

**Approximating Matrices:
One Picture and One Thousand Words**
Gene Golub Memorial Day, February 2017

Dianne P. O'Leary
©2017

In Honor of Gene H. Golub
Approximating Matrices:
One Picture and One Thousand Words

Dianne P. O'Leary

Computer Science Dept. and
Institute for Advanced Computer Studies
University of Maryland



Joint work with
Julianne Chung,
Matthias Chung,
John M. Conroy,
Yi-Kai Liu

Support from NSF.

Introduction: The Problem

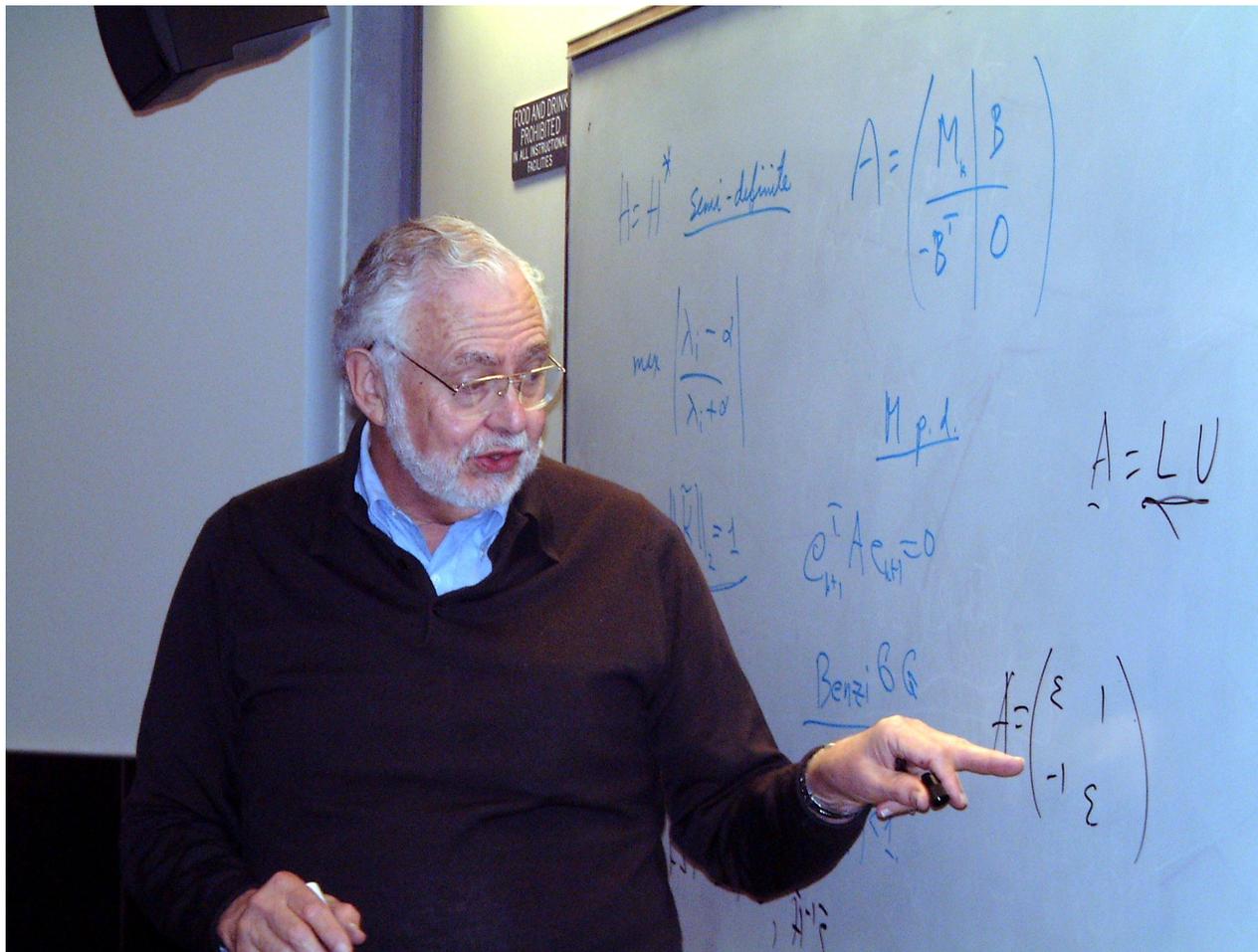
Given: a rank r matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, with $m \geq n$,

Determine: an **approximation** \mathbf{Z} of rank $k \leq r$.

Low rank approximations of matrices and their inverses play an important role in many applications, such as **image reconstruction**, machine learning, **text processing**, matrix completion problems, **signal processing**, optimal control problems, statistics, and mathematical biology.

This work was inspired by Gene Golub

An inspired and inspiring teacher,



A dedicated mentor,



A generous friend to many of us here,



An outstanding researcher and catalyst.



Major contributions to

- numerical solution of eigenvalue problems,
- solution of least squares and total least squares problems,
- iterative solution of linear systems,
- understanding of orthogonal polynomials and quadrature,
- computation and applications of matrix factorizations.

Low-rank approximation was a focus of Gene's work for over 40 years

From his post-doctoral work, partially motivated by the use of the SVD in low-rank matrix approximation...

Calculating the Singular Values and Pseudo-Inverse of a Matrix
Gene Golub and W Kahan, SIAM Journal 1965

... To one of his last publications:

Rank-one Approximation to High Order Tensors
T Zhang, GH Golub, SIMAX 2001

With many important contributions in between, such as:

- [Tracking a Few Extreme Singular Values and Vectors in Signal Processing, P Comon, GH Golub - Proceedings of the IEEE, 1990](#)
Maintaining a low-rank approximation to a covariance matrix slowly changing over time, based on Gene's updating methods for matrix factorizations.
- [Fast Algorithms for Updating Signal Subspaces, G Xu, H Zha, G Golub, T Kailath - IEEE Transactions on Circuits, 1994](#)
Tracking a low-dimensional subspace slowly changing over time using the Lanczos algorithm.
- [The Restricted Singular Value Decomposition: Properties and Applications, Bart LR De Moor and Gene H Golub, SIMAX 2006](#)
Low-rank approximations to partitioned matrices.

Gene emphasized two methods for computing low-rank approximations

- The Singular Value Decomposition (SVD)
- The Conjugate Gradient / Lanczos algorithm



Low-rank approximations from the SVD: best in the Frobenius norm

$$\min_{\text{rank}(\mathbf{Z}) \leq k} f(\mathbf{A}, \mathbf{Z}), \quad f(\mathbf{A}, \mathbf{Z}) = \|\mathbf{A} - \mathbf{Z}\|_F^2$$

Solution: best rank k approximation of an $m \times n$ matrix \mathbf{A} Let

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$$

be the singular value decomposition of \mathbf{A} , with $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_m]$ and $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_n]$ and $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$.

The **Eckart-Young Theorem**, Schmidt 1907, Eckart and Young 1936, Mirsky 1960 (unitarily invariant norms)

$$\hat{\mathbf{Z}} = \mathbf{U}_k \mathbf{\Sigma}_k \mathbf{V}_k^\top$$

with

$$\mathbf{U}_k = [\mathbf{u}_1, \dots, \mathbf{u}_k], \quad \mathbf{V}_k = [\mathbf{v}_1, \dots, \mathbf{v}_k], \\ \mathbf{\Sigma}_k = \text{diag}(\sigma_1, \dots, \sigma_k).$$

The solution is unique if and only if $\sigma_k > \sigma_{k+1}$.

Low-rank approximations from conjugate gradients

In the 1975 Concus-Golub-O'Leary paper on the conjugate gradient algorithm, Gene made sure that the matrix approximation viewpoint was included: If an $n \times n$ symmetric positive definite matrix \mathbf{A} could be expressed as

$$\mathbf{A} = \mathbf{M} - \mathbf{N},$$

where \mathbf{M} is a preconditioner (close to \mathbf{A} in some sense and easily inverted), then the matrix

$$\mathbf{K} = \mathbf{I} - \mathbf{M}^{-1}\mathbf{N} = \bar{\mathbf{Z}}_n \mathbf{J}_n \bar{\mathbf{Z}}_n^{-1},$$

where

- the columns of $\bar{\mathbf{Z}}_n$ are the preconditioned residuals and are \mathbf{M} -orthogonal,
- \mathbf{J}_n is a tridiagonal matrix containing the cg parameters.

Therefore, after $k < n$ steps, we obtain a matrix approximation

$$\mathbf{K} \approx \bar{\mathbf{Z}}_k \mathbf{J}_k \bar{\mathbf{Z}}_k^\dagger.$$

Is that the end of the story, or just the beginning?

In approximating \mathbf{A} by \mathbf{Z} , how we

- choose $f(\mathbf{A}, \mathbf{Z})$ to measure the **goodness** of the approximation, and
- what additional constraints we put on \mathbf{Z}

determine how hard our problem is.

How hard can it be?

How hard can it be?

Today we'll consider two other ways to form low-rank approximations:

- Inverse approximation.
Joint work with Julianne Chung and Matthias Chung.
- Interval approximation.
Joint work with John M. Conroy and Yi-Kai Liu.

And we'll discuss why these viewpoints are necessary and useful.

Just as in the Eckart-Young theorem, our approximations take the form

$$\mathbf{Z} = \mathbf{X}\mathbf{Y}, \quad \mathbf{X} : m \times k, \mathbf{Y} : k \times n.$$

First we'll give two examples of why we might want constraints on \mathbf{X} and \mathbf{Y} , not just on \mathbf{Z} .

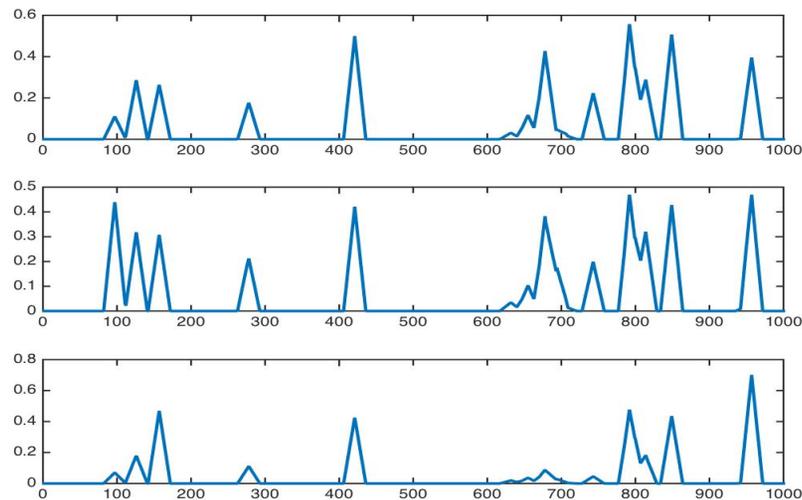
This seems odd.

Why put constraints on \mathbf{X} and \mathbf{Y} when it is $\mathbf{Z} = \mathbf{X}\mathbf{Y}$ that is supposed to approximate \mathbf{A} ?

An application: spectroscopy

Suppose we are given 50 samples, each made up of various proportions of 5 unknown substances that don't react with each other.

And suppose we use spectroscopy (IR, for example) to measure a (noisy) spectrum ($m = 1000$ numbers) for each sample. Here are 3 of the 50 (simulated data):



These are 3 columns of an $m \times 50$ data matrix \mathbf{A} .

Each spectrum measures a response at 1000 frequencies, and they are nonnegative.

We've been told that each spectrum \mathbf{a}_i is of the form

$$\mathbf{a}_i = y_{1i}\mathbf{x}_1 + y_{2i}\mathbf{x}_2 + y_{3i}\mathbf{x}_3 + y_{4i}\mathbf{x}_4 + y_{5i}\mathbf{x}_5, \quad i = 1, \dots, n$$

\mathbf{x} vectors: the spectra for the five unknown substances

y coefficients: how much of each substance is in the i th sample.

So we want a good approximation $\mathbf{Z} = \mathbf{X}\mathbf{Y}$ to the data matrix \mathbf{A} , but to make it physically meaningful, we need both \mathbf{X} and \mathbf{Y} to be nonnegative.

Also, the vectors \mathbf{X} should be sparse, because the spectrum for each substance has a small number of **peaks** (nonzero responses).

Another application: document classification

Suppose we are given 500 documents (e.g., articles from the math literature) and we want to classify them into k subjects.

Let's gather the distinct **terms** in the documents ('decomposition', 'topology', 'category theory', 'derivative', etc.) and create a matrix \mathbf{A} that has one row for each term.

Set a_{ij} to be a measure of the **importance** of term i in document j , for $i = 1, \dots, m$ = the number of terms and $j = 1, \dots, 500$.

If we can determine a nonnegative $m \times k$ matrix \mathbf{X} and a nonnegative $k \times 500$ matrix \mathbf{Y} with $\mathbf{A} \approx \mathbf{X}\mathbf{Y}$, then each column of \mathbf{A} is approximately a combination of the columns of \mathbf{X} . The coefficients in the i th column of \mathbf{Y} tell us how important each of these **dictionary** columns is to document i .

We would expect both \mathbf{X} and \mathbf{Y} to be nonnegative.

And we would also expect sparsity in each of these matrices.

Do we want to approximate a matrix, or its inverse?

A different measure of goodness: Inverse approximation

Instead of

$$\min_{\text{rank}(\mathbf{Z}) \leq k} f(\mathbf{A}, \mathbf{Z}), \quad f(\mathbf{A}, \mathbf{Z}) = \|\mathbf{A} - \mathbf{Z}\|_F^2$$

let's consider

$$\min_{\text{rank}(\mathbf{Z}) \leq k} f(\mathbf{A}, \mathbf{Z}) \quad f(\mathbf{A}, \mathbf{Z}) = \|\mathbf{Z}\mathbf{A} - \mathbf{I}_n\|_F^2$$

A solution:

$$\hat{\mathbf{Z}} = \mathbf{V}_k \Sigma_k^{-1} \mathbf{U}_k^\top.$$

But **any** choice of k nonzero singular values and the corresponding vectors gives a global minimizer, so the problem is not well-posed.

Adding constraints can make the problem easier.

For example, if we take our inverse approximation problem

$$\min_{\text{rank}(\mathbf{Z}) \leq k} f(\mathbf{A}, \mathbf{Z}) \quad \text{where} \quad f(\mathbf{A}, \mathbf{Z}) = \|\mathbf{Z}\mathbf{A} - \mathbf{I}_n\|_F^2,$$

and add the constraint

$$\|\mathbf{Z}\|_F \leq c$$

for a small enough constant c , then we will see later that we have made the solution unique.

Adding constraints can make the problem harder.

Constraints on the **factors** $Z = XY$, where X is $m \times k$ and Y is $k \times n$ are useful but complicate the problem.

We could ask that

- X and Y be nonnegative, or
- X and Y be sparse.

Inverse Approximation: Finding a solution

$$\min_{\text{rank}(\mathbf{Z}) \leq k} \|(\mathbf{Z}\mathbf{A} - \mathbf{I}_n)\|_{\text{F}}^2 + \alpha^2 \|\mathbf{Z}\|_{\text{F}}^2,$$

- \mathbf{A} has dimension $m \times n$, with $m \geq n$.
- $k \leq \text{rank}(\mathbf{A})$ is a given positive integer.
- α is a given parameter, nonzero if $\text{rank}(\mathbf{A}) < m$.
- $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^{\top}$.

$$\min_{\text{rank}(\mathbf{Z}) \leq k} \|(\mathbf{Z}\mathbf{A} - \mathbf{I}_n)\|_{\text{F}}^2 + \alpha^2 \|\mathbf{Z}\|_{\text{F}}^2,$$

Theorem: (Chung, Chung, O'Leary) A global minimizer $\hat{\mathbf{Z}} \in \mathbb{R}^{n \times m}$ is

$$\hat{\mathbf{Z}} = \mathbf{V}_k \mathbf{\Psi}_k \mathbf{U}_k^{\top},$$

where \mathbf{V}_k contains the first k columns of \mathbf{V} , \mathbf{U}_k contains the first k columns of \mathbf{U} , and

$$\mathbf{\Psi}_k = \text{diag}\left(\frac{\sigma_1}{\sigma_1^2 + \alpha^2}, \dots, \frac{\sigma_k}{\sigma_k^2 + \alpha^2}\right).$$

Moreover, if $\alpha \neq 0$, this $\hat{\mathbf{Z}}$ is the *unique* global minimizer if and only if $\sigma_k > \sigma_{k+1}$

What if $\alpha = 0$?

$$\min_{\text{rank}(\mathbf{Z}) \leq k} \|(\mathbf{Z}\mathbf{A} - \mathbf{I}_n)\|_{\text{F}}^2 + \alpha^2 \|\mathbf{Z}\|_{\text{F}}^2,$$

A solution still exists, but, as noted by Friedland and Torokhti (2007), it is not unique if $k < \text{rank}(\mathbf{A})$. In fact,

$$\|\mathbf{Z}\mathbf{A} - \mathbf{I}_n\|_{\text{F}}^2 \geq n - k.$$

We can achieve this lower bound by choosing $\mathbf{Z} = \mathbf{V}_k \Sigma_k^{-1} \mathbf{U}_k^{\top}$, or a matrix of this form constructed using *any* choice of k singular values and corresponding singular vectors.

These are the only global minimizers.

What if we use a different norm?

Consider, for example,

$$\min_{\text{rank}(\mathbf{Z}) \leq k} \|(\mathbf{Z}\mathbf{A} - \mathbf{I}_n)\|_2^2 + \alpha^2 \|\mathbf{Z}\|_2^2,$$

where $\|\mathbf{F}\|_2^2$ is the largest eigenvalue of $\mathbf{F}^\top \mathbf{F}$.

- For any \mathbf{Z} with $\text{rank } k < n$, $\|\mathbf{Z}\mathbf{A} - \mathbf{I}_n\|_2^2 \geq 1$.
- Therefore, $\mathbf{Z} = \mathbf{0}$ is a global minimizer, unique if $\alpha \neq 0$, so the problem is not interesting.

How does Z relate to well-known approximate inverses $V\Phi\Sigma^\dagger U^\top$?

Φ is a diagonal matrix of **filter factors** ϕ_j .

- Tikhonov regularization:

$$\phi_j = \frac{\sigma_j^2}{\sigma_j^2 + \alpha^2}, \quad j = 1, \dots, n.$$

- Truncated SVD (TSVD) regularization:

$$\phi_j = \begin{cases} 1, & \text{for } j \leq k, \\ 0, & \text{for } j > k. \end{cases}$$

How does Z relate to well-known approximate inverses $V\Phi\Sigma^\dagger U^\top$?

Φ is a diagonal matrix of **filter factors** ϕ_j .

- **Tikhonov** regularization:

$$\phi_j = \frac{\sigma_j^2}{\sigma_j^2 + \alpha^2}, \quad j = 1, \dots, n.$$

- **Truncated SVD (TSVD)** regularization:

$$\phi_j = \begin{cases} 1, & \text{for } j \leq k, \\ 0, & \text{for } j > k. \end{cases}$$

- Our approximate inverse: **Truncated Tikhonov regularization**

$$\phi_j = \begin{cases} \frac{\sigma_j^2}{\sigma_j^2 + \alpha^2}, & \text{for } j \leq k, \\ 0, & \text{for } j > k. \end{cases}$$

How the problem arises in minimizing Bayes risk

Given

- an image ξ , drawn according to a probability distribution of images,
- an operator $A \in \mathbb{R}^{m \times n}$ that distorts the observed image,
- noise δ , drawn according to a probability distribution of noise samples,
- an observation

$$b = A\xi + \delta.$$

It would be nice to have a matrix Z that minimizes error

$$Zb - \xi = (ZA - I_n)\xi + Z\delta$$

averaged over the sampling of images and noise.

This would allow fast reconstruction of the signals using the precomputed matrix Z .

We choose to find \mathbf{Z} by minimizing the *Bayes risk* $\tilde{f}(\mathbf{Z})$ defined, using our model

$$\mathbf{b} = \mathbf{A}\boldsymbol{\xi} + \boldsymbol{\delta}$$

and a quadratic loss function, to be

$$\tilde{f}(\mathbf{Z}) = \mathcal{E} \left(\left\| (\mathbf{Z}\mathbf{A} - \mathbf{I}_n)\boldsymbol{\xi} + \mathbf{Z}\boldsymbol{\delta} \right\|_2^2 \right),$$

where \mathcal{E} denotes expected value.

Assumptions

Assume:

- The random variables ξ and δ are statistically independent.
- The probability distribution for ξ has mean μ_ξ and variance $\eta^2 \mathbf{I}$.
- The probability distribution for δ has mean $\mu_\delta = \mathbf{0}$ and covariance matrix $\beta^2 \mathbf{I}_m$.

(General covariance matrices can be handled.)

Minimizing the Bayes risk

Lemma: Under these assumptions, the Bayes risk is

$$\tilde{f}(\mathbf{Z}) = \|(\mathbf{Z}\mathbf{A} - \mathbf{I}_n)\boldsymbol{\mu}_\xi\|_2^2 + \|(\mathbf{Z}\mathbf{A} - \mathbf{I}_n)\|_F^2 + \alpha^2 \|\mathbf{Z}\|_F^2.$$

Theorem: Consider the problem

$$\hat{\mathbf{Z}} = \arg \min_{\mathbf{Z}} \tilde{f}(\mathbf{Z}).$$

If either $\alpha \neq 0$ or \mathbf{A} has rank m , then the unique global minimizer is

$$\hat{\mathbf{Z}} = (\boldsymbol{\mu}_\xi \boldsymbol{\mu}_\xi^\top + \mathbf{I}) \mathbf{A}^\top [\mathbf{A} (\boldsymbol{\mu}_\xi \boldsymbol{\mu}_\xi^\top + \mathbf{I}) \mathbf{A}^\top + \alpha^2 \mathbf{I}_m]^{-1}.$$

where $\alpha = \eta/\beta$.

This result is interesting, but a full-rank (and probably dense) solution $\hat{\mathbf{Z}}$ is impractical for large-scale problems.

A more practical formulation

$$\tilde{f}(\mathbf{Z}) = \|(\mathbf{Z}\mathbf{A} - \mathbf{I}_n)\boldsymbol{\mu}_\xi\|_2^2 + \|(\mathbf{Z}\mathbf{A} - \mathbf{I}_n)\|_F + \alpha^2 \|\mathbf{Z}\|_F^2.$$

What happens if we add the assumption that $\boldsymbol{\mu}_\xi = \mathbf{0}$, and we seek a matrix \mathbf{Z} of rank at most k ?

We now have the problem solved by our inverse approximation theorem:

If $\sigma_k > \sigma_{k+1}$, the unique global minimizer $\hat{\mathbf{Z}} \in \mathbb{R}^{n \times m}$ is

$$\hat{\mathbf{Z}} = \mathbf{V}_k \boldsymbol{\Psi}_k \mathbf{U}_k^\top,$$

where \mathbf{V}_k contains the first k columns of \mathbf{V} , \mathbf{U}_k contains the first k columns of \mathbf{U} , and

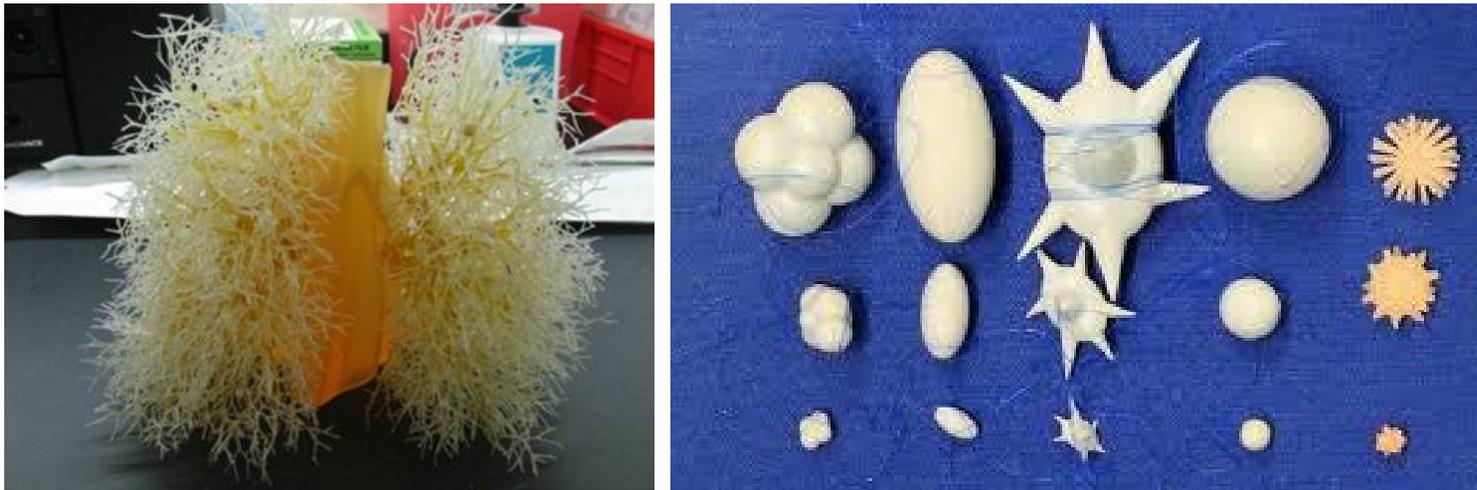
$$\boldsymbol{\Psi}_k = \text{diag}\left(\frac{\sigma_1}{\sigma_1^2 + \alpha^2}, \dots, \frac{\sigma_k}{\sigma_k^2 + \alpha^2}\right).$$

For general covariance matrices, the solution is similar but written using a generalized SVD.

Experiment 3: A deconvolution example

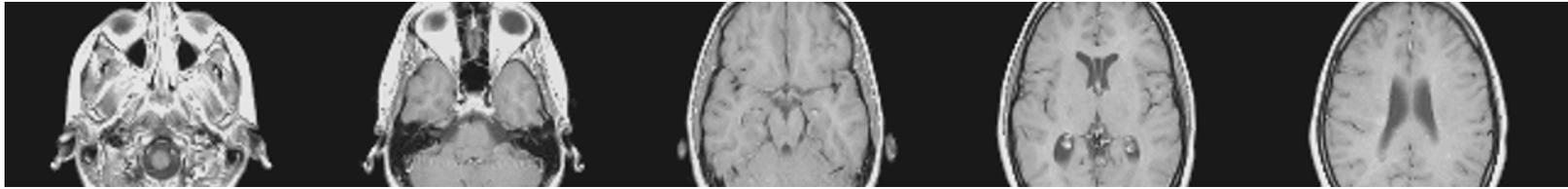
In many imaging applications, calibration data is available.

For example, an imaging device might record an image of a **phantom** with known properties:

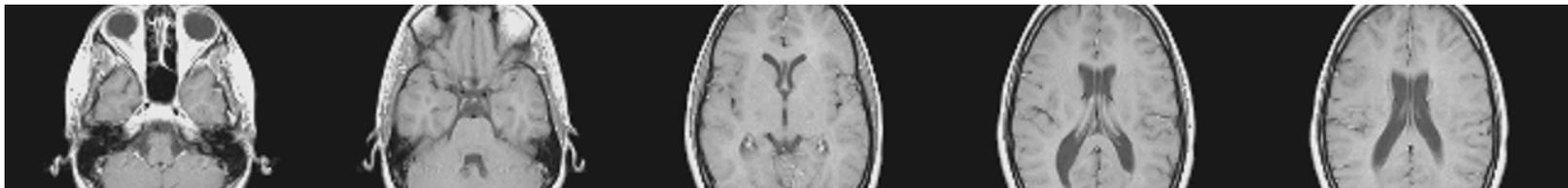


<http://www.fda.gov/>

Our calibration data

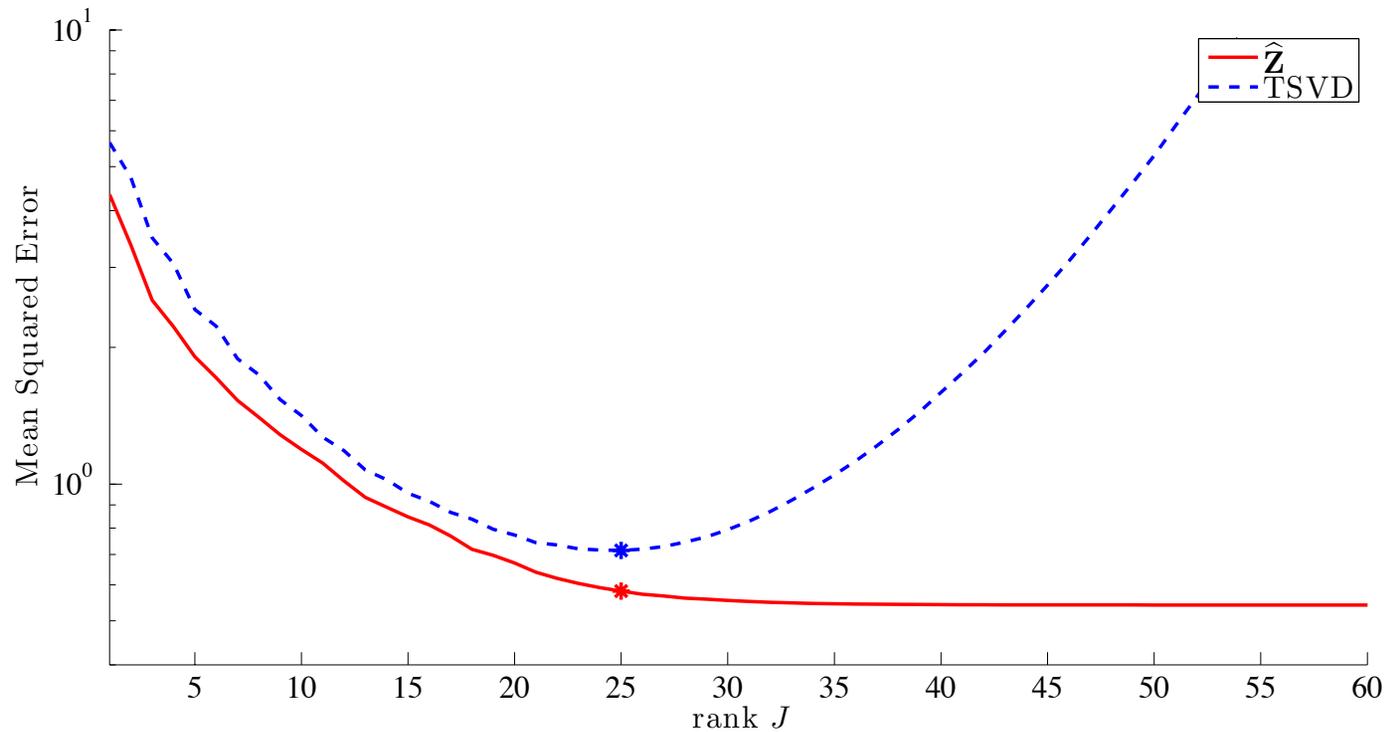


- We use columns from images above to compute a sample covariance matrix.
- We construct our approximate inverse.
- We use columns from images below to evaluate how well it behaves.

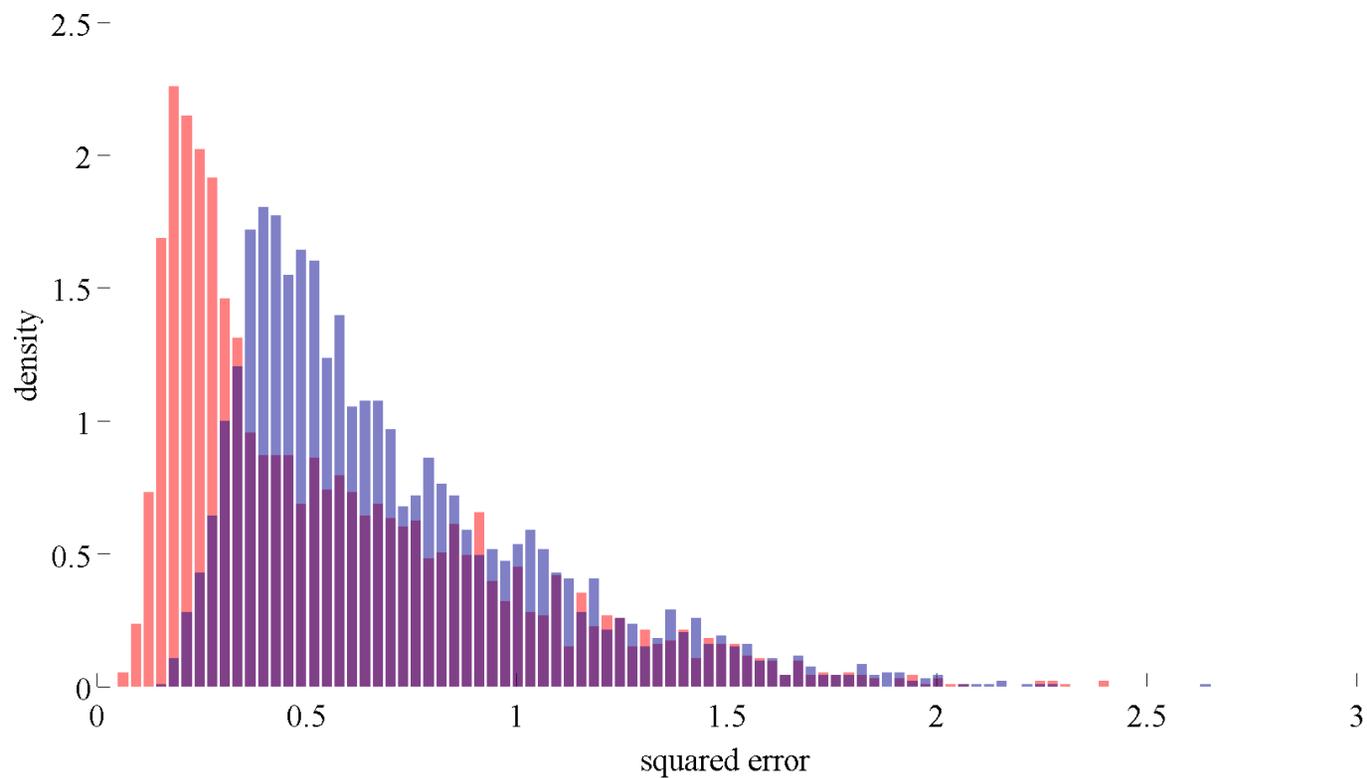


The experiment:

- 1280 possible signals, each 150×1 .
- Compute the sample covariance matrix.
- Scale its diagonal by 1.001 to ensure positive definiteness.
- Matrix $\mathbf{A} \in \mathbb{R}^{150 \times 150}$ represents a Gaussian convolution kernel with variance 4.
- For $\alpha = 0.1$, we compute the optimal low-rank inverse approximation $\widehat{\mathbf{Z}}$ for various ranks k .
- Construct 3072 signals, $\boldsymbol{\xi}^{(k)}$, by convolving columns of the (b) images with the Gaussian kernel, and adding noise $\boldsymbol{\delta}^{(k)}$ (normal distribution with covariance matrix $\alpha^2 \mathbf{I}_m$).
- For various ranks k , calculate the mean squared error $\frac{1}{K} \sum_{k=1}^K e^{(k)}$, where $e^{(k)} = \left\| \mathbf{Z}\mathbf{b}^{(k)} - \boldsymbol{\xi}^{(k)} \right\|_2^2$ and $\mathbf{b}^{(k)} = \mathbf{A}\boldsymbol{\xi}^{(k)} + \boldsymbol{\delta}^{(k)}$, for both $\widehat{\mathbf{Z}}$ and the TSVD reconstruction matrix, \mathbf{A}_k^\dagger .



Comparison of mean squared errors for reconstructions using the optimal approximate inverse matrix $\hat{\mathbf{Z}}$ and the TSVD reconstruction matrix \mathbf{A}_k^\dagger for various ranks.

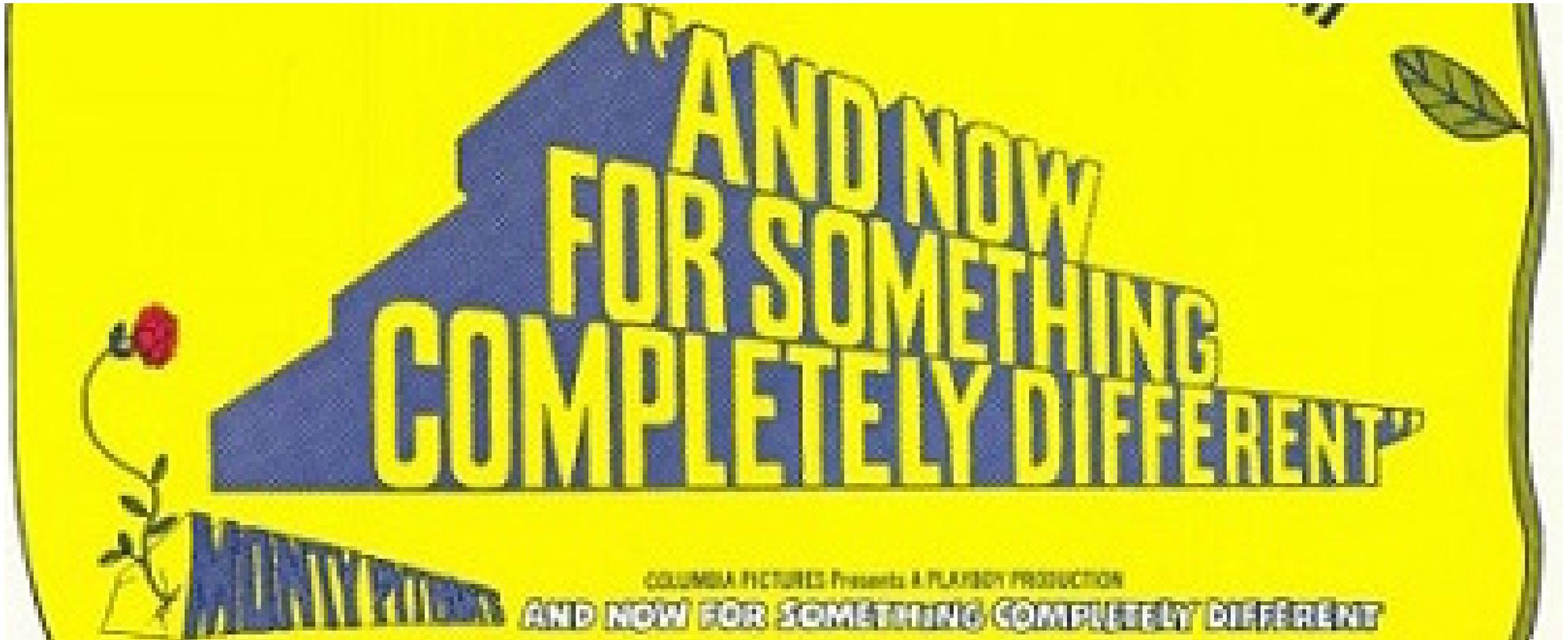


Histogram comparison of squared errors for $\hat{\mathbf{Z}}$ (red bars) and TSVD (blue bars) for rank $k = 25$.

Summary: so far

- We can calculate low-rank approximations to matrix (pseudo)inverses.
- These approximations give us useful reconstructions of images blurred by measurement.

Interval approximation



We have discussed some great ways to approximate matrices...

... very useful for the particular situations above.

But sometimes noise is quite **abnormal**.

Often our data matrix A stores **counts**.

We might count more or less than the true value, but we can never get a **negative** count.

An alternative framework

What if we have an **uncertainty interval** for each matrix element? We would be given matrices \mathbf{L} and \mathbf{U} satisfying

$$\mathbf{L} \leq \mathbf{A}_{true} \leq \mathbf{U},$$

where \mathbf{A}_{true} is unknown.

- We can explicitly impose a nonnegativity assumption, choosing $\mathbf{L} \geq \mathbf{0}$.
- We can easily accommodate different uncertainties in different matrix elements.
- We can even accommodate **missing observations** by setting elements of \mathbf{L} to zero and elements of \mathbf{U} to $+\infty$.
- Our original formulation is a special case, with $\mathbf{L} = \mathbf{U}$.

Our problem becomes

$$\begin{aligned} \min_{\mathbf{X}, \mathbf{Y}, \mathbf{Z}} \quad & f(\mathbf{XY} - \mathbf{Z}) \\ & \ell \leq \mathbf{z} \leq \mathbf{u}, \end{aligned}$$

For f , any matrix norm can be used, and we can also include a weight matrix \mathbf{W} if we care more about keeping some particular elements within their bounds.

So far, our problem is ill-posed.

If $\mathbf{X}\mathbf{Y} = \mathbf{Z}$ solve the problem,

then (for example) so do $\gamma\mathbf{X}, \gamma^{-1}\mathbf{Y}, \mathbf{Z}$ for any positive γ .

Added constraints

- We add constraints:

$$\mathbf{X} \geq \mathbf{0} \quad \text{and} \quad \mathbf{Y} \geq \mathbf{0}.$$

Still ill-posed.

- It is also useful to add terms to the minimization function to “encourage” sparsity in the factors:

$$e^T \mathbf{x} + e^T \mathbf{y}.$$

Now the γ non-uniqueness goes away.

BUT the sparsest \mathbf{X} and \mathbf{Y} are the zero matrices, so we can't weight these terms too heavily without risking a rank-deficient solution to our problem.

So to balance these terms, we also can include these:

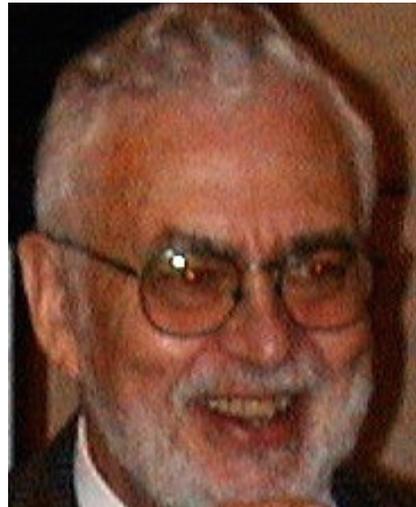
$$-\log \det(\mathbf{X}^T \mathbf{X}) - \log \det(\mathbf{Y} \mathbf{Y}^T).$$

log det is scary!

$$\hat{F}(\mathbf{X}) = -\log \det(\mathbf{X}^T \mathbf{X}),$$



- Both $\mathbf{X}^T \mathbf{X}$ and $\mathbf{Y} \mathbf{Y}^T$ are small (5×5 or 50×50 in our examples, so evaluation is inexpensive.
- The partial derivatives of $\hat{F}(\mathbf{X})$, arranged as a matrix, are $2\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1}$, so $\hat{F}(\mathbf{X})$ and its gradient are easily calculated using (compact) QR factorization of \mathbf{X} .
- The Hessian matrix with respect to entries in each **row** of \mathbf{X} is also easy to calculate given the (compact) QR factors.



Z is scary!

Up until now, we have had only $k(m + n)$ variables.

But Z adds mn more!!

Z is scary!

Really, really scary!

Up until now, we have had only $k(m + n)$ variables.

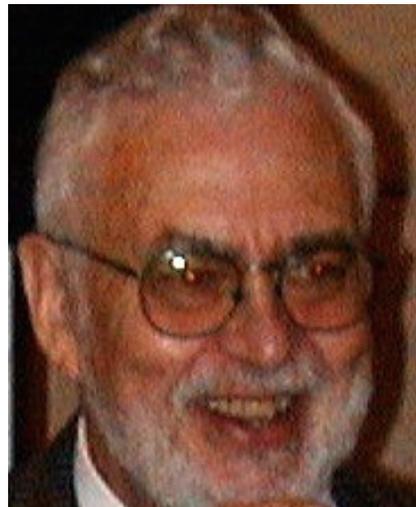
But Z adds mn more!!



Fortunately, I can tell you an optimal choice of \mathbf{Z} for any given \mathbf{X} and \mathbf{Y} :

$$z_{ij}(\mathbf{X}, \mathbf{Y}) = \begin{cases} l_{ij}, & (\mathbf{XY})_{ij} \leq l_{ij} \\ (\mathbf{XY})_{ij}, & l_{ij} \leq (\mathbf{XY})_{ij} \leq u_{ij}, \\ u_{ij}, & u_{ij} \leq (\mathbf{XY})_{ij}. \end{cases}$$

So algorithms that alternate updating \mathbf{X} and updating \mathbf{Y} can be used, and we update \mathbf{Z} using this formula.



Alternating algorithms

$$\min_{\mathbf{X}, \mathbf{Y}, \mathbf{Z}} \|\mathbf{X}\mathbf{Y} - \mathbf{Z}\|_F^2 + \alpha_1 \mathbf{e}^T \mathbf{x} + \alpha_2 \mathbf{e}^T \mathbf{y} - \alpha_3 (\log \det(\mathbf{X}^T \mathbf{X}) - \log \det(\mathbf{Y}\mathbf{Y}^T))$$

Repeat

Update \mathbf{X} .

Determine the optimal \mathbf{Z} based on the current \mathbf{X} and \mathbf{Y} .

Update \mathbf{Y} .

Determine the optimal \mathbf{Z} based on the current \mathbf{X} and \mathbf{Y} .

Until convergence.

- **Blue steps:** not seen in previous algorithms, because now $\mathbf{L} \neq \mathbf{U}$.
- **Red steps:** various choices in the literature.

Updating \mathbf{X} , Option 1

Relevant terms:

$$h(\mathbf{X}) = \|\mathbf{X}\mathbf{Y} - \mathbf{Z}\|_F^2 + \alpha_1 \mathbf{e}^T \mathbf{x} - \alpha_3 \log \det(\mathbf{X}^T \mathbf{X})$$

Option 1 Advocated by Lee & Seung (but without $\log \det$):

Minimize $h(\mathbf{X})$ (or at least take a step that reduces $h(\mathbf{X})$) and then set negative entries in \mathbf{X} to zero.

- Advantages: The computation is inexpensive!
- Disadvantages: No guarantee that the new \mathbf{X} reduces $h(\mathbf{X})$.

Updating \mathbf{X} , Option 2

Relevant terms:

$$h(\mathbf{X}) = \|\mathbf{X}\mathbf{Y} - \mathbf{Z}\|_F^2 + \alpha_1 \mathbf{e}^T \mathbf{x} - \alpha_3 \log \det(\mathbf{X}^T \mathbf{X})$$

Option 2 Advocated by Kim, Sra, & Dhillon (but without $\log \det$):
Minimize $h(\mathbf{X})$ (or at least take a step that reduces $h(\mathbf{X})$) while maintaining nonnegativity of \mathbf{X} .

- Advantages: The new \mathbf{X} reduces $h(\mathbf{X})$ and maintains nonnegativity.
- Advantages: Convergence proof for the alternating algorithm, but under an obnoxious assumption: \mathbf{X} and \mathbf{Y} never become rank deficient. Unfortunately, in practice, they do (without the red term)!
- Disadvantages: The computation is expensive!

Notice that (without the red term) each row of \mathbf{X} can be updated independently.

New option: Downhill, constraint satisfying, and full rank

New option: $h(\mathbf{X}) = \|\mathbf{X}\mathbf{Y} - \mathbf{Z}\|_F^2 + \alpha_1 e^T \mathbf{x} - \alpha_3 \log \det(\mathbf{X}^T \mathbf{X})$

- Search direction: $\mathbf{s} = -\mathbf{PDPg}$, where
 - \mathbf{g} is the gradient of h with respect to \mathbf{x} ,
 - \mathbf{D} is a positive definite scaling matrix,
 - \mathbf{P} is a projection matrix that sets $(\mathbf{P}\mathbf{t})_j$ to zero if $x_j = 0$ and $t_j > 0$.
- Update $\mathbf{x} \leftarrow \mathbf{x} + \nu \mathbf{s}$, where ν is determined by a line search.

This is the same search direction used by Kim, Sra, & Dhillon, except that we include the $\log \det$ term in our problem formulation.

So now we retain full rank in the iterates.

What can we prove?

Assume that we never take an uphill step. So

$$F(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) = \|\mathbf{X}\mathbf{Y} - \mathbf{Z}\|_F^2 + \alpha_1 \mathbf{e}^T \mathbf{x} + \alpha_2 \mathbf{e}^T \mathbf{y} - \alpha_3 (\log \det(\mathbf{X}^T \mathbf{X}) + \log \det(\mathbf{Y}\mathbf{Y}^T))$$

is never larger than its initial value f_0 .

- The log det term is (usually) small.
- The linear term gives us a bound on $\|\mathbf{Y}\|$.
- The log det term then gives us a nonzero bound on the smallest eigenvalue of $\mathbf{Y}\mathbf{Y}^T$.
- So we can guarantee that our iterates are full rank!

$$F(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) = \|\mathbf{X}\mathbf{Y} - \mathbf{Z}\|_F^2 + \alpha_1 \mathbf{e}^T \mathbf{x} + \alpha_2 \mathbf{e}^T \mathbf{y} - \alpha_3 (\log \det(\mathbf{X}^T \mathbf{X}) + \log \det(\mathbf{Y}\mathbf{Y}^T))$$

Notice:

- We have guaranteed that $(\mathbf{Y}\mathbf{Y}^T)^{-1}$ is uniformly bounded throughout our iteration.
- $\mathbf{Y}\mathbf{Y}^T$ is the Hessian matrix w.r.t. \mathbf{X} variables for the quadratic terms.
- Linear systems involving $\mathbf{Y}\mathbf{Y}^T$ are easy to solve using a Cholesky factorization, or a QR factorization of \mathbf{Y}^T , and we already needed a factorization in order to evaluate the $\log \det$ term.

Therefore, $(\mathbf{Y}\mathbf{Y}^T)^{-1}$ is an ideal candidate for the scaling matrix \mathbf{D} in the \mathbf{X} iteration and makes the iteration **Newton-like**.

Interchange \mathbf{X} and \mathbf{Y} in this discussion to get the same result for the \mathbf{Y} iteration.

Results: Document classification

DUC 2004 data: Term-document matrix for 500 documents (some repeats) that are “correctly” classified in 50 classes of 10 documents each.

Measures of correctness: Let

A = number of agreements: document i and j in same/different class

D = number of disagreements,

so that $A+D = \text{comb}(n,2)$. Then

.9846 = RI = Rand index (1971) = $A / \text{comb}(n,2) = \text{prob agree}$

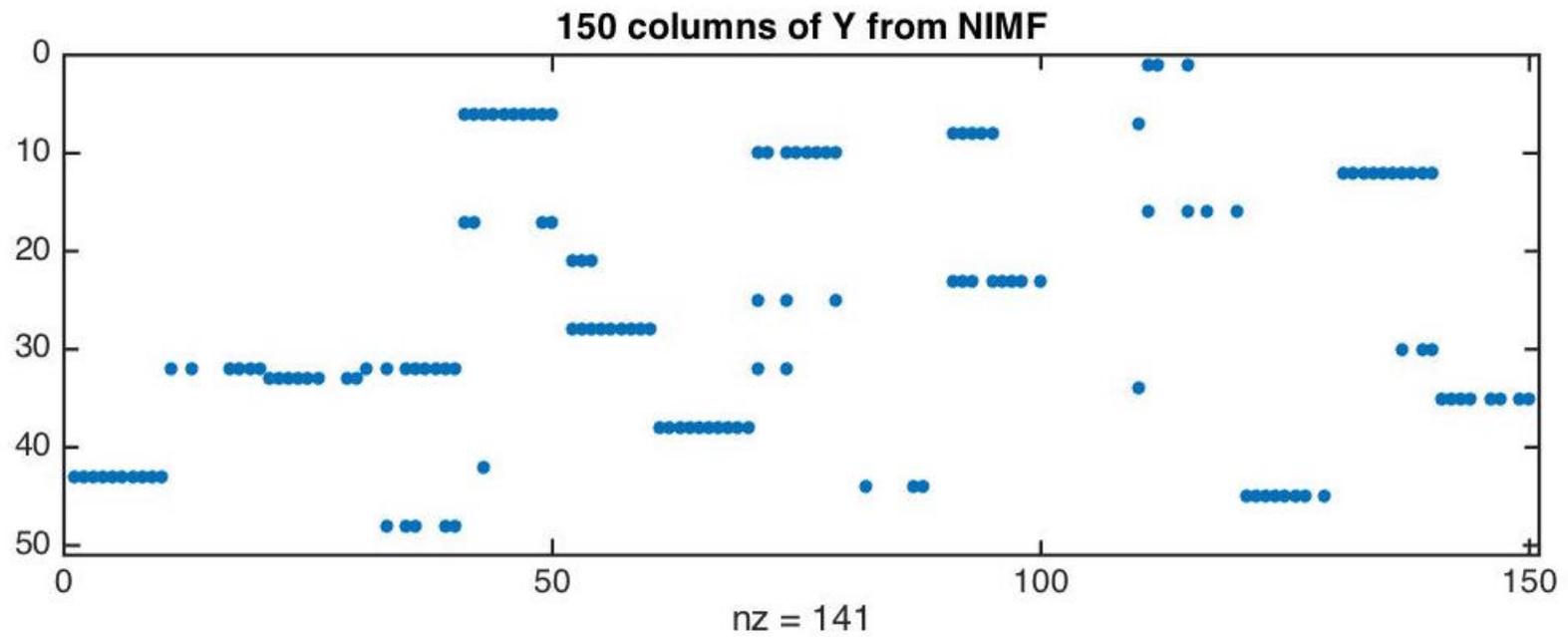
.9692 = HI = Hubert index (1977) = $(A - D) / \text{comb}(n,2) = \text{RI} - \text{MI}$

.6252 = AR = adjusted RI, corrected for chance Hubert+Arabic (1985)

.7500 = AMI = adjusted mutual index

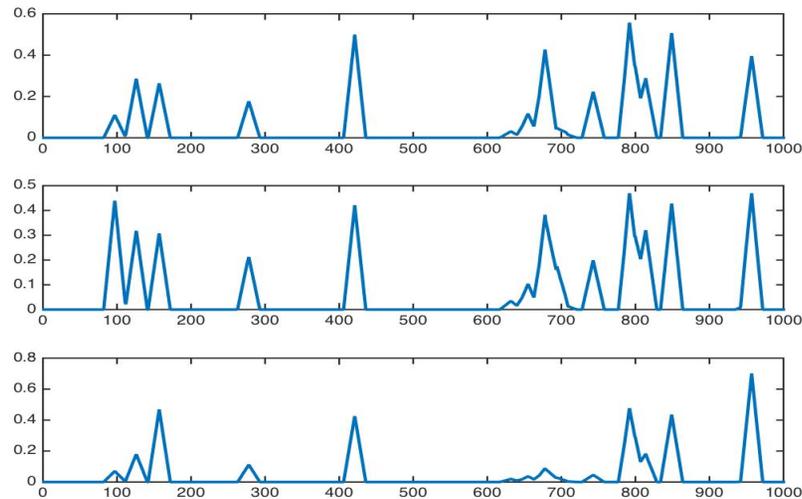
Comparable to competing methods.

The Computed Y matrix



Results: Spectroscopy

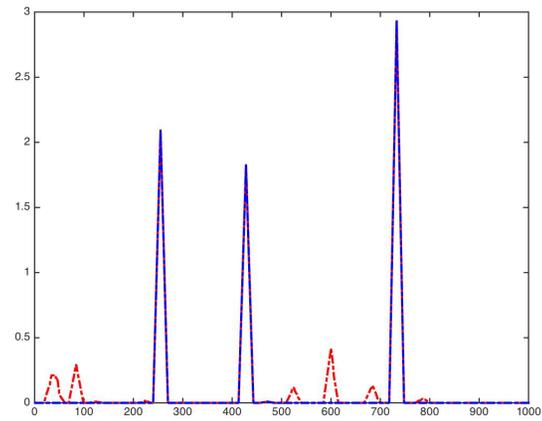
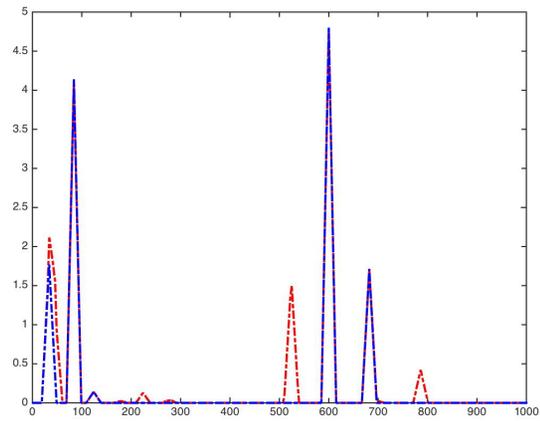
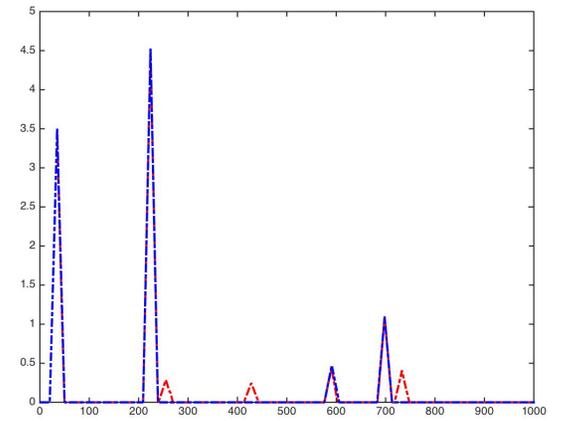
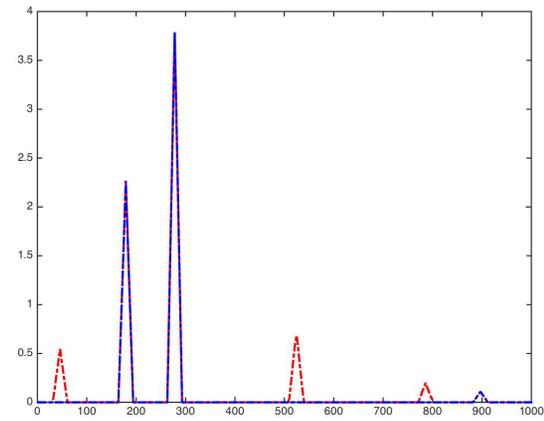
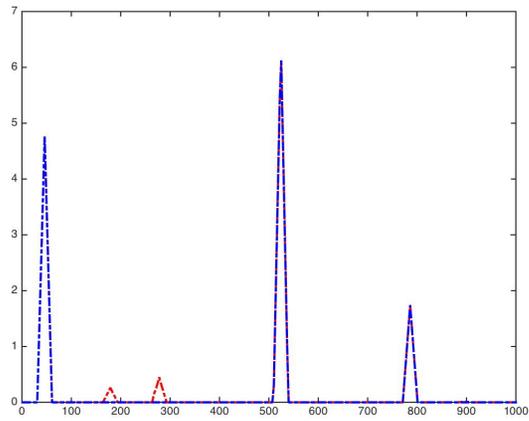
Given 50 **noisy combinations** of the 5 underlying spectra



(and 47 more)

Reconstruct the 5 spectra.

Blue = True, Red = Computed



Conclusions

A fresh look at low-rank matrix approximation is productive and useful!

- **Inverse Approximation of Matrices**
 - We can calculate low-rank approximations to matrix (pseudo)inverses.
 - These approximations give us useful reconstructions of images blurred by measurement when calibration data is available.
- **Approximation of Interval Matrices**
 - We have a new formulation of matrix approximation that allows error bounds and weights on each matrix entry.
 - We have a descent algorithm that takes Newton-like steps to improve the free variables.
 - We have demonstrated promising results for document classification and signal recovery.

References for Inverse Approximation:

- Julianne Chung, Matthias Chung, and Dianne P. O’Leary
“Optimal Regularized Low-Rank Inverse Approximation,” *Linear Algebra and Its Applications*, 468(1) (2015) 260-269.
- Julianne M. Chung, Matthias Chung, and Dianne P. O’Leary, “Optimal Filters from Calibration Data for Image Deconvolution with Data Acquisition Errors,” *Journal of Mathematical Imaging and Vision*, 44(3), pp. 366–374 (2012)

Thank you, Gene!



And thank you to my collaborators



Julianne Chung



Matthias Chung



John Conroy



Yi-Kai Liu

www.cs.umd.edu/users/oleary/talkview.pdf