

For  $A$  with SVD

$$A = U \Sigma V^T$$

we call

$$A^+ = V \Sigma^{-1} U^T$$

The pseudo-inverse of  $A$

Other names  $A^+$  goes by:

- "Natural inverse"
- "Lanczos inverse"
- "Moore-Penrose inverse"

---

Least-squares example: simple extrapolation

$$\begin{bmatrix} y[0] \\ y[1] \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} x[0] \\ x[1] \\ x[2] \\ x[3] \end{bmatrix}$$

We see the first two samples, but not the second two.

It is clear that  $x[2]$  &  $x[3]$  can be anything, and our observations  $y[0], y[1]$  will not change.

What does the least squares solution

$$\hat{x} = V \Sigma^{-1} U^T y + V_0 \kappa_0$$

look like?

Note that

$$A^T A = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

so we can take

$$V = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \quad V_0 = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}$$

(these choices are not unique, though)

$$A A^T = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \text{so } U = [u_1 \ u_2] = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

and there is no  $V_0$  (full row rank)  
( $\Rightarrow \beta_0 = 0$ )

$$\begin{aligned} \hat{x} &= V \Sigma^{-1} U^T y + V_0 \kappa_0 \\ &= \begin{bmatrix} y[0] \\ y[1] \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ \kappa_0[2] \\ \kappa_0[3] \end{bmatrix} \end{aligned}$$

for any  $\kappa_0[2], \kappa_0[3]$

The min-norm solution has  $\kappa_0[2] = \kappa_0[3] = 0$   
(usual case)

## Pseudo-inverse summary

Given  $y$ , we set

$$\hat{x} = A^+ y = V \Sigma^{-1} U^T y$$

We know that

- there is no  $x$  whose residual  $y - Ax$  is smaller than  $y - A\hat{x}$

In other words, there is no  $x$  such that  $Ax$  is closer to  $y$  than  $AA^+y$ .

- Of all the  $x$  with minimum residual, none has smaller norm than  $A^+y$ .

Also note that the pseudo-inverse always exists (every matrix has an SVD).

### Explicit formulas for $A^+$ when $A$ is full rank:

• If  $A$  is square & invertible, then

$$A^+ = A^{-1}$$

Check:

$$A = U \Sigma V^T$$

$U$  is  $N \times N$ ,

$V$  is  $N \times N$ ,

$$U^T U = U U^T = I$$

$$V^T V = V V^T = I$$

$$A^+ A = V \Sigma^{-1} \underbrace{U^T U}_{=I} \Sigma V^T = V \underbrace{\Sigma^{-1} \Sigma}_{=I} V^T = V V^T = I$$

- If  $A^T A$  is invertible (full column rank), then

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

This occurs when  $A$  is "tall & skinny" ( $M > N$ )

$$\begin{bmatrix} A \end{bmatrix}$$

and all the columns are linearly independent.

We will have  $p = N$  non-zero singular values

Check:

$$A = U \Sigma V^T$$

Since  $p = N$ ,  $V$  is  $N \times N$   $\perp$  matrix  
 $V V^T = V^T V = I$

$U$  is  $M \times N$   $\perp$  matrix  
 $U^T U = I$

$$A^T A = V \Sigma U^T U \Sigma V^T = V \Sigma^2 V^T$$

$$(A^T A)^{-1} A^T = V \Sigma^{-2} V^T V \Sigma U^T = V \Sigma^{-1} U^T = A^+ \quad \checkmark$$

- If  $AA^T$  is invertible, then (full row rank)  

$$A^+ = A^T (AA^T)^{-1}$$

This occurs when  $A$  is "short & fat"

$$\begin{bmatrix} A \end{bmatrix}$$

and all the rows are linearly independent.

We will have  $p=M$  non-zero singular values.

Check:

$$A = U \Sigma V^T$$

$U$  is  $M \times M, \perp$

$V$  is  $N \times M, \perp$ ,

$$U U^T = U^T U = I$$

$$V^T V = I$$

$$AA^T = U \Sigma V^T V \Sigma U^T$$

$$= U \Sigma^2 U^T$$

$$A^T (AA^T)^{-1} = V \Sigma U^T U \Sigma^{-2} U^T = V \Sigma^{-1} U^T = A^+$$

$A^+$  is as close to an inverse of  $A$  as possible

Left inverse:

Given  $y = Ax$ , we'd like  $A^+y = A^+Ax \approx x$   
for any  $x$ . That is, we'd like  $A^+A = I$ .

Of course this is not always possible.

But is there an  $H$  such that  $HA$   
is closer to the identity than  $A^+A$ ?

The answer is no. To see this, we

consider the problem  
(\*)  $\min_H \|HA - I\|_F^2$  ( $\|\cdot\|_F$  = "Frobenius norm")

(Here the  $F$ -norm<sup>2</sup> of a matrix  $Q$

$$\|Q\|_F^2 = \sum_{k,e} Q_{k,e}^2$$

= sum of the squares of the entries)

The solution to this is the  $H$  which make  
 $HA$  closest to the identity.

It is a fact (which comes from setting derivatives = 0) that the solution to (\*) will obey

$$AA^T H^T = A$$

So, does  $H = A^+ = V\Sigma^{-1}U^T$  work?

$$AA^T H^T = U \underbrace{\Sigma V^T V \Sigma^{-1}}_{\mathbf{I}} \underbrace{U^T U}_{\mathbf{I}} \Sigma^{-1} V^T = U \Sigma V^T = A$$

yes.

Right Inverse:

Our solution  $\hat{x}$  should come as close to explaining the data as possible. We'd like

$$A \hat{x} = y$$

That is,

$$AA^+ y = y \quad (\text{i.e. } AA^+ = \mathbf{I})$$

Of course, this might not be possible for all  $y$ .  
But is there an  $H$  such that  $AH$  is closer to the identity than  $AA^+$ ?

Again, we look at

$$\min_H \|AH - I\|_2^2$$

After a little bit of work, we have that the solution to the above will obey

$$HH^T A^T = H$$

Does  $H = V \Sigma^{-1} U^T$  work?

$$HH^T A^T = V \underbrace{\Sigma^{-1} U^T U}_{I} \underbrace{\Sigma^{-1} V^T V}_{I} \Sigma U^T = V \Sigma^{-1} U^T = H$$

so yes.

Moral:  $A^+ = V \Sigma^{-1} U^T$  is as close to an inverse of  $A$  as you can possibly have.



## Stability

What happens when there is noise?

$$y = Ax + e \quad A = U\Sigma V^T$$

Say we apply  $A^+$  to  $y$  above to get  $\hat{x}$ :

$$\hat{x} = A^+ y = \underbrace{A^+ A x}_{= VV^T x} + A^+ e = \tilde{x} = \text{"optimal estimate" when there is no noise.}$$

How far is  $\hat{x}$  from  $\tilde{x}$ ?

(What is the reconstruction error  $\|\hat{x} - \tilde{x}\|_2^2$ ?)

$$\text{It is } \|\hat{x} - \tilde{x}\|_2^2 = \|A^+ e\|_2^2$$

Suppose for a minute that  $\|e\|_2^2 = 1$   
What is the "worst case" for  $\|A^+ e\|_2^2$

$$\max_{\|e\|_2^2=1} \|A^+ e\|_2^2 = \max_{\|e\|_2^2=1} \|V\Sigma^{-1}U^T e\|_2^2$$

Since  $\|U^T e\|_2^2 \leq \|e\|_2^2$  (the columns of  $U$  are  $\perp$ )

the above is equivalent to

$$\max_{\|\beta\|_2^2=1} \|V\Sigma^{-1}\beta\|_2^2$$

Also, for any vector  $z$

$$\|Vz\|_2^2 = \langle Vz, Vz \rangle = \langle z, V^T V z \rangle = \langle z, z \rangle = \|z\|_2^2$$

(the columns of  $V$  are  $\perp$ )

So we can further simplify to

$$\max_{\|\beta\|_2^2=1} \|\Sigma^{-1}\beta\|_2^2$$

The worst case  $\beta$  will have a 1 in the entry corresp. to the largest value on the diagonal of  $\Sigma^{-1}$ , and zero everywhere else.

Thus

$$\max_{\|\beta\|_2^2=1} \|\Sigma^{-1}\beta\|_2^2 = \max_{k=1,\dots,p} \frac{1}{\sigma_k^2} = \frac{1}{\sigma_p^2}$$

So

$$\|\hat{x} - \tilde{x}\|_2^2 = \|V\Sigma^{-1}U^T e\|_2^2 \leq \frac{1}{\sigma_p^2} \cdot \|e\|_2^2$$

Similarly,  $\|\hat{x} - \tilde{x}\|_2^2 \geq \frac{1}{\sigma_1^2} \cdot \|e\|_2^2$ , so

$$\frac{1}{\sigma_1^2} \cdot \|e\|_2^2 \leq \|\hat{x} - \tilde{x}\|_2^2 \leq \frac{1}{\sigma_p^2} \cdot \|e\|_2^2$$

If  $\sigma_p$  is small, the worst case error can be very bad.

"Average case" error can also be related to the singular values.

Say  $e$  is iid Gaussian noise (AWGN)

$$e_k \sim \text{Normal}(0, \sigma^2)$$

The average measurement error is

$$E\|e\|_2^2 = M\sigma^2$$

The average reconstruction error is

$$E\|A^+e\|_2^2 = \frac{1}{M} \cdot \left( \frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2} + \dots + \frac{1}{\sigma_p^2} \right) \cdot M\sigma^2$$

error magnification  $\sim$  average of the singular values

If  $\sigma_p$  is tiny,  $\frac{1}{\sigma_p^2}$  will dominate the rest of the terms.

Moral: Applying the pseudo-inverse  $A^+$  can be disastrous in the presence of noise.

One way to protect against this is to truncate the SVD.

We have

$$A = U \Sigma V^T = \sum_{k=1}^p \sigma_k u_k v_k^T$$

Set a threshold  $\sigma_{\min}$ . Find  $p' \leq p$

so that  $\sigma_{p'} \geq \sigma_{\min}$ ,  $\sigma_{p'+1} < \sigma_{\min}$

i.e.  $\sigma_1, \sigma_2, \dots, \sigma_{p'}$  are above  $\sigma_{\min}$   
 $\sigma_{p'+1}, \dots, \sigma_p$  are below  $\sigma_{\min}$

The truncated sum

$$A' = \sum_{k=1}^{p'} \sigma_k u_k v_k^T$$

is a good approximation to  $A$  (especially if  $\sigma_{\min}$  is small compared to  $\sigma_1$ ).

Then, instead of applying the pseudo inverse of  $A$ , we use  $A'^{\dagger}$ .

$$\hat{x} = A'^{\dagger} y = A'^{\dagger} A x + A'^{\dagger} e$$

The reconstruction error

$$\hat{X} - \tilde{X} = A'^+ A x - A^+ A x + A'^+ e$$

$$= \underbrace{(A'^+ - A^+) A x}_{\text{truncation error}} + \underbrace{A'^+ e}_{\text{noise error}}$$

now has two parts

To see what the truncation error is, write

$$A'^+ - A^+ = \sum_{k=p'+1}^p \frac{1}{\sigma_k} V_k U_k^T$$

then

$$(A'^+ - A^+) A x = \left( \sum_{k=p'+1}^p \frac{1}{\sigma_k} V_k U_k^T \right) \cdot \left( \sum_{\ell=1}^p \sigma_\ell U_\ell V_\ell^T x \right)$$

$$= \sum_{k=p'+1}^p V_k \cdot \langle V_k, x \rangle$$

= portion of  $x$  in the space  
we are cutting out

$$\| (A'^+ - A^+) A x \|_2^2 =$$

Also,

$$A^{++}e = \sum_{k=1}^{p'} \frac{1}{\sigma_k} \langle u_k, e \rangle v_k$$

and

$$\|A^{++}e\|_2^2 = \sum_{k=1}^{p'} \frac{1}{\sigma_k^2} |\langle u_k, e \rangle|^2$$

Note that

$$\underbrace{(A^{++} - A^+)Ax}_{\text{truncation error}} \perp \underbrace{A^{++}e}_{\text{noise error}}$$

Since  $A^{++}e \in \text{span}\{v_1, \dots, v_{p'}\}$

$(A^{++} - A^+)Ax \in \text{span}\{v_{p'+1}, \dots, v_p\}$

So the total error is

$$\|\hat{x} - \tilde{x}\|_2^2 = \sum_{k=p'+1}^p |\langle v_k, x \rangle|^2 + \sum_{k=1}^{p'} \frac{1}{\sigma_k^2} |\langle u_k, e \rangle|^2$$

## Tikhonov Regularization

There is another way (other than truncation) to get stable solutions to least squares problems.

Recall that the pseudo-inverse will give us a solution that solves

$$(*) \quad \min_x \|y - Ax\|_2^2$$

If  $A$  has singular values which are very small,  $(*)$  could make  $x$  huge (in the direction of the  $v_k$  corresponding to tiny  $\sigma_k$ ) in its effort to match  $y$  as closely as possible.

We can modify  $(*)$  by adding a regularization term that penalizes  $x$  with large norm.

Instead of  $(*)$ , we solve

$$(**) \quad \min_x \|y - Ax\|_2^2 + \delta \cdot \|x\|_2^2$$

for some small  $\delta > 0$ .

← keeps  $\|x\|_2^2$  from getting too big

It is not terribly hard to show (on homework) that  $(**)$  has a unique solution given by

$$\hat{x}_{\text{Tik}} = (A^T A + \delta I)^{-1} A^T y$$

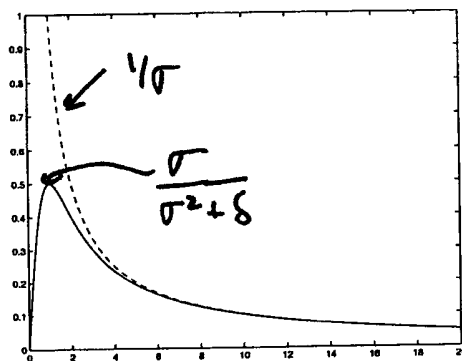
For small  $\delta$ ,  $A^T A + \delta I$  will be close to  $A^T A$ , but it is guaranteed to be invertible.

In terms of the SVD of  $A$ , we can write

$$\hat{x}_{\text{Tik}} = \sum_{k=1}^p \frac{\sigma_k}{\sigma_k^2 + \delta} \cdot \langle u_k, y \rangle \cdot v_k$$

↳ This looks like  $1/\sigma_k$  for  $\sigma_k \gg \delta$  and goes to 0 (dampens) for  $\sigma_k \ll \delta$ .

$$\delta = 1$$



$\sigma \rightarrow$



## Summary:

Pseudo-Inverse (least-squares):

$$\bullet \hat{X}_{LS} = V \Sigma^{-1} U^T y = \sum_{k=1}^p \frac{1}{\sigma_k} \langle u_k, y \rangle v_k$$

- If  $A$  has full rank, we can compute without SVD using

$$\hat{X}_{LS} = (A^T A)^{-1} A^T y \quad (M \geq N)$$

or

$$\hat{X}_{LS} = A^T (A A^T)^{-1} y \quad (M \leq N)$$

- The solution can "blow up" (and probably will) if  $\sigma_p$  is very small

Truncated Pseudo-Inverse

$$\bullet \hat{X}_{\text{Trunc}} = \sum_{k=1}^{p'} \frac{1}{\sigma_k} \langle u_k, y \rangle v_k$$

- simply throw away the terms in the sum for  $\hat{X}_{LS}$  that correspond to small singular values

- need the SVD to compute this

- stable if  $p'$  is chosen appropriately

## Tikhonov Regularized Least-Squares

$$\begin{aligned} \hat{X}_{\text{Tik}} &= V \Sigma (\Sigma^2 + \delta I)^{-1} U^T y \\ &= \sum_{k=1}^p \frac{\sigma_k}{\sigma_k^2 + \delta} \langle u_k, y \rangle v_k \end{aligned}$$

- gradually "dampens" the components corresponding to very small singular values
- also has the advantage that we don't need the SVD to compute it:

$$\hat{X}_{\text{Tik}} = (A^T A + \delta I)^{-1} A^T y$$

## Extensions

There are natural generalizations to both Tikhonov regularization and the truncated SVD.

For Tikhonov, in adding the regularization term  $\delta \|x\|_2^2$  we are implicitly assuming some type of structure on  $x$  — that is, that its energy is not too large. (\*\*) on page 39 favors signals with smaller energy.

We can encourage different types of structure as well. For any matrix  $D$  w/  $N$  columns, we can solve

$$\min_x \|y - Ax\|_2^2 + \delta \|Dx\|_2^2$$

this encourages  $Dx$  to have energy which is not too large.

A typical example is if  $D$  is an approximation to a derivative.

i.e.  $D = \begin{bmatrix} 1 & & & & \\ & -1 & & & \\ & & 1 & & \\ & & & -1 & \\ & & & & 1 \\ & & & & & -1 \\ & & & & & & \dots \end{bmatrix}$

Roughly speaking, this keeps the solution from "wiggling around too much". Another model which has a similar effect is

$$D = WF,$$

where  $F$  is a Fourier transform matrix and  $W$  is a diagonal matrix with positive entries which are larger in places (i.e. rows of  $F$ ) corresponding to high frequencies and smaller at indices corresponding to low frequencies.

This penalizes high-frequency components of the solution more than low-frequency components.

The program

$$\min_x \|y - Ax\|_2^2 + \delta \|Dx\|_2^2$$

has closed-form solution

$$\hat{x} = (A^T A + \delta D^T D)^{-1} A^T y$$

when  $A^T A + \delta D^T D$  is invertible.

---

### Least-Squares with Linear Constraints

Another way to incorporate a model for  $x$  is to force the solution to lie in a particular subspace of  $\mathbb{R}^N$ .

If this  $d$ -dimensional subspace is spanned by the columns of the  $N \times d$  matrix  $G$ , we solve

$$\begin{aligned} & \min_{x \in \mathbb{R}^N} \|y - Ax\|_2^2 \\ & \text{subject to } x = G\alpha \text{ for some } \alpha \in \mathbb{R}^d \\ & = \min_{\alpha \in \mathbb{R}^d} \|y - A G \alpha\|_2^2 \end{aligned}$$

If  $AG$  has full column rank, then the solution is

$$\hat{\alpha} = (G^T A^T A G)^{-1} G^T A^T y$$

and the corresponding estimate of  $x$  is

$$\hat{x} = G\hat{\alpha} = G(G^T A^T A G)^{-1} G^T A^T y.$$

(If  $AG$  is not full column-rank, then we can take  $\hat{\alpha} = (AG)^+ y$ .)

The basic idea is that while the  $N \times N$  matrix  $A^T A$  may not be well conditioned, the  $d \times d$  matrix  $G^T A^T A G$  may be — of course, establishing this requires a careful study of the relationship between  $G$  &  $A$ .

The truncated SVD is actually an example of constrained least-squares. There we implicitly took

$$G = \left[ v_1 \mid v_2 \mid \dots \mid v_p \right]$$

where the  $\{v_i\}$  are the "right singular vectors" of  $A$ .

## Example: circular deconvolution

Let  $h \in \mathbb{R}^{1024}$  be defined by

$$h_n = \begin{cases} 1 & 0 \leq n \leq 31 \\ 0 & 32 \leq n \leq 1023 \end{cases},$$

and let  $H$  be the corresponding matrix generated by  $H$ :

$$H = \begin{bmatrix} h_0 & h_{N-1} & h_{N-2} & \cdots & h_1 \\ h_1 & h_0 & h_{N-1} & \cdots & h_2 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ h_{N-1} & h_{N-2} & h_{N-3} & \cdots & h_0 \end{bmatrix}.$$

(Here  $N = 1024$ .) As we saw in the previous homework, we can write

$$H = FDF^H = \sum_{k=0}^{N-1} d_k f_k f_k^H$$

where the matrix  $F$  has columns

$$f_k = \frac{1}{\sqrt{N}} \begin{bmatrix} 1 \\ e^{j2\pi k/N} \\ e^{j2\pi 2k/N} \\ \vdots \\ e^{j2\pi(N-1)k/N} \end{bmatrix}, \quad k = 0, \dots, N-1$$

and the  $d_k$  are the DFT coefficients of  $h$ . This is not quite an SVD of  $H$  (since  $H$  is real-valued, and both  $F$  and  $D$  are complex-valued), but it has the same essential properties — just be careful below about conjugates and the difference between  $z^2$  and  $|z|^2$ .

1. In MATLAB, construct  $H$ . If you like, you can use the diagonalization formula above to do this; it is easy to construct  $F$  using `F = 1/sqrt(N)*fft(eye(N));`. Turn in your code and a plot of `imagesc(H)`.

2. Suppose we observe

$$y = Hx + e,$$

and attempt to recover  $x$  using Tikhonov regularization

$$\min_x \|y - Hx\|_2^2 + \delta \|x\|_2^2.$$

Show that we can write  $\hat{x} = Gy$ , where  $G$  is also a circulant matrix. What are the Fourier coefficients for the  $g \in \mathbb{R}^N$  whose circular shifts generate  $G$ ?

3. The file **hw7prob5.mat** contains a vector  $x$  of length 1024 and a noisy observation through  $H$ ,  $y = Hx + e$ , where  $e_n \sim \text{Normal}(0, 1)$ . Estimate  $x$  using Tikhonov regularization for  $\delta = 1e - 4, 1e - 2, 1$  and  $5$ . For each value of  $\delta$ , plot  $x$  and  $\hat{x}$  overlaid on one another. Comment on how the estimate changes as  $\delta$  gets larger. Turn in your plots and your code.