \times R$  matrix, it is possible when  $R < M$  that the reconstruction error is zero. This happens when  $\mathbf{e}$  is orthogonal to every column of  $\mathbf{U}$ , i.e.  $\mathbf{U}^T\mathbf{e} = \mathbf{0}$ . Putting this together with the work above means

$$0 \leq \frac{1}{\sigma_1^2}\|\mathbf{U}^T\mathbf{e}\|_2^2 \leq \|\hat{\mathbf{x}}_{\text{ls}} - \mathbf{x}_{\text{pinv}}\|_2^2 \leq \frac{1}{\sigma_R^2}\|\mathbf{U}^T\mathbf{e}\|_2^2 \leq \frac{1}{\sigma_R^2}\|\mathbf{e}\|_2^2.$$

Notice that if  $\sigma_R$  is small, the worst case reconstruction error can be **very bad**.

We can also relate the “average case” error to the singular values. Say that  $\mathbf{e}$  is additive Gaussian white noise, that is each entry  $e[m]$  is a random variable independent of all the other entries, and distributed

$$e[m] \sim \text{Normal}(0, \nu^2).$$

Then, as we have argued before, the average measurement error is

$$\mathbb{E}[\|\mathbf{e}\|_2^2] = M\nu^2,$$

and the average reconstruction error<sup>1</sup> is

$$\begin{aligned} \mathbb{E} \left[ \|\mathbf{A}^\dagger \mathbf{e}\|_2^2 \right] &= \nu^2 \cdot \text{trace}(\mathbf{A}^{\dagger T} \mathbf{A}^\dagger) = \nu^2 \cdot \left( \frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2} + \cdots + \frac{1}{\sigma_R^2} \right) \\ &= \frac{1}{M} \left( \frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2} + \cdots + \frac{1}{\sigma_R^2} \right) \cdot \mathbb{E}[\|\mathbf{e}\|_2^2]. \end{aligned}$$

Again, if  $\sigma_R$  is tiny,  $1/\sigma_R^2$  will dominate the sum above, and the average reconstruction error will be quite large.

**Exercise:** Let  $\mathbf{D}$  be a diagonal  $R \times R$  matrix whose diagonal elements are positive. Show that the maximizer  $\hat{\boldsymbol{\beta}}$  to

$$\underset{\boldsymbol{\beta} \in \mathbb{R}^R}{\text{maximize}} \quad \|\mathbf{D}\boldsymbol{\beta}\|_2^2 \quad \text{subject to} \quad \|\boldsymbol{\beta}\|_2 = 1$$

has a 1 in the entry corresponding to the largest diagonal element of  $\mathbf{D}$ , and is 0 elsewhere.

---

<sup>1</sup>We are using the fact that if  $\mathbf{e}$  is vector of iid Gaussian random variables,  $\mathbf{e} \sim \text{Normal}(\mathbf{0}, \nu^2 \mathbf{I})$ , then for any matrix  $\mathbf{M}$ ,  $\mathbb{E}[\|\mathbf{M}\mathbf{e}\|_2^2] = \nu^2 \text{trace}(\mathbf{M}^T \mathbf{M})$ . We leave this as an exercise (or maybe a homework!)

## Stable Reconstruction with the Truncated SVD

We have seen that if  $\mathbf{A}$  has very small singular values and we apply the pseudo-inverse in the presence of noise, the results can be disastrous. But it doesn't have to be this way. There are several ways to stabilize the pseudo-inverse. We start by discussing the simplest one, where we simply “cut out” the part of the reconstruction which is causing the problems.

As before, we are given noisy indirect observations of a vector  $\mathbf{x}$  through a  $M \times N$  matrix  $\mathbf{A}$ :

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{e}. \quad (3)$$

The matrix  $\mathbf{A}$  has SVD  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ , and pseudo-inverse  $\mathbf{A}^\dagger = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T$ . We can rewrite  $\mathbf{A}$  as a sum of rank-1 matrices:

$$\mathbf{A} = \sum_{r=1}^R \sigma_r \mathbf{u}_r \mathbf{v}_r^T,$$

where  $R$  is the rank of  $\mathbf{A}$ , the  $\sigma_r$  are the singular values, and  $\mathbf{u}_r \in \mathbb{R}^M$  and  $\mathbf{v}_r \in \mathbb{R}^N$  are columns of  $\mathbf{U}$  and  $\mathbf{V}$ , respectively. Similarly, we can write the pseudo-inverse as

$$\mathbf{A}^\dagger = \sum_{r=1}^R \frac{1}{\sigma_r} \mathbf{v}_r \mathbf{u}_r^T.$$

Given  $\mathbf{y}$  as above, we can write the least-squares estimate of  $\mathbf{x}$  from the noisy measurements as

$$\hat{\mathbf{x}}_{\text{ls}} = \mathbf{A}^\dagger \mathbf{y} = \sum_{r=1}^R \frac{1}{\sigma_r} \langle \mathbf{y}, \mathbf{u}_r \rangle \mathbf{v}_r. \quad (4)$$

As we can see (and have seen before) if any one of the  $\sigma_r$  are very small, the least-squares reconstruction can be a disaster.

A simple way to avoid this is to simply truncate the sum (4), leaving out the terms where  $\sigma_r$  is too small ( $1/\sigma_r$  is too big). Exactly how many terms to keep depends a great deal on the application, as there are competing interests. On the one hand, we want to ensure that each of the  $\sigma_r$  we include has an inverse of reasonable size, on the other, we want the reconstruction to be accurate (i.e. not to deviate from the noiseless least-squares solution by too much).

We form an approximation  $\mathbf{A}'$  to  $\mathbf{A}$  by taking

$$\mathbf{A}' = \sum_{r=1}^{R'} \sigma_r \mathbf{u}_r \mathbf{v}_r^T,$$

for some  $R' < R$ . Again, our final answer will depend on which  $R'$  we use, and choosing  $R'$  is often times something of an art. It is clear that the approximation  $\mathbf{A}'$  has rank  $R'$ . Note that the pseudo-inverse of  $\mathbf{A}'$  is also a truncated sum

$$\mathbf{A}'^\dagger = \sum_{r=1}^{R'} \frac{1}{\sigma_r} \mathbf{v}_r \mathbf{u}_r^T.$$

Given noisy data  $\mathbf{y}$  as in (3), we reconstruct  $\mathbf{x}$  by applying the truncated pseudo-inverse to  $\mathbf{y}$ :

$$\hat{\mathbf{x}}_{\text{trunc}} = \mathbf{A}'^\dagger \mathbf{y} = \sum_{r=1}^{R'} \frac{1}{\sigma_r} \langle \mathbf{y}, \mathbf{u}_r \rangle \mathbf{v}_r.$$

How good is this reconstruction? To answer this question, we will compare it to the noiseless least-squares reconstruction  $\mathbf{x}_{\text{pinv}} = \mathbf{A}^\dagger \mathbf{y}_{\text{clean}}$ , where  $\mathbf{y}_{\text{clean}} = \mathbf{A}\mathbf{x}$  are “noiseless” measurements of  $\mathbf{x}$ . The difference between these two is the reconstruction error (relative to  $\mathbf{x}_{\text{pinv}}$ ) as

$$\begin{aligned}\hat{\mathbf{x}}_{\text{trunc}} - \mathbf{x}_{\text{pinv}} &= \mathbf{A}'^\dagger \mathbf{y} - \mathbf{A}^\dagger \mathbf{A}\mathbf{x} \\ &= \mathbf{A}'^\dagger \mathbf{A}\mathbf{x} + \mathbf{A}'^\dagger \mathbf{e} - \mathbf{A}^\dagger \mathbf{A}\mathbf{x} \\ &= (\mathbf{A}'^\dagger - \mathbf{A}^\dagger) \mathbf{A}\mathbf{x} + \mathbf{A}'^\dagger \mathbf{e}.\end{aligned}$$

Proceeding further, we can write the matrix  $\mathbf{A}'^\dagger - \mathbf{A}^\dagger$  as

$$\mathbf{A}'^\dagger - \mathbf{A}^\dagger = \sum_{r=R'+1}^R -\frac{1}{\sigma_r} \mathbf{v}_r \mathbf{u}_r^\top,$$

and so the first term in the reconstruction error can be written as

$$\begin{aligned}(\mathbf{A}'^\dagger - \mathbf{A}^\dagger) \mathbf{A}\mathbf{x} &= \sum_{r=R'+1}^R -\frac{1}{\sigma_r} \langle \mathbf{A}\mathbf{x}, \mathbf{u}_r \rangle \mathbf{v}_r \\ &= \sum_{r=R'+1}^R -\frac{1}{\sigma_r} \left\langle \sum_{j=1}^R \sigma_j \langle \mathbf{x}, \mathbf{v}_j \rangle \mathbf{u}_j, \mathbf{u}_r \right\rangle \mathbf{v}_r \\ &= \sum_{r=R'+1}^R -\frac{1}{\sigma_r} \sum_{j=1}^R \sigma_j \langle \mathbf{x}, \mathbf{v}_j \rangle \langle \mathbf{u}_j, \mathbf{u}_r \rangle \mathbf{v}_r \\ &= \sum_{r=R'+1}^R -\langle \mathbf{x}, \mathbf{v}_r \rangle \mathbf{v}_r \quad (\text{since } \langle \mathbf{u}_r, \mathbf{u}_j \rangle = 0 \text{ unless } j = r).\end{aligned}$$

The second term in the reconstruction error can also be expanded against the  $\mathbf{v}_r$ :

$$\mathbf{A}'^\dagger \mathbf{e} = \sum_{r=1}^{R'} \frac{1}{\sigma_r} \langle \mathbf{e}, \mathbf{u}_r \rangle \mathbf{v}_r.$$

Combining these expressions, the reconstruction error can be written

$$\begin{aligned}\hat{\mathbf{x}}_{\text{trunc}} - \mathbf{x}_{\text{pinv}} &= \underbrace{\sum_{r=1}^{R'} \frac{1}{\sigma_r} \langle \mathbf{e}, \mathbf{u}_r \rangle \mathbf{v}_r}_{\text{Noise error}} + \underbrace{\sum_{k=R'+1}^R -\langle \mathbf{x}, \mathbf{v}_k \rangle \mathbf{v}_k}_{\text{Approximation error}} \\ &= \text{Noise error} + \text{Approximation error}.\end{aligned}$$

Since the  $\mathbf{v}_r$  are mutually orthogonal, and the two sums run over disjoint index sets, the noise error and the approximation error will be orthogonal. Also

$$\begin{aligned}\|\hat{\mathbf{x}}_{\text{trunc}} - \mathbf{x}_{\text{pinv}}\|_2^2 &= \|\text{Noise error}\|_2^2 + \|\text{Approximation error}\|_2^2 \\ &= \sum_{r=1}^{R'} \frac{1}{\sigma_r^2} |\langle \mathbf{e}, \mathbf{u}_r \rangle|^2 + \sum_{r=R'+1}^R |\langle \mathbf{x}, \mathbf{v}_r \rangle|^2.\end{aligned}$$

The reconstruction error, then, is signal dependent and will depend on how much of the vector  $\mathbf{x}$  is concentrated in the subspace spanned by  $\mathbf{v}_{R'+1}, \dots, \mathbf{v}_R$ . We will lose everything in this subspace; if it contains a significant part of  $\mathbf{x}$ , then there is not much least-squares can do for you.

The worst-case noise error occurs when  $\mathbf{e}$  is aligned with  $\mathbf{u}_{R'}$ :

$$\|\text{Noise error}\|_2^2 = \sum_{r=1}^{R'} \frac{1}{\sigma_r^2} |\langle \mathbf{e}, \mathbf{u}_r \rangle|^2 \leq \frac{1}{\sigma_{R'}^2} \cdot \|\mathbf{e}\|_2^2.$$

As seen before, if the error  $\mathbf{e}$  is random, this bound is a bit pessimistic. Specifically, if each entry of  $\mathbf{e}$  is an independent identically distributed Normal random variable with mean zero and variance  $\nu^2$ , then the expected noise error in the reconstruction will be

$$\mathbb{E}[\|\text{Noise error}\|_2^2] = \frac{1}{M} \left( \frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2} + \dots + \frac{1}{\sigma_{R'}^2} \right) \cdot \mathbb{E}[\|\mathbf{e}\|_2^2].$$



## Stable Reconstruction using Tikhonov Regularization

Tikhonov<sup>2</sup> regularization is another way to stabilize the least-squares recovery. It has the nice features that: 1) it can be interpreted using optimization, and 2) it can be computed without direct knowledge of the SVD of  $\mathbf{A}$ .

Recall that we motivated the pseudo-inverse by showing that  $\hat{\mathbf{x}}_{LS} = \mathbf{A}^\dagger \mathbf{y}$  is a solution to

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2. \quad (5)$$

When  $\mathbf{A}$  has full column rank,  $\hat{\mathbf{x}}_{LS}$  is the unique solution, otherwise it is the solution with smallest energy. When  $\mathbf{A}$  has full column rank but has singular values which are very small, huge variations in  $\mathbf{x}$  (in directions of the singular vectors  $\mathbf{v}_k$  corresponding to the tiny  $\sigma_k$ ) can have very little effect on the residual  $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$ . As such, the solution to (5) can have wildly inaccurate components in the presence of even mild noise.

One way to counteract this problem is to modify (5) with a **regularization** term that penalizes the size of the solution  $\|\mathbf{x}\|_2^2$  as well as the residual error  $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$ :

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \delta \|\mathbf{x}\|_2^2. \quad (6)$$

The parameter  $\delta > 0$  gives us a trade-off between accuracy and regularization; we want to choose  $\delta$  small enough so that the residual

---

<sup>2</sup>Andrey Tikhonov (1906-1993) was a 20<sup>th</sup> century Russian mathematician.

for the solution of (6) is close to that of (5), and large enough so that the problem is well-conditioned.

Just as with (5), which is solved by applying the pseudo-inverse to  $\mathbf{y}$ , we can write the solution to (6) in closed form. To see this, recall that we can decompose any  $\mathbf{x} \in \mathbb{R}^N$  as

$$\mathbf{x} = \mathbf{V}\boldsymbol{\alpha} + \mathbf{V}_0\boldsymbol{\alpha}_0,$$

where  $\mathbf{V}$  is the  $N \times R$  matrix (with orthonormal columns) used in the SVD of  $\mathbf{A}$ , and  $\mathbf{V}_0$  is a  $N \times N - R$  matrix whose columns are an orthogonal basis for the null space of  $\mathbf{A}$ . This means that the columns of  $\mathbf{V}_0$  are orthogonal to each other and all of the columns of  $\mathbf{V}$ . Similarly, we can decompose  $\mathbf{y}$  as

$$\mathbf{y} = \mathbf{U}\boldsymbol{\beta} + \mathbf{U}_0\boldsymbol{\beta}_0,$$

where  $\mathbf{U}$  is the  $M \times R$  matrix used in the SVD of  $\mathbf{A}$ , and the columns of  $\mathbf{U}_0$  are an orthogonal basis for the left null space of  $\mathbf{A}$  (everything in  $\mathbb{R}^M$  that is not in the range of  $\mathbf{A}$ ).

For any  $\mathbf{x}$ , we can write

$$\begin{aligned} \mathbf{y} - \mathbf{A}\mathbf{x} &= \mathbf{U}\boldsymbol{\beta} + \mathbf{U}_0\boldsymbol{\beta}_0 - \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T(\mathbf{V}\boldsymbol{\alpha} + \mathbf{V}_0\boldsymbol{\alpha}_0) \\ &= \mathbf{U}(\boldsymbol{\beta} - \boldsymbol{\Sigma}\boldsymbol{\alpha}) + \mathbf{U}_0\boldsymbol{\beta}_0. \end{aligned}$$

Since the columns of  $\mathbf{U}$  are orthonormal,  $\mathbf{U}^T\mathbf{U} = \mathbf{I}$ , and also  $\mathbf{U}_0^T\mathbf{U}_0 = \mathbf{I}$ , and  $\mathbf{U}^T\mathbf{U}_0 = \mathbf{0}$ , we have

$$\begin{aligned} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 &= \langle \mathbf{U}(\boldsymbol{\beta} - \boldsymbol{\Sigma}\boldsymbol{\alpha}) + \mathbf{U}_0\boldsymbol{\beta}_0, \mathbf{U}(\boldsymbol{\beta} - \boldsymbol{\Sigma}\boldsymbol{\alpha}) + \mathbf{U}_0\boldsymbol{\beta}_0 \rangle \\ &= \|\boldsymbol{\beta} - \boldsymbol{\Sigma}\boldsymbol{\alpha}\|_2^2 + \|\boldsymbol{\beta}_0\|_2^2, \end{aligned}$$

and

$$\|\mathbf{x}\|_2^2 = \|\boldsymbol{\alpha}\|_2^2 + \|\boldsymbol{\alpha}_0\|_2^2.$$

Using these facts, we can write the functional in (6) as

$$\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \delta\|\mathbf{x}\|^2 = \|\boldsymbol{\beta} - \boldsymbol{\Sigma}\boldsymbol{\alpha}\|_2^2 + \delta\|\boldsymbol{\alpha}\|_2^2 + \delta\|\boldsymbol{\alpha}_0\|_2^2. \quad (7)$$

We want to choose  $\boldsymbol{\alpha}$  and  $\boldsymbol{\alpha}_0$  that minimize (7). It is clear that, just as in the standard least-squares problem, we need  $\boldsymbol{\alpha}_0 = \mathbf{0}$ . The part of the functional that depends on  $\boldsymbol{\alpha}$  can be rewritten as

$$\|\boldsymbol{\beta} - \boldsymbol{\Sigma}\boldsymbol{\alpha}\|_2^2 + \delta\|\boldsymbol{\alpha}\|_2^2 = \sum_{k=1}^R (\beta[k] - \sigma_k\alpha[k])^2 + \delta\alpha[k]^2. \quad (8)$$

We can minimize this sum simply by minimizing each term independently. Since

$$\frac{d}{d\alpha[k]} [(\beta[k] - \sigma_k\alpha[k])^2 + \delta\alpha[k]^2] = -2\beta[k]\sigma_k + 2\sigma_k^2\alpha[k] + 2\delta\alpha[k],$$

we need

$$\alpha[k] = \frac{\sigma_k}{\sigma_k^2 + \delta} \beta[k].$$

Putting this back in vector form, (8) is minimized by

$$\hat{\boldsymbol{\alpha}}_{\text{tik}} = (\boldsymbol{\Sigma}^2 + \delta\mathbf{I})^{-1}\boldsymbol{\Sigma}\boldsymbol{\beta},$$

and so the minimizer to (6) is

$$\begin{aligned} \hat{\mathbf{x}}_{\text{tik}} &= \mathbf{V}\hat{\boldsymbol{\alpha}}_{\text{tik}} \\ &= \mathbf{V}(\boldsymbol{\Sigma}^2 + \delta\mathbf{I})^{-1}\boldsymbol{\Sigma}\mathbf{U}^T\mathbf{y}. \end{aligned} \quad (9)$$

We can get a better feel for what Tikhonov regularization is doing by comparing it directly to the pseudo-inverse. The least-squares

reconstruction  $\hat{\mathbf{x}}_{\text{ls}}$  can be written as

$$\begin{aligned}\hat{\mathbf{x}}_{\text{ls}} &= \mathbf{V}\boldsymbol{\Sigma}^{-1}\mathbf{U}^T\mathbf{y} \\ &= \sum_{r=1}^R \frac{1}{\sigma_r} \langle \mathbf{y}, \mathbf{u}_r \rangle \mathbf{v}_r,\end{aligned}$$

while the Tikhonov reconstruction  $\hat{\mathbf{x}}_{\text{tik}}$  derived above is

$$\hat{\mathbf{x}}_{\text{tik}} = \sum_{r=1}^R \frac{\sigma_r}{\sigma_r^2 + \delta} \langle \mathbf{y}, \mathbf{u}_r \rangle \mathbf{v}_r. \quad (10)$$

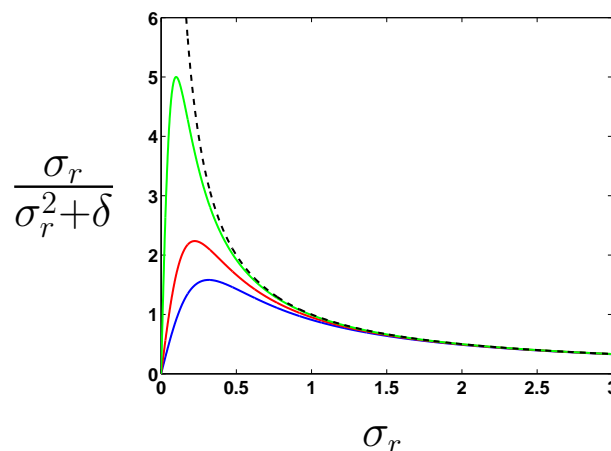
Notice that when  $\sigma_r$  is much larger than  $\delta$ ,

$$\frac{\sigma_r}{\sigma_r^2 + \delta} \approx \frac{1}{\sigma_r}, \quad \sigma_r \gg \delta,$$

but when  $\sigma_r$  is small

$$\frac{\sigma_r}{\sigma_r^2 + \delta} \approx 0, \quad \sigma_r \ll \delta.$$

Thus the Tikhonov reconstruction modifies the important parts (components where the  $\sigma_r$  are large) of the pseudo-inverse very little, while ensuring that the unimportant parts (components where the  $\sigma_r$  are small) affect the solution only by a very small amount. This **damp-**  
**ing** of the singular values, is illustrated below.



Above, we see the damped multipliers  $\sigma_r/(\sigma_r^2 + \delta)$  versus  $\sigma_r$  for  $\delta = 0.1$  (blue),  $\delta = 0.05$  (red), and  $\delta = 0.01$  (green). The black dotted line is  $1/\sigma_r$ , the least-squares multiplier. Notice that for large  $\sigma_r$  ( $\sigma_r > 2\sqrt{\delta}$ , say), the damping has almost no effect.

This damping makes the Tikhonov reconstruction exceptionally stable; large multipliers never appear in the reconstruction (10). In fact it is easy to check that

$$\frac{\sigma_r}{\sigma_r^2 + \delta} \leq \frac{1}{2\sqrt{\delta}}$$

no matter the value of  $\sigma_r$ .

## Tikhonov Error Analysis

Given noisy observations  $\mathbf{y} = \mathbf{A}\mathbf{x}_0 + \mathbf{e}$ , how well will Tikhonov regularization work? The answer to this question depends on multiple factors including the choice of  $\delta$ , the nature of the perturbation  $\mathbf{e}$ , and how well  $\mathbf{x}_0$  can be approximated using a linear combination of the singular vectors  $\mathbf{v}_r$  corresponding to the large (relative to  $\delta$ ) singular values. Since a closed-form expression for the solution to (6) exists, we can quantify these trade-offs precisely.

We compare the Tikhonov reconstruction to the reconstruction we would obtain if we used standard least-squares on perfectly noise-free observations  $\mathbf{y}_{\text{clean}} = \mathbf{A}\mathbf{x}_0$ . This noise-free reconstruction can

be written as

$$\begin{aligned}
 \mathbf{x}_{\text{pinv}} &= \mathbf{A}^\dagger \mathbf{y}_{\text{clean}} = \mathbf{A}^\dagger \mathbf{A} \mathbf{x}_0 \\
 &= \mathbf{V} \boldsymbol{\Sigma}^{-1} \mathbf{U}^T \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T \mathbf{x}_0 \\
 &= \mathbf{V} \mathbf{V}^T \mathbf{x}_0 \\
 &= \sum_{r=1}^R \langle \mathbf{x}_0, \mathbf{v}_r \rangle \mathbf{v}_r.
 \end{aligned}$$

The vector  $\mathbf{x}_{\text{pinv}}$  is the orthogonal projection of  $\mathbf{x}_0$  onto the row space (everything orthogonal to the null space) of  $\mathbf{A}$ . If  $\mathbf{A}$  has full column rank, then  $\mathbf{x}_{\text{pinv}} = \mathbf{x}_0$ . If not, then the application of  $\mathbf{A}$  destroys the part of  $\mathbf{x}_0$  that is not in  $\mathbf{x}_{\text{pinv}}$ , and so we only attempt to recover the “visible” components. In some sense,  $\mathbf{x}_{\text{pinv}}$  contains all of the components of  $\mathbf{x}_0$  that we could ever hope to recover, and has them preserved perfectly.

The Tikhonov regularized solution is given by

$$\begin{aligned}
 \hat{\mathbf{x}}_{\text{tik}} &= \sum_{r=1}^R \frac{\sigma_r}{\sigma_r^2 + \delta} \langle \mathbf{y}, \mathbf{u}_r \rangle \mathbf{v}_r \\
 &= \sum_{r=1}^R \frac{\sigma_r}{\sigma_r^2 + \delta} \langle \mathbf{e}, \mathbf{u}_r \rangle \mathbf{v}_r + \sum_{r=1}^R \frac{\sigma_r}{\sigma_r^2 + \delta} \langle \mathbf{A} \mathbf{x}_0, \mathbf{u}_r \rangle \mathbf{v}_r \\
 &= \sum_{r=1}^R \frac{\sigma_r}{\sigma_r^2 + \delta} \langle \mathbf{e}, \mathbf{u}_r \rangle \mathbf{v}_r + \sum_{r=1}^R \frac{\sigma_r^2}{\sigma_r^2 + \delta} \langle \mathbf{x}_0, \mathbf{v}_r \rangle \mathbf{v}_r,
 \end{aligned}$$

and so the reconstruction error, relative to the best possible recon-

struction  $\mathbf{x}_{\text{pinv}}$ , is

$$\begin{aligned}\hat{\mathbf{x}}_{\text{tik}} - \mathbf{x}_{\text{pinv}} &= \sum_{r=1}^R \frac{\sigma_r}{\sigma_r^2 + \delta} \langle \mathbf{e}, \mathbf{u}_r \rangle \mathbf{v}_r + \sum_{r=1}^R \left( \frac{\sigma_r^2}{\sigma_r^2 + \delta} - 1 \right) \langle \mathbf{x}_0, \mathbf{v}_r \rangle \mathbf{v}_r \\ &= \underbrace{\sum_{r=1}^R \frac{\sigma_r}{\sigma_r^2 + \delta} \langle \mathbf{e}, \mathbf{u}_r \rangle \mathbf{v}_r}_{\text{Noise error}} + \underbrace{\sum_{r=1}^R \frac{-\delta}{\sigma_r^2 + \delta} \langle \mathbf{x}_0, \mathbf{v}_r \rangle \mathbf{v}_r}_{\text{Approximation error}}.\end{aligned}$$

The approximation error is signal dependent, and depends on  $\delta$ . Since the  $\mathbf{v}_r$  are orthonormal,

$$\|\text{Approximation error}\|_2^2 = \sum_{r=1}^R \frac{\delta^2}{(\sigma_r^2 + \delta)^2} |\langle \mathbf{x}_0, \mathbf{v}_r \rangle|^2.$$

Note that for the components much smaller than  $\delta$ ,

$$\sigma_r^2 \ll \delta \quad \Rightarrow \quad \frac{\delta^2}{(\sigma_r^2 + \delta)^2} \approx 1,$$

so this portion of the approximation error will be about the same as if we had simply truncated these components.

For large components,

$$\sigma_r^2 \gg \delta \quad \Rightarrow \quad \frac{\delta^2}{(\sigma_r^2 + \delta)^2} \approx \frac{\delta^2}{\sigma_r^2}$$

and so this portion of the approximation error will be very small.

For the noise error energy, we have

$$\begin{aligned} \|\text{Noise error}\|_2^2 &= \sum_{r=1}^R \left( \frac{\sigma_r}{\sigma_r^2 + \delta} \right)^2 |\langle \mathbf{e}, \mathbf{u}_r \rangle|^2 \\ &\leq \frac{1}{4\delta} \sum_{r=1}^R |\langle \mathbf{e}, \mathbf{u}_r \rangle|^2 \\ &\leq \frac{1}{4\delta} \|\mathbf{e}\|_2^2. \end{aligned}$$

The worst-case error is more or less determined by the choice of  $\delta$ . The regularization makes the effective condition number of  $\mathbf{A}$  about  $1/(2\sqrt{\delta})$ ; no matter how small the smallest singular value is, the noise energy will not increase by more than a factor of  $1/(4\delta)$  during the reconstruction process.

As usual, the average case error is less pessimistic. If each entry of  $\mathbf{e}$  is an independent identically distributed Normal random variable with mean zero and variance  $\nu^2$ , then the expected noise error in the reconstruction will be

$$\begin{aligned} \mathbb{E} [\|\text{Noise error}\|_2^2] &= \sum_{r=1}^R \left( \frac{\sigma_r}{\sigma_r^2 + \delta} \right)^2 \mathbb{E} [|\langle \mathbf{e}, \mathbf{u}_r \rangle|^2] \\ &= \nu^2 \cdot \sum_{r=1}^R \left( \frac{\sigma_r}{\sigma_r^2 + \delta} \right)^2 \\ &= \frac{1}{M} \cdot \left( \sum_{r=1}^R \frac{\sigma_r^2}{(\sigma_r^2 + \delta)^2} \right) \cdot \mathbb{E} [\|\mathbf{e}\|_2^2]. \quad (11) \end{aligned}$$

Note that  $\frac{\sigma_r^2}{(\sigma_r^2 + \delta)^2} \leq \min\left(\frac{1}{\sigma_r^2}, \frac{1}{4\delta}\right)$ , so we can think of the error in (11) as an average of the  $\frac{1}{\sigma_r^2}$ , with the large values simply replaced by  $1/(4\delta)$ .



## A Closed Form Expression

Tikhonov regularization is in some sense very similar to the truncated SVD, but with one significant advantage: we do not need to explicitly calculate the SVD to solve (6). Indeed, the solution to (6) can be written as

$$\hat{\mathbf{x}}_{\text{tik}} = (\mathbf{A}^T \mathbf{A} + \delta \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y}. \quad (12)$$

To see that the expression above is equivalent to (9), note that we can write  $\mathbf{A}^T \mathbf{A}$  as

$$\mathbf{A}^T \mathbf{A} = \mathbf{V} \boldsymbol{\Sigma}^2 \mathbf{V}^T = \mathbf{V}' \boldsymbol{\Sigma}'^2 \mathbf{V}'^T,$$

where  $\mathbf{V}'$  is  $N \times N$ ,

$$\mathbf{V}' = [\mathbf{V} \ \mathbf{V}_0],$$

and the  $N \times N$  diagonal matrix  $\boldsymbol{\Sigma}'$  is simply  $\boldsymbol{\Sigma}$  padded with zeros:

$$\boldsymbol{\Sigma}' = \begin{bmatrix} \boldsymbol{\Sigma} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}.$$

The verification of (12) is now straightforward:

$$\begin{aligned} (\mathbf{A}^T \mathbf{A} + \delta \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y} &= (\mathbf{V}' \boldsymbol{\Sigma}'^2 \mathbf{V}'^T + \delta \mathbf{I})^{-1} \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^T \mathbf{y} \\ &= \mathbf{V}' (\boldsymbol{\Sigma}'^2 + \delta \mathbf{I})^{-1} \mathbf{V}'^T \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^T \mathbf{y} \\ &= \mathbf{V}' (\boldsymbol{\Sigma}'^2 + \delta \mathbf{I})^{-1} \begin{bmatrix} \boldsymbol{\Sigma} \mathbf{U}^T \mathbf{y} \\ \mathbf{0} \end{bmatrix} \\ &= \mathbf{V} (\boldsymbol{\Sigma}^2 + \delta \mathbf{I})^{-1} \boldsymbol{\Sigma} \mathbf{U}^T \mathbf{y} \\ &= \hat{\mathbf{x}}_{\text{tik}}. \end{aligned}$$

The expression (12) holds for all  $M$ ,  $N$ , and  $R$ . We will leave it as an exercise to show that

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

The importance of not needing to explicitly compute the SVD is significant when we are solving large problems. When  $\mathbf{A}$  is large ( $M, N > 10^5$ , say) it may be expensive or even impossible to construct the SVD and compute with it explicitly. However, if it has special structure (if it is sparse, for example), then it may take many fewer than  $MN$  operations to compute a matrix vector product  $\mathbf{Ax}$ .

In these situations, a **matrix free** iterative algorithm can be used to perform the inverse required in (12). A prominent example of such an algorithm is **conjugate gradients**, which we will see later in this course.