

# Algebraic Iterative Methods with Noisy Data

## Semi-Convergence and Stopping Rules

Per Christian Hansen

DTU Compute  
Department of Applied Mathematics and Computer Science  
Technical University of Denmark



# Plan for Today

- 1 Enter the noise  $\rightarrow$  semi-convergence.
- 2 SVD analysis  $\rightarrow$  iteration error and noise error.
- 3 Analysis of Landweber and Kaczmarz with projections.
- 4 The need for stopping rules.
- 5 Fit to noise level; min. of prediction error; extract all information.
- 6 Estimation of trace term and noise level.

## Points to take home today:

- For noisy data we rely on semi-convergence of the iterative methods.
- We have a good theoretical understanding of this phenomenon.
- We need to terminate iterations at smallest reconstruction error.
- Several stopping rules are available.
- We can estimate the crucial parameters in these rules.

## Real Problems Have Noisy Data

So far we have discussed how to solve  $\mathbf{A} \mathbf{x} = \mathbf{b}$  by iterative methods. But when noise is present in the data, we don't quite want to do that!

We assume that the data, in the form of the right-hand side  $\mathbf{b}$ , is a sum of “clean” noise-free data  $\mathbf{A} \bar{\mathbf{x}}$  from the ground-truth image plus a noise component  $\mathbf{e}$ :

$$\mathbf{b} = \mathbf{A} \bar{\mathbf{x}} + \mathbf{e}, \quad \bar{\mathbf{x}} = \text{ground truth}, \quad \mathbf{e} = \text{noise}.$$

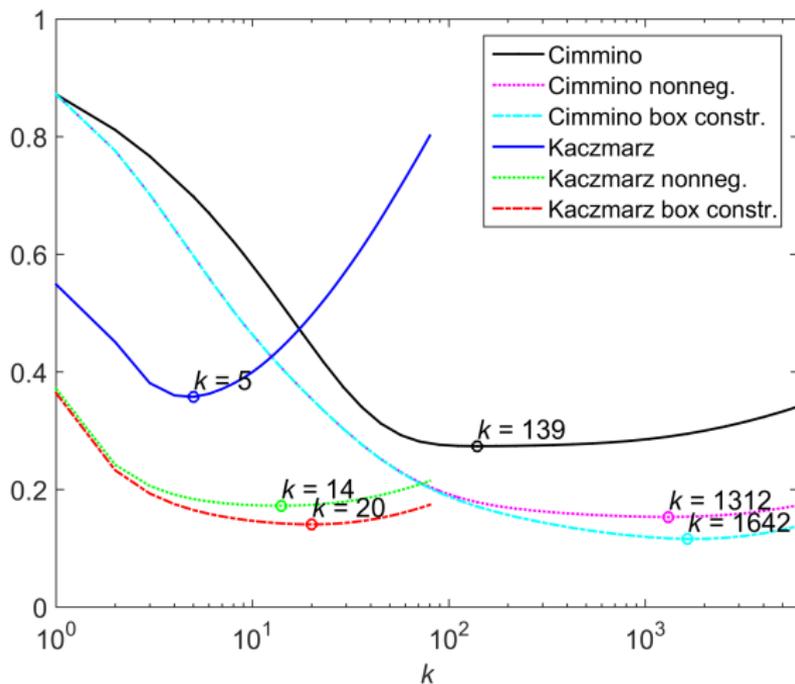
The plain-vanilla or naïve solution  $\mathbf{x}^{\text{naïve}} = \mathbf{A}^{-1} \mathbf{b}$  is undesired, because it has a large component coming from the noise in the data:

$$\mathbf{x}^{\text{naïve}} = \mathbf{A}^{-1} \mathbf{b} = \mathbf{A}^{-1} (\mathbf{A} \bar{\mathbf{x}} + \mathbf{e}) = \bar{\mathbf{x}} + \mathbf{A}^{-1} \mathbf{e}.$$

The component  $\mathbf{A}^{-1} \mathbf{e}$  typically dominates over  $\bar{\mathbf{x}}$ , because  $\mathbf{A}$  is an ill conditioned matrix.

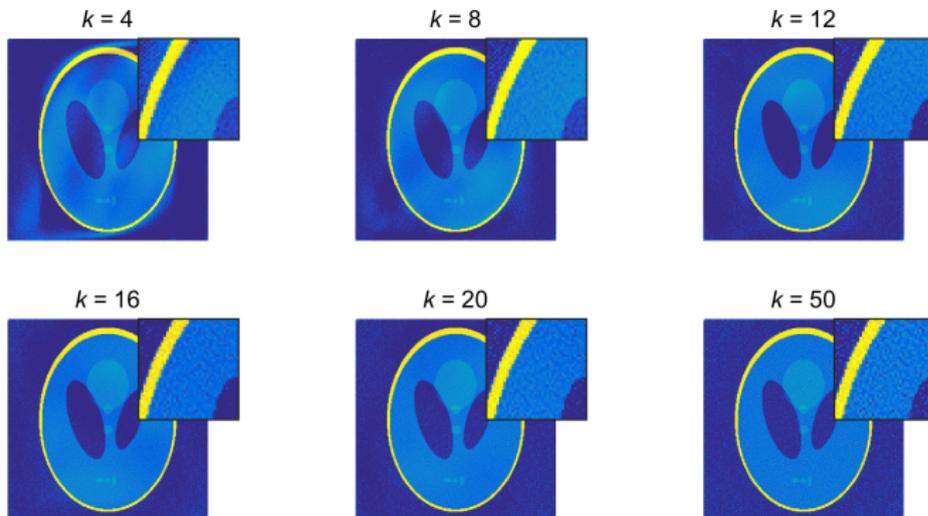
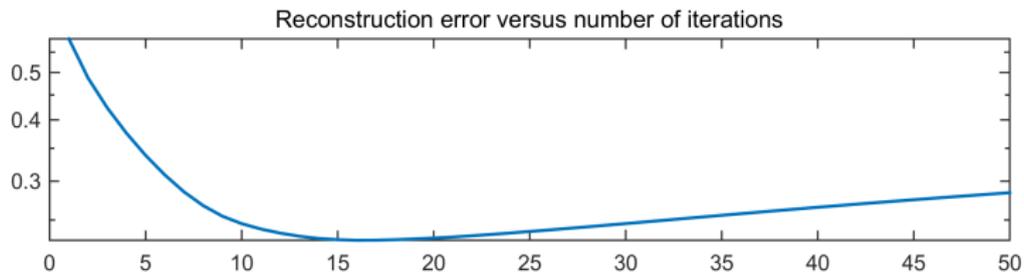
But something interesting happens during the iterations ...

# Convergence for Noisy Data



For all six methods the error  $\|\bar{\mathbf{x}} - \mathbf{x}^{(k)}\|_2$  decreases until it reaches a **minimum**, shown by the circles, after which it starts to **increase again**.

# Semi-Convergence For Kaczmarz's Method



# Semi-Convergence

This behavior is often referred to as **semi-convergence**:

- During the initial iterations, the iteration vector  $\mathbf{x}^{(k)}$  approaches the desired – but un-obtainable – solution  $\bar{\mathbf{x}}$  to the noise-free problem.
- During later iterations,  $\mathbf{x}^{(k)}$  converges to the undesired naïve solution associated with the particular AIR method (i.e.,  $\mathbf{A}^{-1}\mathbf{b}$  if the system matrix is invertible).

We want to stop the iterations just when the convergence behavior changes from the former to the latter.

Then we achieve a **regularized solution** – an approximation to the noise-free solution which is not too perturbed by the noise in the data.

Today we explain

- 1 why we have semi-convergence for noisy data, and
- 2 how to stop the iterations at the right time.

## Analysis of Landweber's Method I

With an arbitrary starting vector  $\mathbf{x}^{(0)}$ , the  $k$ th Landweber iterate is:

$$\begin{aligned}\mathbf{x}^{(k)} &= \mathbf{x}^{(k-1)} + \omega \mathbf{A}^T (\mathbf{b} - \mathbf{A} \mathbf{x}^{(k-1)}) \\ &= (\mathbf{I} - \omega \mathbf{A}^T \mathbf{A}) \mathbf{x}^{(k-1)} + \omega \mathbf{A}^T \mathbf{b} \\ &= (\mathbf{I} - \omega \mathbf{A}^T \mathbf{A}) \left[ (\mathbf{I} - \omega \mathbf{A}^T \mathbf{A}) \mathbf{x}^{(k-2)} + \omega \mathbf{A}^T \mathbf{b} \right] + \omega \mathbf{A}^T \mathbf{b} \\ &= (\mathbf{I} - \omega \mathbf{A}^T \mathbf{A})^2 \mathbf{x}^{(k-2)} + ((\mathbf{I} - \omega \mathbf{A}^T \mathbf{A}) + \mathbf{I}) \omega \mathbf{A}^T \mathbf{b} \\ &= (\mathbf{I} - \omega \mathbf{A}^T \mathbf{A})^3 \mathbf{x}^{(k-3)} + ((\mathbf{I} - \omega \mathbf{A}^T \mathbf{A})^2 + (\mathbf{I} - \omega \mathbf{A}^T \mathbf{A}) + \mathbf{I}) \omega \mathbf{A}^T \mathbf{b} \\ &= \dots \\ &= (\mathbf{I} - \omega \mathbf{A}^T \mathbf{A})^k \mathbf{x}^{(0)} + \\ &\quad \left[ (\mathbf{I} - \omega \mathbf{A}^T \mathbf{A})^{k-1} + (\mathbf{I} - \omega \mathbf{A}^T \mathbf{A})^{k-2} + \dots + \mathbf{I} \right] \omega \mathbf{A}^T \mathbf{b} \\ &= (\mathbf{I} - \omega \mathbf{A}^T \mathbf{A})^k \mathbf{x}^{(0)} + \sum_{j=0}^{k-1} (\mathbf{I} - \omega \mathbf{A}^T \mathbf{A})^j \omega \mathbf{A}^T \mathbf{b}.\end{aligned}$$

## Analysis of Landweber's Method II

For simplicity we now assume that  $\mathbf{x}^{(0)} = 0$ . We insert the SVD of the system matrix  $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$  and use  $\mathbf{I} = \mathbf{V} \mathbf{V}^T$ :

$$\mathbf{x}^{(k)} = \mathbf{V} \sum_{j=0}^{k-1} (\mathbf{I} - \omega \mathbf{\Sigma}^2)^j \omega \mathbf{\Sigma} \mathbf{U}^T \mathbf{b} = \mathbf{V} \mathbf{\Phi}^{(k)} \mathbf{\Sigma}^{-1} \mathbf{U}^T \mathbf{b},$$

where we introduced the  $n \times n$  diagonal matrix

$$\mathbf{\Phi}^{(k)} = \sum_{j=0}^{k-1} (\mathbf{I} - \omega \mathbf{\Sigma}^2)^j \omega \mathbf{\Sigma}^2 = \omega \mathbf{\Sigma}^2 \sum_{j=0}^{k-1} (\mathbf{I} - \omega \mathbf{\Sigma}^2)^j = \begin{pmatrix} \phi_1^{(k)} & & \\ & \phi_2^{(k)} & \\ & & \ddots \end{pmatrix}$$

with diagonal elements

$$\phi_i^{(k)} = \omega \sigma_i^2 \sum_{j=0}^{k-1} (1 - \omega \sigma_i^2)^j, \quad i = 1, 2, \dots, n.$$

## Analysis of Landweber's Method III

The sum  $\sum_{j=0}^{k-1} (1 - \omega \sigma_i^2)^j$  is a geometric series:

$$\sum_{j=0}^{k-1} z^j = (1 - z^k)/(1 - z),$$

and thus for  $i = 1, 2, \dots, n$  we have:

$$\phi_i^{(k)} = \omega \sigma_i^2 \sum_{j=0}^{k-1} (1 - \omega \sigma_i^2)^j = \omega \sigma_i^2 \frac{1 - (1 - \omega \sigma_i^2)^k}{1 - (1 - \omega \sigma_i^2)} = 1 - (1 - \omega \sigma_i^2)^k$$

leading to a simple expression for the diagonal Landweber filter matrix

$$\Phi^{(k)} = \begin{pmatrix} 1 - (1 - \omega \sigma_1^2)^k & & \\ & 1 - (1 - \omega \sigma_2^2)^k & \\ & & \ddots \end{pmatrix}.$$

# SVD Expression of Landweber Iteration Vectors

After  $k$  iterations we have obtained a regularized solution

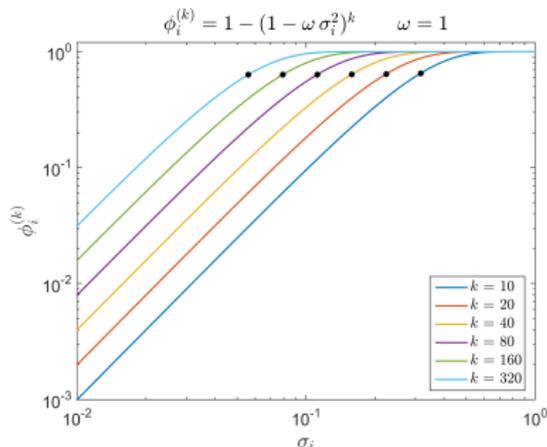
$$\mathbf{x}^{(k)} = \mathbf{V} \Phi^{(k)} \Sigma^{-1} \mathbf{U}^T \mathbf{b} = \sum_{i=1}^n \phi_i^{(k)} \frac{\mathbf{u}_i^T \mathbf{b}}{\sigma_i} \mathbf{v}_i$$

$$\phi_i^{(k)} = 1 - (1 - \omega \sigma_i^2)^k \approx \begin{cases} 1 & \text{for } \sigma_i \gg \omega, \\ k\omega \sigma_i^2 & \text{for } \sigma_i \ll \omega, \end{cases}$$

SVD components for large  $\sigma_i$  are **essentially unfiltered**; those for small  $\sigma_i$  are **damped by a factor  $\propto \sigma_i^2$** .

The breakpoint is approximately for  $\sigma_i \approx 1/\sqrt{k\omega}$ , see black dots  $\bullet \rightarrow$

More SVD components are included as we perform more iterations ( $\sim$  TSVD).



## About The Relaxation Parameter

This analysis also provides an asymptotic convergence analysis for Landweber's method as  $k \rightarrow \infty$ .

For the geometric series to converge we must require that the relaxation parameter  $\omega$  satisfies

$$|1 - \omega \sigma_i^2| < 1, \quad i = 1, 2, \dots, m,$$

which implies that we must have

$$\omega < 2/\sigma_1^2 = 2/\|\mathbf{A}\|_2^2 = 2/\|\mathbf{A}^T \mathbf{A}\|_2.$$

When this is satisfied then  $\phi_i^{(k)} \rightarrow 1$  for all  $i$  and thus  $\Phi^{(k)} \rightarrow \mathbf{I}$  for  $k \rightarrow \infty$ .

Hence  $\mathbf{x}^{(k)}$  converges to the naïve noisy solution  $\mathbf{V} \Sigma^{-1} \mathbf{U}^T \mathbf{b} = \mathbf{A}^{-1} \mathbf{b}$ .

## Iteration Error and Noise Error

Get ready for semi-convergence analysis!

We can always split the **reconstruction error** for  $\mathbf{x}^{(k)}$  into two components:

$$\bar{\mathbf{x}} - \mathbf{x}^{(k)} = (\bar{\mathbf{x}} - \bar{\mathbf{x}}^{(k)}) + (\bar{\mathbf{x}}^{(k)} - \mathbf{x}^{(k)}).$$

The “clean” iteration vector  $\bar{\mathbf{x}}^{(k)}$  is defined as the iteration vector obtained when we apply  $k$  steps of Landweber’s method to the noise-free data  $\mathbf{A} \bar{\mathbf{x}}$ .

- The first component  $\bar{\mathbf{x}} - \bar{\mathbf{x}}^{(k)}$  is the **iteration error** which is an approximation error caused by the finite number of iterations, and which is independent of the noise in the data.
- The second component  $\bar{\mathbf{x}}^{(k)} - \mathbf{x}^{(k)}$  is the **noise error** which is due to the presence of the data errors and causing the actual iteration vector  $\mathbf{x}^{(k)}$  to differ from the “clean” iteration vector  $\bar{\mathbf{x}}^{(k)}$ .

## SVD Analysis of Iteration Error and Noise Error

$$\bar{\mathbf{x}} = \sum_{i=1}^n \mathbf{v}_i^T \bar{\mathbf{x}} \mathbf{v}_i,$$

$$\mathbf{u}_i^T \mathbf{b} = \mathbf{u}_i^T (\mathbf{A} \bar{\mathbf{x}} + \mathbf{e}) = \mathbf{u}_i^T (\mathbf{A} \bar{\mathbf{x}}) + \mathbf{u}_i^T \mathbf{e},$$

$$\mathbf{u}_i^T (\mathbf{A} \bar{\mathbf{x}}) = \mathbf{u}_i^T \sum_{j=1}^n \mathbf{u}_j \sigma_j \mathbf{v}_j^T \bar{\mathbf{x}} = \sigma_i \mathbf{v}_i^T \bar{\mathbf{x}}.$$

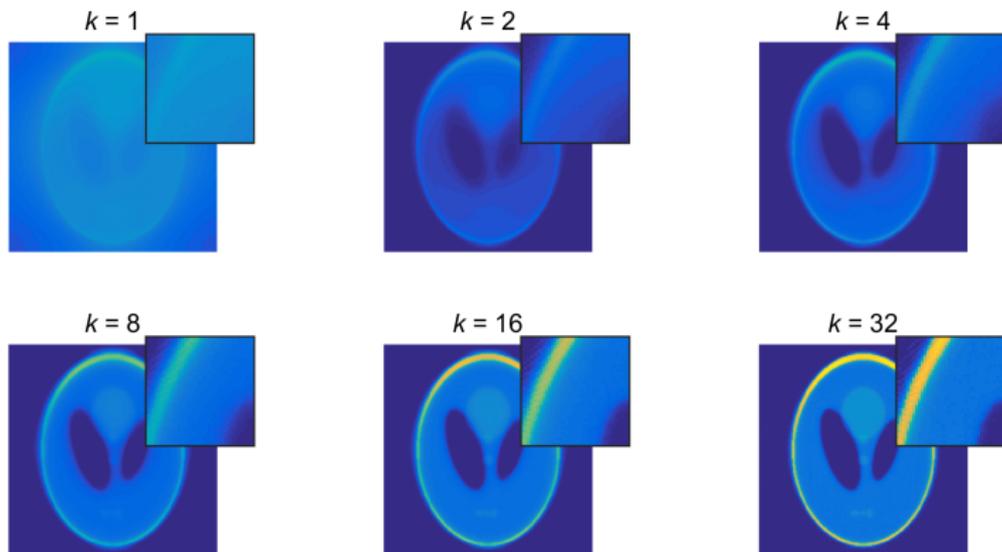
It follows that

$$\bar{\mathbf{x}} - \bar{\mathbf{x}}^{(k)} = \sum_{i=1}^n \mathbf{v}_i^T \bar{\mathbf{x}} \mathbf{v}_i - \sum_{i=1}^n \phi_i^{(k)} \frac{\mathbf{u}_i^T (\mathbf{A} \bar{\mathbf{x}})}{\sigma_i} \mathbf{v}_i = \sum_{i=1}^n (1 - \phi_i^{(k)}) \mathbf{v}_i^T \bar{\mathbf{x}} \mathbf{v}_i,$$

$$\mathbf{x}^{(k)} - \bar{\mathbf{x}}^{(k)} = \sum_{i=1}^n \phi_i^{(k)} \frac{\mathbf{u}_i^T \mathbf{b}}{\sigma_i} \mathbf{v}_i - \sum_{i=1}^n \phi_i^{(k)} \frac{\mathbf{u}_i^T (\mathbf{A} \bar{\mathbf{x}})}{\sigma_i} \mathbf{v}_i = \sum_{i=1}^n \phi_i^{(k)} \frac{\mathbf{u}_i^T \mathbf{e}}{\sigma_i} \mathbf{v}_i.$$

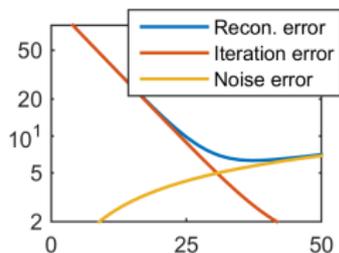
As  $k$  increases, and more  $\phi_i^{(k)}$  approach 1, the iteration error tends to zero while the noise error increases because more noise components are included.

## Progression of the Iteration Error

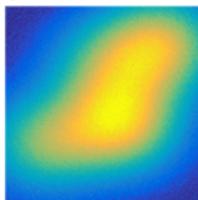


“Clean” iteration vectors  $\bar{\mathbf{x}}^{(k)}$  for Landweber applied to noise-free data  $\mathbf{A} \bar{\mathbf{x}}$ . As we take more iterations we include more SVD components with higher frequencies and we obtain sharper edges in the images.

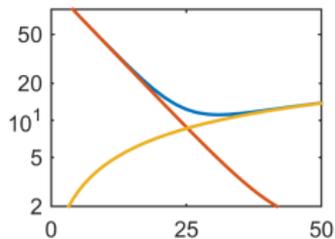
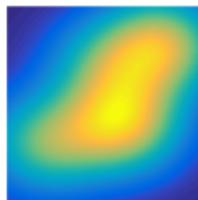
# Progression of All Error Types



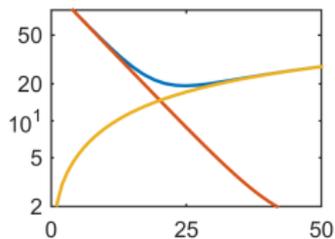
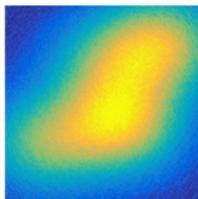
$\rho = 0.01, k = 37$



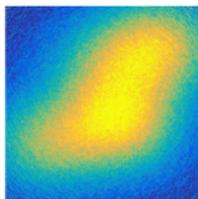
Phantom



$\rho = 0.02, k = 31$



$\rho = 0.04, k = 25$



The iteration error  $\|\bar{\mathbf{x}} - \bar{\mathbf{x}}^{(k)}\|_2$  is independent of the noise.

The noise error  $\|\bar{\mathbf{x}}^{(k)} - \mathbf{x}^{(k)}\|_2$  increases when the noise increases.

The combined reconstruction error  $\|\bar{\mathbf{x}} - \mathbf{x}^{(k)}\|_2$  has a minimum  $\mapsto$  semi-convergence.

## Landweber's Method with Projections

When projections are incorporated in the iterative algorithm, we can no longer perform an SVD analysis. But it can be shown that:

$$\begin{aligned}\|\bar{\mathbf{x}} - \bar{\mathbf{x}}^{(k)}\|_2 &\leq (1 - \omega \sigma_n^2)^k \|\bar{\mathbf{x}}\|_2, \\ \|\bar{\mathbf{x}}^{(k)} - \mathbf{x}^{(k)}\|_2 &\leq \frac{1 - (1 - \omega \sigma_n^2)^k}{\sigma_n^2} \|\mathbf{A}\|_2 \|\mathbf{e}\|_2.\end{aligned}$$

Here  $\sigma_n$  is the smallest singular value of  $\mathbf{A}$ ; hence  $\omega \sigma_n^2$  is quite small, but the iteration error usually decreases faster.

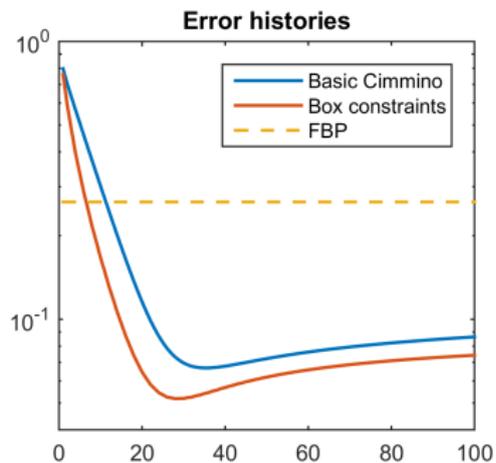
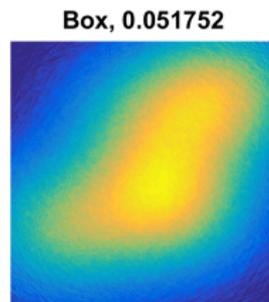
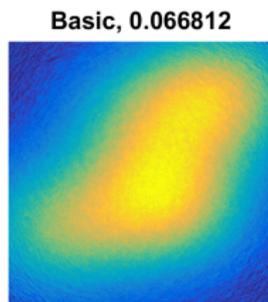
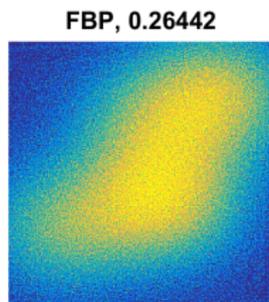
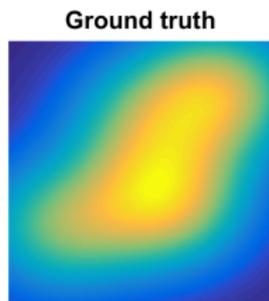
From  $(1 - \epsilon)^k = 1 - k\epsilon + 1/2k(k+1)\epsilon^2 + \dots$  it follows that

$$\frac{1 - (1 - \omega \sigma_n^2)^k}{\sigma_n^2} = \frac{1 - (1 - k\omega \sigma_n^2 + O(\sigma_n^4))}{\sigma_n^2} = k\omega + O(\sigma_n^2)$$

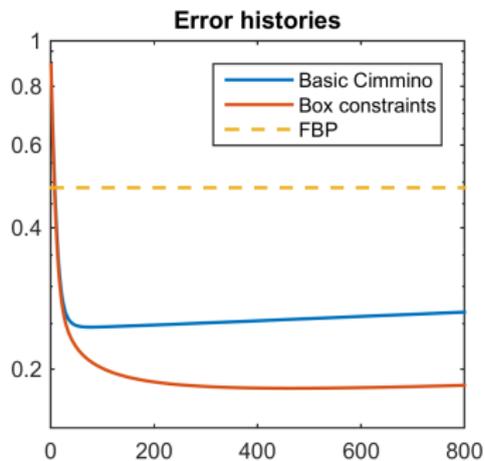
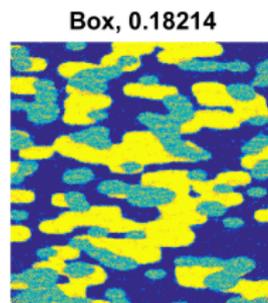
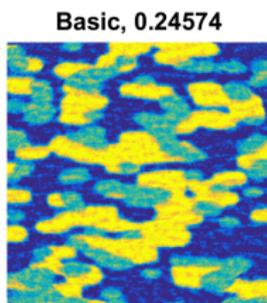
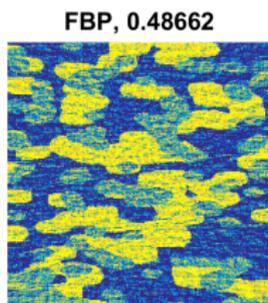
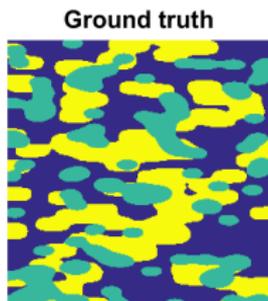
and we obtain the approximate bound for the [noise error](#):

$$\|\bar{\mathbf{x}}^{(k)} - \mathbf{x}^{(k)}\|_2 \lesssim k\omega \|\mathbf{A}\|_2 \|\mathbf{e}\|_2.$$

# Projected Cimmino, Example I



# Projected Cimmino, Example II



## Analysis of Kaczmarz's Method

To study the noise error in Kaczmarz's method we take this approach.

Define the “splitting”

$$\mathbf{A}\mathbf{A}^T = \mathbf{L} + \mathbf{D} + \mathbf{L}^T,$$

$\mathbf{D}$  = diagonal matrix with diagonal elements of  $\mathbf{A}\mathbf{A}^T$

$\mathbf{L}$  = strictly lower triangular matrix (zeros on the diagonal).

Also define the lower triangular matrix

$$\widehat{\mathbf{L}} = (\mathbf{D} + \omega \mathbf{L})^{-1}.$$

Then one iteration of Kaczmarz's method can be written as

$$\mathbf{x}^{(k)} = P_C(\mathbf{x}^{(k-1)} + \omega \mathbf{A}^T \widehat{\mathbf{L}} (\mathbf{b} - \mathbf{A} \mathbf{x}^{(k-1)}))$$

This is for *purely theoretical use*; it should not be used for computations!

## Upper Bound for Iteration Error

Following the same approach as for Landweber, we obtain:

$$\mathbf{x}^{(k)} - \bar{\mathbf{x}}^{(k)} = \sum_{j=0}^{k-1} (\mathbf{I} - \omega \mathbf{A}^T \hat{\mathbf{L}} \mathbf{A})^j \mathbf{A}^T \hat{\mathbf{L}} \mathbf{b},$$

and it can be shown that for both the un-constrained and the constrained problem, the noise error is bounded as

$$\|\bar{\mathbf{x}}^{(k)} - \mathbf{x}^{(k)}\|_2 \leq \frac{1 - (1 - \omega \varsigma^2)^k}{\varsigma^2} \|\mathbf{A}^T \hat{\mathbf{L}} \mathbf{e}\|_2 + O(\varsigma^2),$$

where  $\varsigma$  is the smallest nonzero singular value of the matrix  $\mathbf{D}^{1/2} \hat{\mathbf{L}} \mathbf{A}$ .

Following the previous arguments, we again obtain an approximate upper bound of the form

$$\|\bar{\mathbf{x}}^{(k)} - \mathbf{x}^{(k)}\|_2 \lesssim k\omega \|\mathbf{A}^T \hat{\mathbf{L}} \mathbf{e}\|_2.$$

## Can We Do Something About $\|\mathbf{A}^T \widehat{\mathbf{L}} \mathbf{e}\|_2$ ?

Can we replace  $\alpha = \|\mathbf{A}^T \widehat{\mathbf{L}} \mathbf{e}\|_2$  with the upper bound  $\beta = \|\mathbf{A}\|_2 \|\widehat{\mathbf{L}}\|_2 \|\mathbf{e}\|_2$ ?

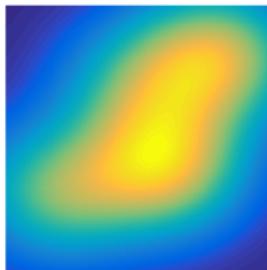
Example:  $N = 64$ , no. detector pixels = 90, length of detector =  $L$  and  $\sqrt{2}L$  where  $L$  is the object's side length.

projection angles	$\ \mathbf{e}\ _2$	short detector		long detector	
		$\alpha$	$\beta$	$\alpha$	$\beta$
$\theta = 1^\circ, 2^\circ, 3^\circ, \dots, 180^\circ$	49	7.6	1251	28.2	$6.4 \cdot 10^7$
$\theta = 3^\circ, 6^\circ, 9^\circ, \dots, 180^\circ$	70	7.4	403	22.6	$2.1 \cdot 10^7$
$\theta = 6^\circ, 12^\circ, 18^\circ, \dots, 180^\circ$	128	6.1	181	7.1	2232

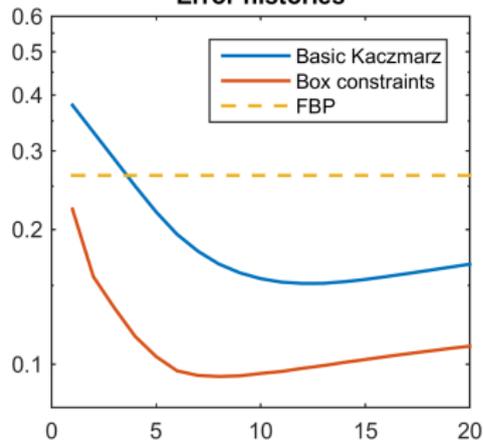
Conclusion:  $\beta \gg \alpha$ ; not a good idea to use  $\beta$ .

# Kaczmarz Semi-Convergence, Example I

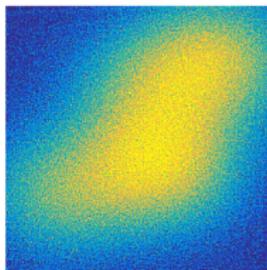
Ground truth



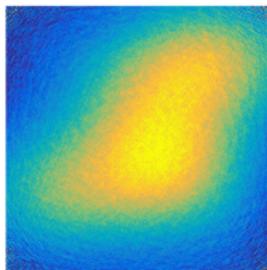
Error histories



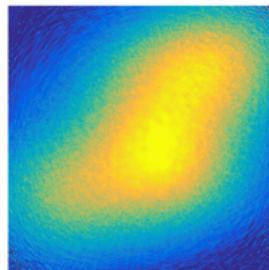
FBP, 0.26442



Basic, 0.15157

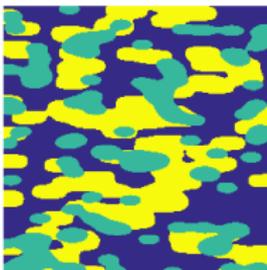


Box, 0.093865

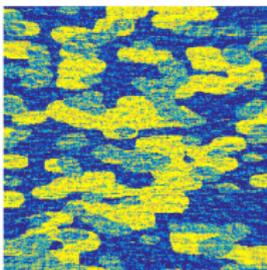


# Kaczmarz Semi-Convergence, Example II

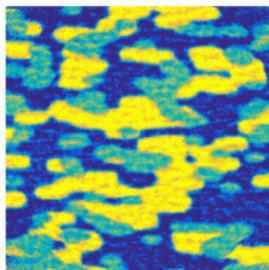
Ground truth



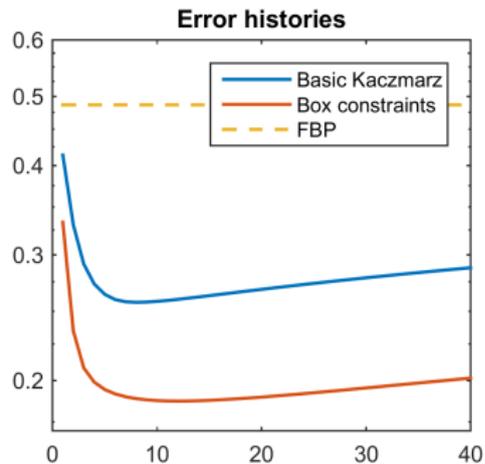
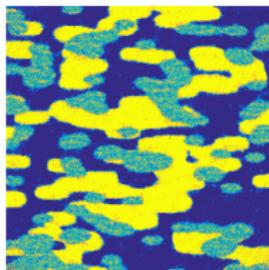
FBP, 0.48662



Basic, 0.25735



Box, 0.18716



## Stopping Rules

To successfully use the iterative methods for noisy data, we obviously need an automatic method – a **stopping rule** – for terminating the iterations at, or near, the point of semi-convergence where the reconstruction error  $\bar{\mathbf{x}} - \mathbf{x}^{(k)}$  is as small as possible.

We must be able to do so without knowing the ground truth  $\bar{\mathbf{x}}$ .

The decision must be made from available information, such as the  $k$ th iterate  $\mathbf{x}^{(k)}$  and/or its corresponding residual  $\boldsymbol{\rho}^{(k)} = \mathbf{b} - \mathbf{A} \mathbf{x}^{(k)}$ .

Such stopping rules are studied frequently by mathematicians, but they do not achieve the same attention in the tomographic reconstruction communities . . .

Our discussion of stopping rules is based on:

P. C. Hansen, J. S. Jørgensen, and P. W. Rasmussen, *Stopping rules for algebraic iterative reconstruction methods in computed tomography*; in 21st International Conference on Computational Science and Its Applications (ICCSA), IEEE (2021), pp. 60–70, doi 10.1109/ICCSA54496.2021.00019.

# On the Need for Stopping Rules

- The noise error  $\bar{\mathbf{x}}^{(k)} - \mathbf{x}^{(k)}$  often grows slowly with the number of iterations  $k$ . Hence the error history exhibits a *flat minimum*, and it is not crucial to stop at a very specific number of iterations.
- There are many applications for which the users have build a very *good intuition* of approximately how many iterations are needed to obtain a satisfactory reconstruction.
- When very many iterations are needed and the minimum is very flat, the iterations are often terminated with one's *patience runs out* – and hence one may not observe the semi-convergence effect.

Developing a robust stopping rule that works on many types of problems and for many kinds of data is difficult/impossible.

Here we give an overview of successful stopping rules, and then the user can try these methods on a given problem.

## A Bit of Statistical Notation

To make precise statements in this section, we need a small amount of statistical framework and notation.

The exact noise-free data  $\bar{\mathbf{b}}$  that correspond to the ground truth image:

$$\bar{\mathbf{b}} = \mathbf{A} \bar{\mathbf{x}}.$$

The elements of the noise vector  $\mathbf{e} \in \mathbb{R}^m$  are random variables, i.e., their values depend on a set of well-defined random events. The vector of expected values  $\mathcal{E}(\mathbf{e})$  and the covariance matrix  $\mathcal{V}(\mathbf{e})$  are defined as

$$\mathcal{E}(\mathbf{e}) = \begin{pmatrix} \mathcal{E}(e_1) \\ \mathcal{E}(e_2) \\ \vdots \end{pmatrix}, \quad \mathcal{V}(\mathbf{e}) = \mathcal{E}\left( (\mathbf{e} - \mathcal{E}(\mathbf{e})) (\mathbf{e} - \mathcal{E}(\mathbf{e}))^T \right).$$

We restrict our analysis to *white Gaussian noise* with zero mean:

$$\mathcal{E}(\mathbf{e}) = \mathbf{0}, \quad \mathcal{V}(\mathbf{e}) = \eta^2 \mathbf{I}, \quad \mathcal{E}(\|\mathbf{e}\|_2^2) = m\eta^2,$$

where  $\eta$  is the standard deviation of the noise.

# The Discrepancy Principle (DP)

DP: a model's output should fit the data “to the noise level.”

This translates into a stopping rule where we choose the number of iterations  $k$  such that the residual  $\mathbf{r}^{(k)} = \mathbf{b} - \mathbf{A} \mathbf{x}^{(k)}$  is “of the same size” as the noise vector  $\mathbf{e}$ :

$$\|\mathbf{r}^{(k)}\|_2^2 \approx m \eta^2.$$

We return to methods for estimating the standard deviation  $\eta$  from data.

Most authors include a constant  $\tau \geq 1$  such that the above condition takes the form

$$\|\mathbf{r}^{(k)}\|_2^2 \approx \tau m \eta^2.$$

The constant  $\tau$  can be useful when we have only a rough estimate of  $\eta$  and there is a risk that we take too many or too few iterations.

As we shall see, this stopping rule is quite dubious.

## How We Really Should Fit to the Noise Level I

To learn more about this principle, consider the TSVD solution

$$\mathbf{x}_k = \sum_{i=1}^k \frac{\mathbf{u}_i^T \mathbf{b}}{\sigma_i} \mathbf{v}_i \quad \text{with} \quad \mathbf{A} = \sum_{i=1}^n \mathbf{u}_i \sigma_i \mathbf{v}_i^T.$$

The corresponding residual takes the form

$$\boldsymbol{\rho}_k \equiv \mathbf{b} - \mathbf{A} \mathbf{x}_k = \sum_{i=k+1}^m (\mathbf{u}_i^T \mathbf{b}) \mathbf{u}_i = \mathbf{P}_k \mathbf{b} = \mathbf{P}_k \bar{\mathbf{b}} + \mathbf{P}_k \mathbf{e}.$$

where the projection matrix  $\mathbf{P}_k = \sum_{i=k+1}^m \mathbf{u}_i \mathbf{u}_i^T$  projects onto the subspace spanned by  $\mathbf{u}_{k+1}, \dots, \mathbf{u}_m$ .

The two components of  $\boldsymbol{\rho}_k$  take the form

$$\mathbf{P}_k \bar{\mathbf{b}} = \sum_{i=k+1}^m (\mathbf{u}_i^T \bar{\mathbf{b}}) \mathbf{u}_i \quad \text{and} \quad \mathbf{P}_k \mathbf{e} = \sum_{i=k+1}^m (\mathbf{u}_i^T \mathbf{e}) \mathbf{u}_i$$

## How We Really Should Fit to the Noise Level II

The noise-free right-hand side's SVD components  $\mathbf{u}_i^T \bar{\mathbf{b}}$  have a decaying magnitude as  $k$  increases. Hence  $\|\mathbf{P}_k \bar{\mathbf{b}}\|_2$ , which always decreases with  $k$ , will decrease quite fast because it's largest SVD components are extracted first (for the small values of  $k$ ).

The noise component's norm is given by

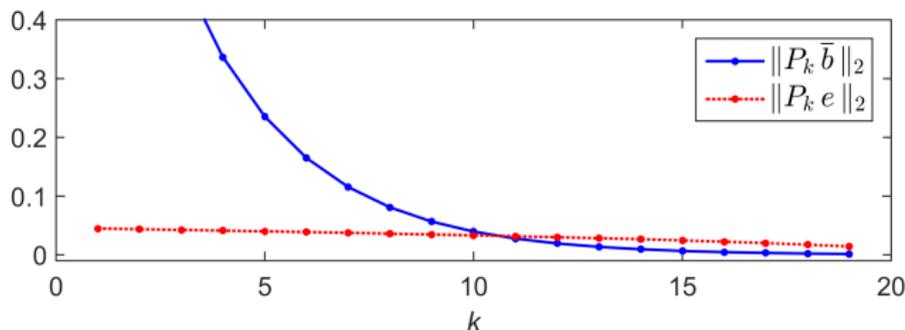
$$\|\mathbf{P}_k \mathbf{e}\|_2^2 = \sum_{i=k+1}^m (\mathbf{u}_i^T \mathbf{e})^2.$$

Since  $\mathbf{e}$  is zero-mean white Gaussian noise, the quantities  $\mathbf{u}_i^T \mathbf{e}$  also follow a Gaussian distribution with standard deviation  $\eta$ ; hence:

$$\mathcal{E}(\|\mathbf{P}_k \mathbf{e}\|_2^2) = \mathcal{E}\left(\sum_{i=k+1}^m (\mathbf{u}_i^T \mathbf{e})^2\right) = \sum_{i=k+1}^m \mathcal{E}((\mathbf{u}_i^T \mathbf{e})^2) = (m - k) \eta^2.$$

The factor  $m - k$  reflects the fact that the vector  $\mathbf{P}_k \mathbf{e}$  lies in a subspace of that dimension and thus has  $m - k$  degrees of freedom. The norm  $\|\mathbf{P}_k \mathbf{e}\|_2$  also decays with  $k$ , and compared to  $\|\mathbf{P}_k \bar{\mathbf{b}}\|_2$  it decays rather slowly.

## How We Really Should Fit to the Noise Level III



- $k$  too small: we have not captured enough SVD components;  $\mathbf{A} \mathbf{x}_k$  is not a good approximation the exact data  $\bar{\mathbf{b}}$ ,  $\mathbf{e}_k$  is dominated by  $P_k \bar{\mathbf{b}}$ , and  $\|P_k \bar{\mathbf{b}}\|_2$  is larger than  $\|P_k \mathbf{e}\|_2$ .
- $k$  “just about right”:  $\mathbf{A} \mathbf{x}_k$  approximates  $\bar{\mathbf{b}}$  as well as possible; the norm  $\|P_k \bar{\mathbf{b}}\|_2$  has now become smaller and is of the same size as  $\|P_k \mathbf{e}\|_2$ .
- $k$  too large: the residual  $\mathbf{e}_k$  is still dominated by the noise component  $P_k \mathbf{e}$ , and hence  $\|P_k \mathbf{e}\|_2$  dominates the residual norm.

Strategy: choose  $k$  such that  $\|P_k \bar{\mathbf{b}}\|_2 \approx \|P_k \mathbf{e}\|_2$ . But both are unknown, so in practise we should choose  $k$  such that  $\|\mathbf{e}_k\|_2^2 \approx (m - k) \eta^2$ .

## Formalization of The Heuristic

This heuristic reasoning was formalized by several authors; we summarize the main results our iterative methods. From the previous analysis it follows that we can write the  $k$ th Landweber iterate as

$$\mathbf{x}^{(k)} = \mathbf{A}_k^\# \mathbf{b} \quad \text{with} \quad \mathbf{A}_k^\# = \mathbf{V} \Phi^{(k)} \Sigma^{-1} \mathbf{U}^T.$$

The data predicted by the  $\mathbf{x}^{(k)}$  is given by  $\mathbf{b}_k = \mathbf{A} \mathbf{x}^{(k)} = \mathbf{A} \mathbf{A}_k^\# \mathbf{b}$ . The matrix  $\mathbf{A} \mathbf{A}_k^\#$  that transform the given, noisy data into this prediction is called the **influence matrix**.

With white Gaussian noise, it can be shown that at the optimal  $k$  we have

$$\mathcal{E}(\|\mathbf{e}^{(k)}\|_2^2) = \eta^2 (m - t_k), \quad t_k = \text{trace}(\mathbf{A} \mathbf{A}_k^\#) = \sum_{i=1}^n \phi_i^{(k)},$$

where  $\phi_i^{(k)}$  are the filter factors. The real number  $m - t_k$  is the effective (or equivalent) degrees of freedom in the residual.

For TSVD the filter factors are 0's and 1's, we simply have  $t_k = k$ .

# Stop Rule: Fit to Noise Level

## Stop Rule: Fit to Noise Level

Stop at the smallest  $k$  for which  $\|\mathbf{e}^{(k)}\|_2^2 \leq \eta^2 (m - t_k)$ .

For our iterative methods, this is particularly convenient because the residual norm  $\|\mathbf{e}^{(k)}\|_2$  decreases monotonically with  $k$ . To see this, we write the residual vector in terms of the SVD:

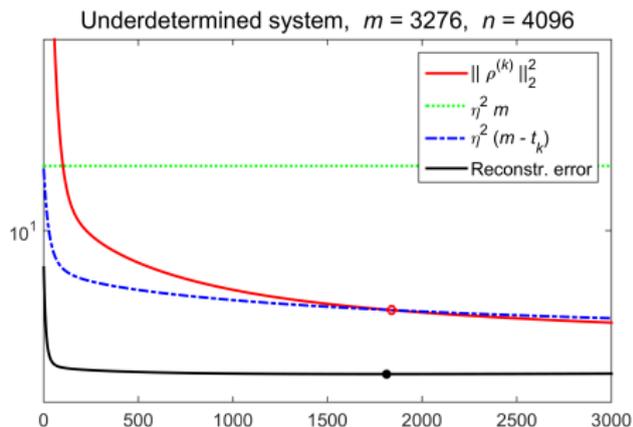
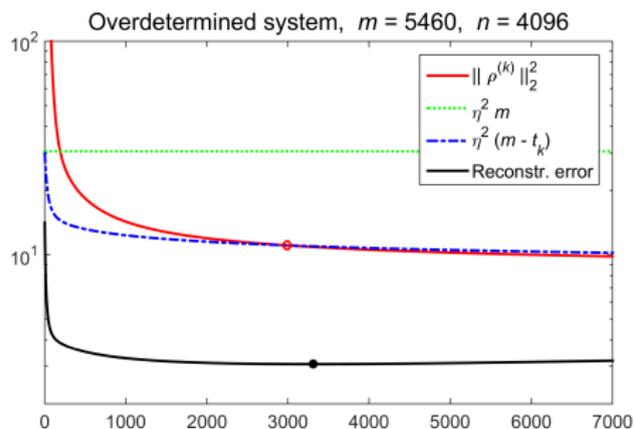
$$\mathbf{e}^{(k)} = \mathbf{b} - \mathbf{A} \mathbf{x}^{(k)} = \mathbf{U} (\mathbf{I} - \mathbf{\Phi}^{(k)}) \mathbf{U}^T \mathbf{b} = \mathbf{U} \begin{pmatrix} (1 - \phi_1^{(k)}) \mathbf{u}_1^T \mathbf{b} \\ (1 - \phi_2^{(k)}) \mathbf{u}_2^T \mathbf{b} \\ \vdots \end{pmatrix}.$$

Hence, for Landweber's method,

$$\|\mathbf{e}^{(k)}\|_2^2 = \sum_{i=1}^m (1 - \phi_i^{(k)})^2 (\mathbf{u}_i^T \mathbf{b})^2 = \sum_{i=1}^m (1 - \omega \sigma_i^2)^{2k} (\mathbf{u}_i^T \mathbf{b})^2.$$

Since  $\omega$  is always chosen such that  $|1 - \omega \sigma_i^2| < 1$  the factors  $(1 - \omega \sigma_i^2)^{2k}$  – and hence the squared residual norm – decrease monotonically with  $k$ .

# Illustration of Fit-to-Noise-Level Rule for Landweber



Parallel-beam example: image size =  $64 \times 64$ ; no. detector pixels = 91; projection angles =  $3^\circ, 6^\circ, 9^\circ, \dots, 180^\circ$  and  $8^\circ, 16^\circ, 24^\circ, \dots, 180^\circ$ .

Figures show reconstruction error  $\|\bar{\mathbf{x}} - \mathbf{x}^{(k)}\|_2$  and residual norm  $\|\rho^{(k)}\|_2$  versus  $k$ , together with threshold  $\eta\sqrt{m}$  and the function  $\eta\sqrt{m - t_k}$ .

We stop near the optimal number of iterations. DP stops much too early.

## Minimization of the Prediction Error

Instead of fitting to the noise level (as described above) we can find the number of iterations that minimizes the **prediction error**, i.e., the difference between the noise-free data  $\bar{\mathbf{b}} = \mathbf{A} \bar{\mathbf{x}}$  and the predicted data  $\mathbf{A} \mathbf{x}^{(k)}$ :

$$\min_k \|\mathbf{A} \mathbf{x}^{(k)} - \bar{\mathbf{b}}\|_2.$$

Statisticians refer to various measures of this difference as the **predictive risk**, and the resulting method for choosing  $k$  is often called the *unbiased predictive risk estimation (UPRE)* method.

Again we present the results specifically in the framework of iterative reconstruction methods.

## Derivation of the UPRE Rule

The expected squared norm of the prediction error (the risk) is

$$\mathcal{E}(\|\bar{\mathbf{b}} - \mathbf{A} \mathbf{x}^{(k)}\|_2^2) = \|(\mathbf{I} - \mathbf{A} \mathbf{A}_k^\#) \bar{\mathbf{b}}\|_2^2 + \eta^2 \text{trace}((\mathbf{A} \mathbf{A}_k^\#)^2)$$

while the expected squared norm of the residual can be written as

$$\begin{aligned} \mathcal{E}(\|\mathbf{b} - \mathbf{A} \mathbf{x}^{(k)}\|_2^2) &= \|(\mathbf{I} - \mathbf{A} \mathbf{A}_k^\#) \bar{\mathbf{b}}\|_2^2 + \\ &\quad \eta^2 \text{trace}((\mathbf{A} \mathbf{A}_k^\#)^2) - 2\eta^2 \text{trace}(\mathbf{A} \mathbf{A}_k^\#) + \eta^2 m. \end{aligned}$$

Combining these two equations we can eliminate one of the trace terms and arrive at the following expression for the risk:

$$\mathcal{E}(\|\bar{\mathbf{b}} - \mathbf{A} \mathbf{x}^{(k)}\|_2^2) = \mathcal{E}(\|\mathbf{b} - \mathbf{A} \mathbf{x}^{(k)}\|_2^2) + 2\eta^2 \text{trace}(\mathbf{A} \mathbf{A}_k^\#) - \eta^2 m.$$

Substituting the actual squared residual norm  $\|\mathbf{e}^{(k)}\|_2^2 = \|\mathbf{b} - \mathbf{A} \mathbf{x}^{(k)}\|_2^2$  for its expected value, we define the **UPRE risk** as a function of  $k$ :

$$U_k = \|\mathbf{e}^{(k)}\|_2^2 + 2\eta^2 t_k - \eta^2 m.$$

Minimizer of  $U_k \rightarrow$  approximation to a minimizer of the prediction error.

## Stop Rule: UPRE

Note that  $U_k$  may not have a unique minimizer, and we therefore choose the smallest  $k$  at which  $U_k$  has a local minimum.

### Stop Rule: UPRE

Find the smallest  $k$  that minimizes  $U_k = \|\mathbf{e}^{(k)}\|_2^2 + 2\eta^2 t_k - \eta^2 m$ .

This rule also depends on an estimate of the standard deviation  $\eta$  of the noise – which may or may not be a problem in practise.

We shall therefore describe an alternative method for minimization of the prediction error that does not depend on knowledge of  $\eta$ .

## Cross Validation

Assume that we remove the  $i$ th element  $b_i$  from the right-hand side (the noisy data), compute a reconstruction  $\mathbf{x}_{[i]}^{(k)}$ , and then use this vector to compute a prediction  $\hat{b}_i = \mathbf{r}_i^T \mathbf{x}_{[i]}^{(k)}$  of the missing data  $b_i$ .

The goal is then to choose the  $k$  that minimizes the following measure of all the prediction errors:

$$\hat{G}_k = \frac{1}{m} \sum_{i=1}^m (b_i - \hat{b}_i)^2 = \frac{1}{m} \sum_{i=1}^m \left( b_i - \mathbf{r}_i^T \mathbf{x}_{[i]}^{(k)} \right)^2.$$

We can avoid the vectors  $\mathbf{x}_{[i]}^{(k)}$  and write  $\hat{G}_k$  directly in terms of  $\mathbf{x}^{(k)}$ :

$$\hat{G}_k = \frac{1}{m} \sum_{i=1}^m \left( \frac{b_i - \mathbf{r}_i^T \mathbf{x}^{(k)}}{1 - \alpha_i^{(k)}} \right)^2,$$

where  $\alpha_i^{(k)}$  =  $i$ th diagonal element of the influence matrix  $\mathbf{A}\mathbf{A}^\#$  for  $\mathbf{x}^{(k)}$ .

## Generalized Cross Validation (GCV)

The minimizer of  $\widehat{G}_k$  depends on the particular ordering of the data. ☹

Generalized cross validation (GCV) circumvents this problem by replacing all  $\alpha_i^{(k)}$  with their average

$$\mu^{(k)} = \frac{1}{m} \sum_{i=1}^m \alpha_i^{(k)} = \frac{1}{m} \text{trace}(\mathbf{A}\mathbf{A}_k^\#) = \frac{t_k}{m},$$

leading to the modified measure

$$\widetilde{G}_k = \frac{1}{m} \frac{1}{(1 - \mu^{(k)})^2} \sum_{i=1}^m (b_i - \mathbf{r}_i^T \mathbf{x}^{(k)})^2 = \frac{\|\mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}\|_2^2}{m(1 - t_k/m)^2} = m \frac{\|\boldsymbol{\varrho}^{(k)}\|_2^2}{(m - t_k)^2}.$$

The minimizer of  $\widetilde{G}_k$  is, of course, independent of the factor  $m$  and hence we choose to define the GCV risk as a function of  $k$  as

$$G_k = \|\boldsymbol{\varrho}^{(k)}\|_2^2 / (m - t_k)^2.$$

## Stop Rule: GCV

We arrived at an  $\eta$ -free stopping rule:

### Stop Rule: GCV

Find the  $k$  that minimizes  $G_k = \|\mathbf{e}^{(k)}\|_2^2 / (m - t_k)^2$ .

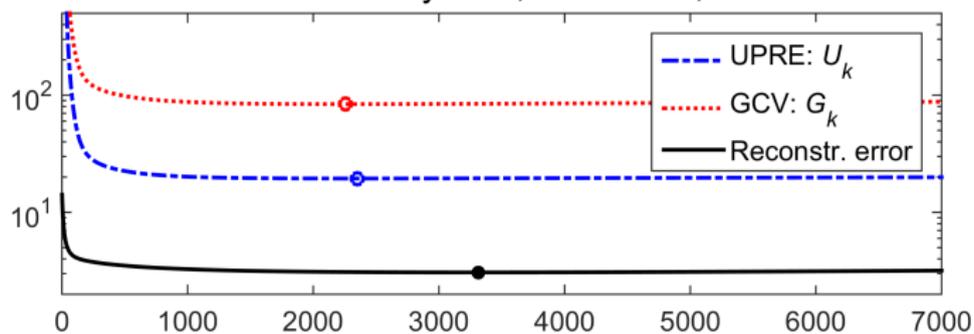
The value of  $k$  which minimizes  $G_k$  is also an estimate of the value that minimizes the prediction error. Specifically, if  $k_{\text{GCV}}$  minimizes the GCV risk  $G_k$  and  $k_{\text{PE}}$  minimizes the prediction error  $\|\bar{\mathbf{b}} - \mathbf{A} \mathbf{x}^{(k)}\|_2^2$ , then

$$\mathcal{E}(\|\bar{\mathbf{b}} - \mathbf{A} \mathbf{x}^{(k_{\text{GCV}})}\|_2^2) \rightarrow \mathcal{E}(\|\bar{\mathbf{b}} - \mathbf{A} \mathbf{x}^{(k_{\text{PE}})}\|_2^2) \quad \text{for } m \rightarrow \infty.$$

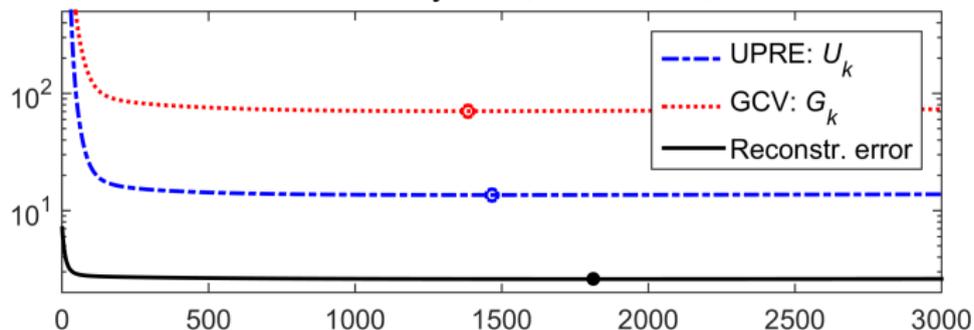
**A practical note.** UPRE and GCV need a few iterations too many detect a minimum of  $U_k$  and  $G_k$ . But the minimum of  $\|\bar{\mathbf{x}} - \mathbf{x}^{(k)}\|_2$  is usually very flat, and it hardly makes any difference if terminate the algorithm a few iterations after the minimum of  $U_k$  or  $G_k$ .

# Illustration of UPRE and GCV Rules for Landweber

Overdetermined system,  $m = 5460$ ,  $n = 4096$



Underdetermined system,  $m = 3276$ ,  $n = 4096$



## Estimation of the Trace Term I

We must estimate the trace term  $t_k$  efficiently – without the SVD of  $\mathbf{A}$  or the influence matrix  $\mathbf{A} \mathbf{A}_k^\#$ . The most common way to compute this estimate is via a Monte Carlo approach.

If  $\bar{\mathbf{w}} \in \mathbb{R}^m$  is a random vector with  $\bar{w}_i \sim \mathcal{N}(0, 1)$ , and  $\mathbf{S} \in \mathbb{R}^{m \times m}$  is a symmetric matrix, then  $\bar{\mathbf{w}}^T \mathbf{S} \bar{\mathbf{w}}$  is an unbiased estimate of  $\text{trace}(\mathbf{S})$ .

Hence  $\bar{t}_k^{\text{est}} = \bar{\mathbf{w}}^T \mathbf{A} \mathbf{A}_k^\# \bar{\mathbf{w}}$  is an unbiased estimator of  $t_k = \text{trace}(\mathbf{A} \mathbf{A}_k^\#)$ .

We need to compute the matrix-vector product  $\mathbf{A}_k^\# \bar{\mathbf{w}}$  efficiently.

Apply the algebraic iterative method to the system  $\mathbf{A} \bar{\boldsymbol{\xi}} = \bar{\mathbf{w}}$  which, after  $k$  iterations, produces the vector  $\bar{\boldsymbol{\xi}}^{(k)} = \mathbf{A}_k^\# \bar{\mathbf{w}}$ . The resulting estimate

$$\bar{t}_k^{\text{est}} = \bar{\mathbf{w}}^T \mathbf{A} \bar{\boldsymbol{\xi}}^{(k)},$$

is the standard [Monte Carlo trace estimate](#).

# Code for Trace Term Estimator I

Basic Landweber algorithm with MC trace estimator

$\bar{\mathbf{w}}$  = random  $m$ -vector for trace estimation

$\mathbf{x}^{(0)}$  = initial vector

$\bar{\boldsymbol{\xi}}^{(0)}$  = 0 initial zero vector for trace estimation

$\mathbf{z} = \mathbf{A}^T \bar{\mathbf{w}}$

for  $k = 0, 1, 2, \dots$

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

$$\bar{\boldsymbol{\xi}}^{(k+1)} = \bar{\boldsymbol{\xi}}^{(k)} + \omega \mathbf{A}^T (\bar{\mathbf{w}} - \mathbf{A} \bar{\boldsymbol{\xi}}^{(k)})$$

$$\bar{t}_{k+1}^{\text{est}} = \mathbf{z}^T \bar{\boldsymbol{\xi}}^{(k+1)} \text{ trace estimate}$$

stopping rule goes here

end

## Estimation of the Trace Term II

We can avoid the expensive multiplication with  $\mathbf{A}$  for each  $k$  with an approach that applies to unprojected iterative methods of the general form

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \omega \mathbf{A}^T \mathbf{B} (\mathbf{b} - \mathbf{A} \mathbf{x}^{(k)}),$$

where  $\mathbf{B}$  is a general  $m \times m$  matrix (it is not required to be symmetric). This includes Landweber's, Cimmino's and Kaczmarz's methods.

When we apply such a method with arbitrary nonzero starting vector  $\boldsymbol{\xi}^{(0)}$  to the system  $\mathbf{A} \boldsymbol{\xi} = \mathbf{0}$  then the iterates are  $\boldsymbol{\xi}^{(k)} = (\mathbf{I} - \omega \mathbf{A}^T \mathbf{B} \mathbf{A})^k \boldsymbol{\xi}^{(0)}$ .

Specifically, if we use a random starting vector  $\boldsymbol{\xi}^{(0)} = \mathbf{w} \in \mathbb{R}^n$  with elements  $w_i \sim \mathcal{N}(0, 1)$ , and if  $\boldsymbol{\xi}^{(k)}$  denote the corresponding iterates for  $\mathbf{A} \boldsymbol{\xi} = \mathbf{0}$ , then it can be shown that  $\mathbf{w}^T \boldsymbol{\xi}^{(k)}$  is an unbiased estimator of  $n - \text{trace}(\mathbf{A} \mathbf{A}_k^\#)$ .

This leads to the [alternative Monte Carlo trace estimate](#)

$$t_k^{\text{est}} = n - \mathbf{w}^T \boldsymbol{\xi}^{(k)}.$$

## Code for Trace Term Estimator II

Basic Landweber algorithm with alternative MC trace estimator

$\mathbf{w}$  = random  $n$ -vector

$\mathbf{x}^{(0)}$  = initial vector

$\boldsymbol{\xi}^{(0)}$  =  $\mathbf{w}$  initial vector for for trace estimation

for  $k = 0, 1, 2, \dots$

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

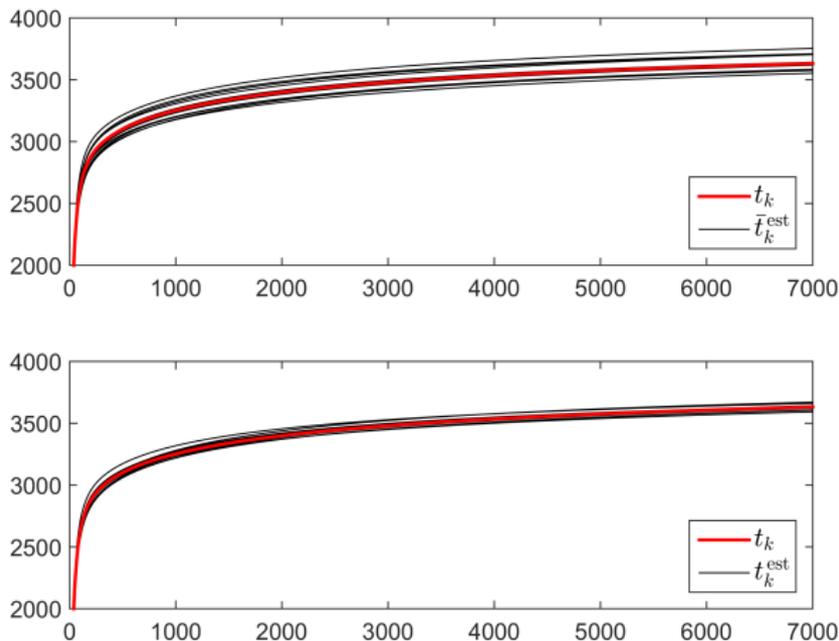
$$\boldsymbol{\xi}^{(k+1)} = \boldsymbol{\xi}^{(k)} + \omega \mathbf{A}^T (0 - \mathbf{A} \boldsymbol{\xi}^{(k)})$$

$$t_{k+1}^{\text{est}} = n - \mathbf{w}^T \boldsymbol{\xi}^{(k+1)} \text{ trace estimate}$$

stopping rule goes here

end

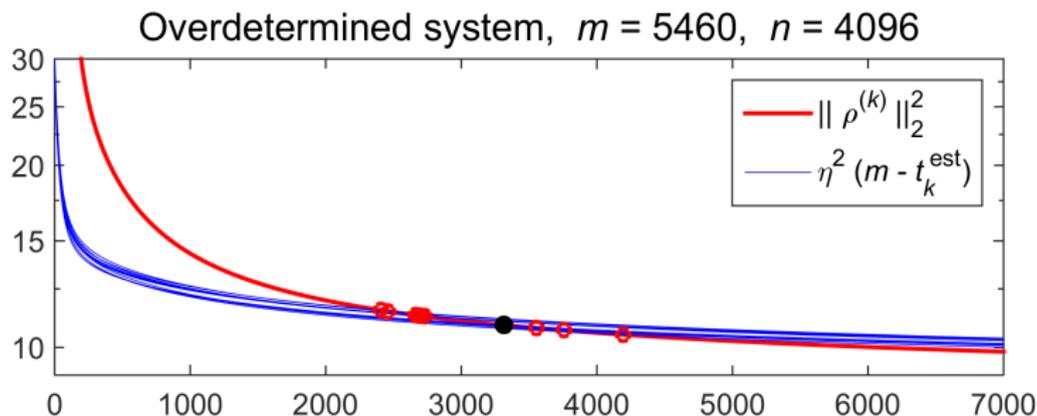
# Comparison of the Two Trace Estimates for Landweber



Thick **red line**: exact trace  $t_k$ .

Thin black lines: trace estimates for 10 different random vectors  $\bar{\mathbf{w}}$  and  $\mathbf{w}$ .

## Applying $t_k^{\text{est}}$ to the Fit-to-Noise-Level Stopping Rule



We used 10 different random vectors  $\mathbf{w}$ . The corresponding 10 intersections between  $\|\rho^{(k)}\|_2^2$  (thick **red line**) and  $\eta^2 (m - t_k^{\text{est}})$  (thin **blue lines**) are shown by the **red circles**.

The black dot  $\bullet$  shows the intersection with the exact  $\eta^2 (m - t_k)$ .

## Estimation of the Noise Level, In Our Framework

When we use the trace estimate  $t_k^{\text{est}}$  in the GCV stopping rule, then we seek a minimum for the approximate GCV risk given by

$$G_k^{\text{est}} = \|\boldsymbol{q}^{(k)}\|_2^2 / (m - t_k^{\text{est}})^2.$$

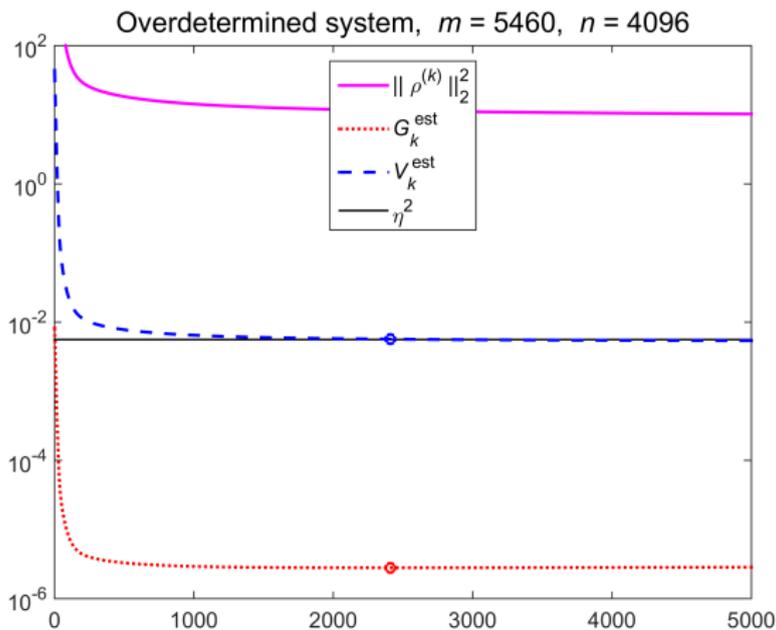
Now define the function

$$V_k^{\text{est}} = \|\boldsymbol{q}^{(k)}\|_2^2 / (m - t_k^{\text{est}}) = G_k^{\text{est}} (m - t_k^{\text{est}}).$$

According to the fit-to-noise-level stopping rule it follows that when we stop the iterations, the ratio  $\|\boldsymbol{q}^{(k)}\|_2^2 / (m - t_k)$  is approximately equal to the noise variance  $\eta^2$ .

Hence, if we terminate at iteration  $k = \hat{k}$  for which  $G_k^{\text{est}}$  is minimum, then the corresponding value  $V_{\hat{k}}^{\text{est}}$  is an inexpensive estimate of  $\eta^2$ .

# Illustration of $G_k^{\text{est}}$ and $V_{\hat{k}}^{\text{est}}$



The approximate GCV risk  $G_k^{\text{est}}$  has a minimum at  $\hat{k} = 2410$ .

Circles represent  $G_{\hat{k}}^{\text{est}}$  and  $V_{\hat{k}}^{\text{est}}$ , the latter being a good estimate of  $\eta^2$ .