

# Grouped Coordinate Descent Algorithms for Robust Edge-Preserving Image Restoration

Jeffrey A. Fessler

EECS Department  
The University of Michigan

SPIE'97: Image Reconstruction and Restoration

July 29, 1997

# OUTLINE

- Problem Description
- Huber Algorithm
- Optimization Transfer
- Convex Algorithm (ala Lange / De Pierro)
- Grouped Coordinate Descent (GCD) Algorithm
- Anecdotal Results
- Summary

# “LINEAR” INVERSE PROBLEM

$$\underline{y} = \mathbf{A}\underline{x} + \text{noise}$$

- $\underline{y}$ : noisy measurements (blurred image or sinogram)
- $\underline{x}$ : unknown object (true image)
- $\mathbf{A}$ : known system model  
(each column is a point response function)
- Errors in  $\mathbf{A}$  partially motivate robust methods

Goal: recover an estimate  $\hat{\underline{x}}$  of  $\underline{x}$  from  $\underline{y}$ .

# DATA-FIT COST FUNCTION

Want  $\hat{\underline{x}}$  to “fit the data,” i.e.  $\underline{y} \approx \mathbf{A}\hat{\underline{x}}$

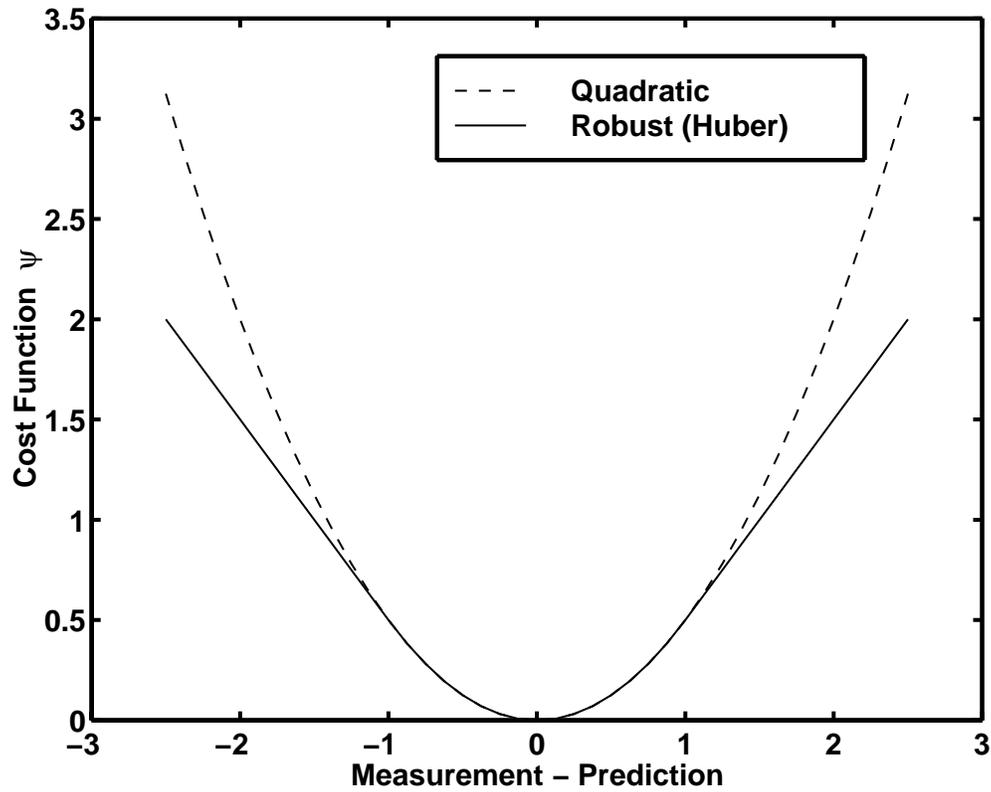
Natural cost function for independent measurement errors:

$$\Phi^{\text{data}}(\underline{x}) = \sum_{i=1}^{m_1} \psi_i^{\text{data}}([\underline{y} - \mathbf{A}\underline{x}]_i)$$

- $[\underline{y} - \mathbf{A}\underline{x}]_i = y_i - \sum_{j=1}^p a_{ij}x_j$
- $m_1$ : length of  $\underline{y}$
- $\psi_i$ : convex function.

Traditional choice:  $\psi_i(t) = t^2/2$ , which is appropriate for Gaussian noise, but is not robust to noise with heavy-tailed distributions.

# ROBUST DATA-FIT COST FUNCTION



Example - Huber function:

$$\psi(t) = \begin{cases} t^2/2, & |t| \leq \delta, \\ \delta|t| - \delta^2/2, & |t| > \delta \end{cases}$$

# ROBUST ESTIMATORS

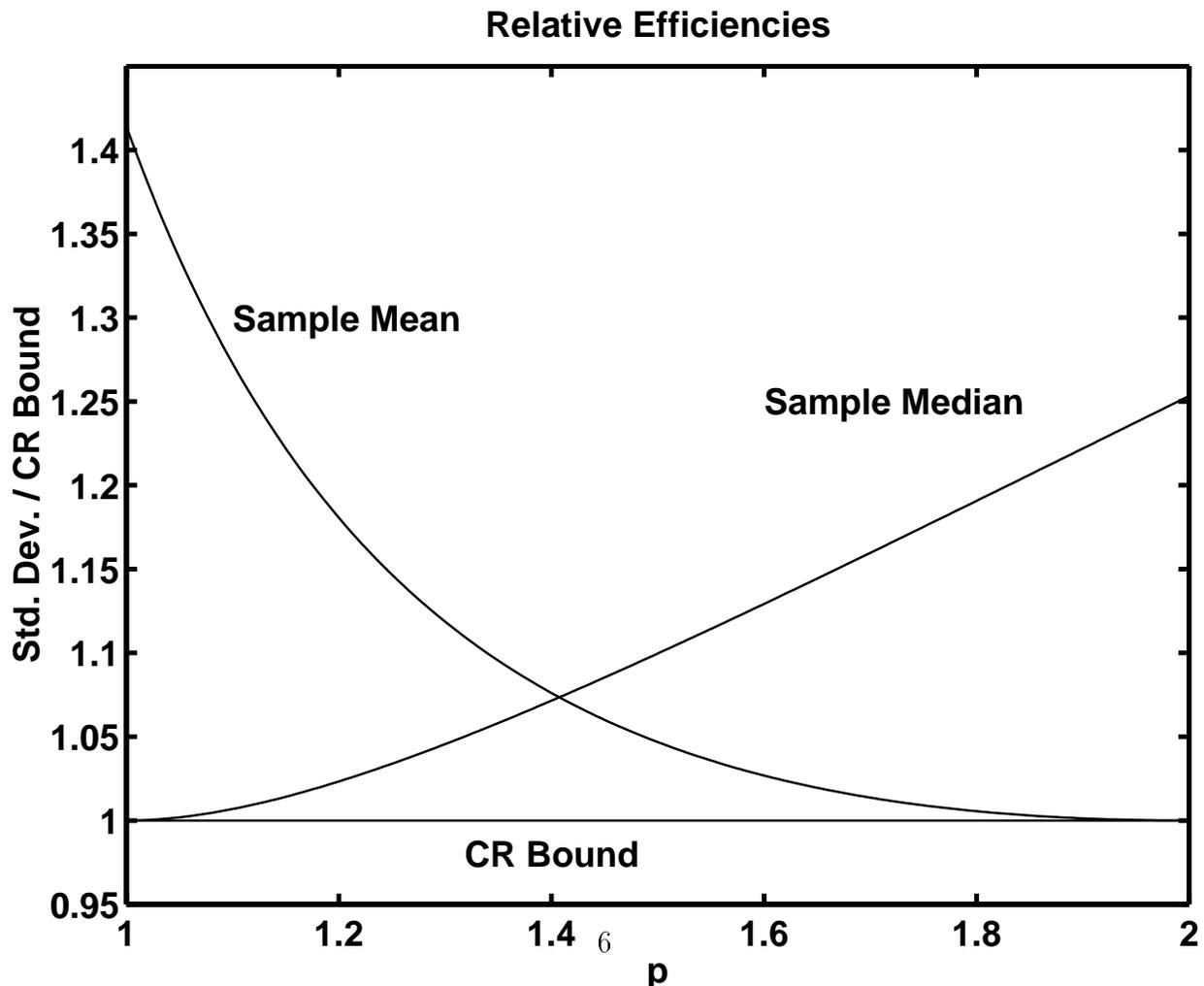
Generalized-Gaussian family of pdfs with unit variance:

$$f_X(x; \mu, p) = \frac{p}{2} \frac{1}{\Gamma(1/p)} \sqrt{r_p} \exp(-|x - \mu|^p r_p^{p/2}) \quad \text{where } r_p = \frac{\Gamma(3/p)}{\Gamma(1/p)}.$$

Asymptotic variance of the sample median estimator for  $\mu$  is:

$$\frac{1}{4n f^2(\mu)} = \frac{1}{n} \frac{\Gamma^2(1/p)}{p^2 r_p} \quad (\text{cf } 1/n \text{ for the sample mean}).$$

$$\text{CR bound for estimating } \mu: \sigma_{\hat{\mu}}^2 \geq \frac{1}{n} \frac{1}{p^2 r_p} \frac{\Gamma(1/p)}{\Gamma(2 - 1/p)}.$$



# REGULARIZATION

Minimizing  $\Phi^{\text{data}}$  is inadequate for ill-conditioned inverse problems.

Prior “knowledge” of piece-wise smoothness:

- $x_j - x_{j-1} \approx 0$  (piece-wise constant)
- $x_{j-1} - 2x_j + x_{j+1} \approx 0$  (piece-wise linear)
- $x_j \approx 0$  (support constraints)
- ... Combining:  $\mathbf{C}\underline{x} \approx \underline{z}$

Regularized cost function:  $\Phi(\underline{x}) = \Phi^{\text{data}}(\underline{x}) + \Phi^{\text{penalty}}(\underline{x}),$

$$\Phi^{\text{penalty}}(\underline{x}) = \sum_{i=1}^{m_2} \psi_i^{\text{penalty}}([\mathbf{C}\underline{x} - \underline{z}]_i)$$

## EXAMPLE: ROUGHNESS PENALTY (AKA GIBBS PRIOR)

$$\mathbf{D}_n = \begin{bmatrix} 1 & -1 & 0 & 0 & 0 \\ 0 & 1 & -1 & 0 & 0 \\ & & \ddots & \ddots & \\ 0 & 0 & 0 & 1 & -1 \end{bmatrix} \quad \mathbf{C} = \begin{bmatrix} \mathbf{I}_{n_y} \otimes \mathbf{D}_{n_x} \\ \mathbf{D}_{n_y} \otimes \mathbf{I}_{n_x} \end{bmatrix}$$

where  $\otimes$  denotes the Kronecker matrix product.

If  $\underline{z} = \underline{0}$  and  $\mathcal{N}_j$  is the four pixel neighborhood of pixel  $j$ , then

$$\Phi^{\text{penalty}}(\underline{x}) = \sum_j \sum_{k \in \mathcal{N}_j} \psi_{j,k}(x_j - x_k)$$

Conventional (Tikhonov-Miller) regularization:  $\psi(t) = t^2/2$ .  
(Gaussian prior)

For edge-preserving image recovery, need non-quadratic  $\psi(\cdot)$ ,  
such as Huber function.

# UNIFIED COST FUNCTION

$$\Phi(\underline{x}) = \sum_{i=1}^m \psi_i([\mathbf{B}\underline{x} - \underline{c}]_i)$$

Regularized edge-preserving cost function is a special case:

$$\Phi(\underline{x}) = \Phi^{\text{data}}(\underline{x}) + \Phi^{\text{penalty}}(\underline{x}), \quad \mathbf{B} = \begin{bmatrix} \mathbf{A} \\ \mathbf{C} \end{bmatrix}, \quad \underline{c} = \begin{bmatrix} \underline{y} \\ \underline{z} \end{bmatrix}$$

Optimization problem:

$$\hat{\underline{x}} = \arg \min_{\underline{x}} \Phi(\underline{x}) \quad \text{or} \quad \hat{\underline{x}} = \arg \min_{\underline{x} \geq \underline{0}} \Phi(\underline{x}).$$

# OPTIMIZATION

Simple in quadratic case where  $\psi_i(t) = t^2/2 \forall i$

$$\hat{\underline{x}} = (\mathbf{B}'\mathbf{B})^{-1}\mathbf{B}'\underline{c}$$

Good algorithms:

- Preconditioned conjugate gradients
- Coordinate descent (Gauss-Siedel)

Challenging for non-quadratic  $\psi_i$ 's

Very challenging for non-convex  $\psi_i$ 's

Proposition: algorithms tailored to structure of  $\Phi$  can outperform general purpose optimization methods.

but cannot solve it all...

# ASSUMPTIONS

$B$  has full column rank, so  $M > \mathbf{0} \Rightarrow B'MB > \mathbf{0}$   
(Easily achieved with sensible regularization design)

- $\psi$  is symmetric
- $\psi$  is everywhere differentiable (and therefore continuous)
- $\dot{\psi}(t) = d/dt \psi(t)$  is non-decreasing (and hence  $\psi$  is convex)
- $\omega_\psi(t) = \dot{\psi}(t)/t$  is non-increasing for  $t \geq 0$
- $\omega_\psi(0) = \lim_{t \rightarrow 0} \dot{\psi}(t)/t$  is finite and nonzero, i.e.  $0 < \omega_\psi(0) < \infty$

$\Phi$  has a unique minimizer  
(Easily ensured with perturbation of regularizer)

rules out entropy,  $|t|^p$

to understand  $\omega$ , look at...

# UNCONSTRAINED SOLUTION

$$\Phi(\underline{x}) = \sum_{i=1}^m \psi_i([\mathbf{B}\underline{x} - \underline{c}]_i)$$

Column gradient:

$$\nabla\Phi(\underline{x}) = \mathbf{B}'\boldsymbol{\Omega}(\underline{x})(\mathbf{B}\underline{x} - \underline{c}), \quad \nabla\Phi(\underline{x})|_{\underline{x}=\hat{\underline{x}}} = \mathbf{0}$$

$$\text{where } \boldsymbol{\Omega}(\underline{x}) = \text{diag}\{\omega_{\psi_i}([\mathbf{B}\underline{x} - \underline{c}]_i)\}$$

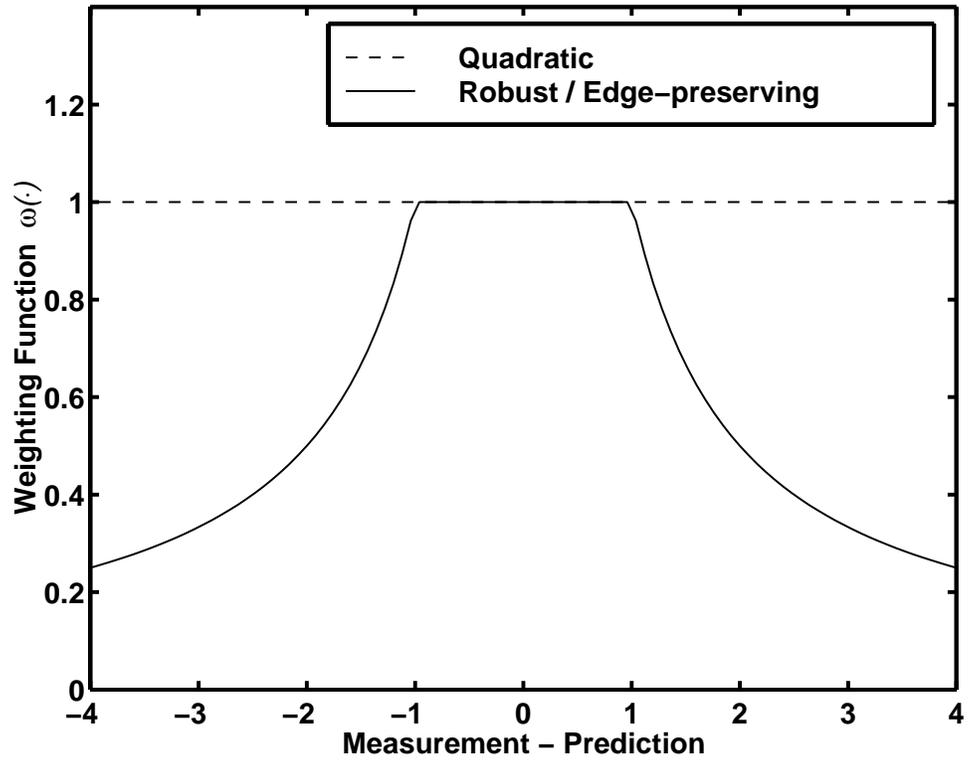
Unconstrained solution:

$$\begin{aligned}\hat{\underline{x}} &= [\mathbf{B}'\boldsymbol{\Omega}(\hat{\underline{x}})\mathbf{B}]^{-1}\mathbf{B}'\boldsymbol{\Omega}(\hat{\underline{x}})\underline{c} \\ &= \arg \min_{\underline{x}} \frac{1}{2}(\underline{c} - \mathbf{B}\underline{x})'\boldsymbol{\Omega}(\hat{\underline{x}})(\underline{c} - \mathbf{B}\underline{x})\end{aligned}$$

(ala WLS, but weights depend on estimate  $\hat{\underline{x}}$ , hence nonlinear)

Therefore need iterative algorithm...

# WEIGHTING FUNCTIONS $\omega_\psi$



# NEWTON-RAPHSON ALGORITHM

$$\underline{\mathbf{x}}^{n+1} = \underline{\mathbf{x}}^n - [\mathbf{B}'\mathbf{\Lambda}(\underline{\mathbf{x}}^n)\mathbf{B}]^{-1}\nabla\Phi(\underline{\mathbf{x}}^n)$$

where

$$\mathbf{\Lambda}(\underline{\mathbf{x}}^n) = \text{diag}\{\ddot{\psi}_i([\mathbf{B}\underline{\mathbf{x}} - \underline{\mathbf{c}}]_i)\}$$

Advantage:

- Super-linear convergence rate (if convergent)

Disadvantages:

- Requires twice-differentiable  $\psi_i$ 's
- Not guaranteed to converge
- Not guaranteed to monotonically decrease  $\Phi$
- Does not enforce nonnegativity constraint
- Impractical for image recovery due to matrix inverse

General purpose remedy: bound-constrained Quasi-Newton algorithms

## HUBER ALGORITHM (1981)

Recall  $\hat{\underline{x}} = [\mathbf{B}'\Omega(\hat{\underline{x}})\mathbf{B}]^{-1}\mathbf{B}'\Omega(\hat{\underline{x}})\underline{c} = \hat{\underline{x}} - [\mathbf{B}'\Omega(\hat{\underline{x}})\mathbf{B}]^{-1}\nabla\Phi(\hat{\underline{x}})$

Successive Substitutions:

$$\underline{x}^{n+1} = \underline{x}^n - [\mathbf{B}'\Omega(\underline{x}^n)\mathbf{B}]^{-1}\nabla\Phi(\underline{x}^n)$$

Advantages:

- Monotonically decreases  $\Phi$
- Converges globally to unique minimizer (not shown by Huber)

Disadvantages:

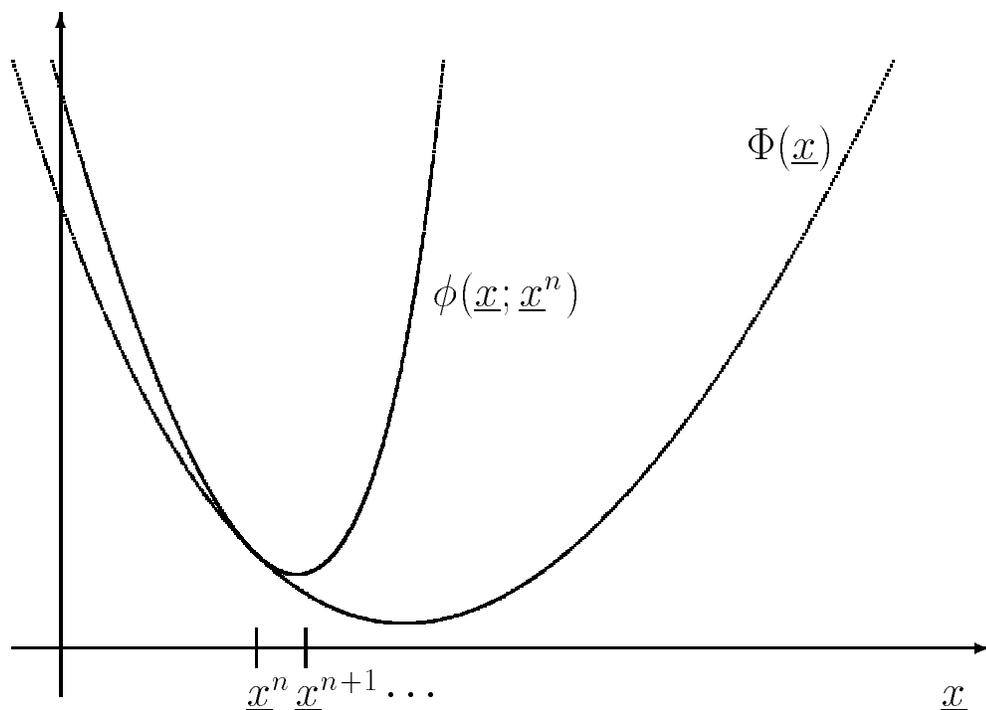
- Does not enforce nonnegativity constraint
- Impractical for image recovery due to matrix inverse

Successive substitutions is often not convergent. Why here?

# OPTIMIZATION TRANSFER

$$\underline{x}^{n+1} = \arg \min_{\underline{x}} \phi^{\text{Huber}}(\underline{x}; \underline{x}^n)$$

$$\phi^{\text{Huber}}(\underline{x}; \underline{x}^n) = \frac{1}{2}(\underline{c} - \mathbf{B}\underline{x})' \mathbf{\Omega}(\underline{x}^n) (\underline{c} - \mathbf{B}\underline{x})$$



Minimizing surrogate function  $\phi$  ensures a monotone decrease in  $\Phi$  if:

- $\phi(\underline{x}^n; \underline{x}^n) = \Phi(\underline{x}^n)$
- $\nabla_{\underline{x}} \phi(\underline{x}; \underline{x}^n)|_{\underline{x}=\underline{x}^n} = \nabla \Phi(\underline{x})|_{\underline{x}=\underline{x}^n}$
- $\Phi(\underline{x}) \leq \phi(\underline{x}; \underline{x}^n)$ .

These 3 (sufficient) conditions are satisfied by  $\phi^{\text{Huber}}$

# OPTIMIZATION TRANSFER IN 2D

# GENERALIZED HUBER ALGORITHM

$$\underline{x}^{n+1} = \underline{x}^n - \mathbf{M}_n^{-1} \nabla \Phi(\underline{x}^n)$$

where

$$\mathbf{M}_n \geq \mathbf{B}' \Omega(\underline{x}^n) \mathbf{B}$$

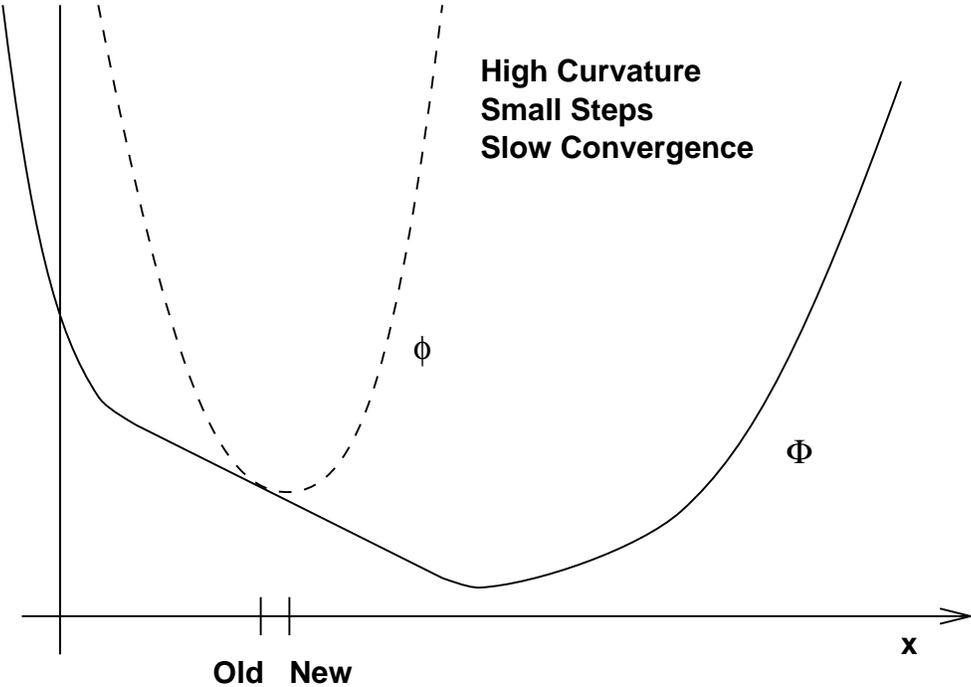
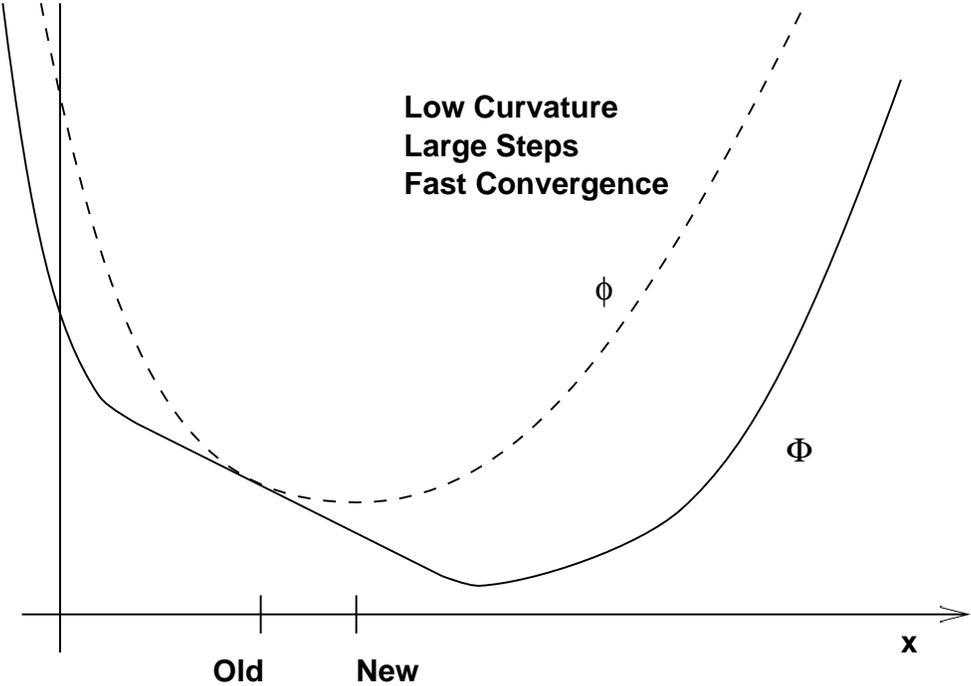
Advantages:

- Monotonically decreases  $\Phi$
- Converges globally to unique minimizer
- Can choose  $\mathbf{M}_n$  to be easily invertible, e.g. diagonal.  
(Or splitting matrices more generally)

Disadvantages:

- Does not enforce nonnegativity constraint
- Converges slower than Huber algorithm

# CONVERGENCE RATE



can we beat this tradeoff?

## USING THE STRUCTURE OF $\Phi$

De Pierro's decomposition (uses form of argument of  $\psi_i$ ):

$$\mathbf{B}\underline{x} - \underline{c} = \sum_{j=1}^p \alpha_{ij} \left[ \frac{b_{ij}}{\alpha_{ij}}(x_j - x_j^n) + \mathbf{B}\underline{x}^n - \underline{c} \right]$$

provided  $\alpha_{ij} \geq 0$  and  $\sum_{j=1}^p \alpha_{ij} = 1, \forall i$ .

The  $\alpha_{ij}$ 's are algorithm design factors.

Natural choice is  $\alpha_{ij} = |b_{ij}| / \sum_{k=1}^p |b_{ik}|$ .

By convexity of  $\psi_i$ :

$$\psi_i([\mathbf{B}\underline{x} - \underline{c}]_i) \leq \sum_{j=1}^p \alpha_{ij} \psi_i \left( \frac{b_{ij}}{\alpha_{ij}}(x_j - x_j^n) + \mathbf{B}\underline{x}^n - \underline{c} \right)$$

Construct surrogate function:

$$\Phi(\underline{x}) = \sum_{i=1}^m \psi_i([\mathbf{B}\underline{x} - \underline{c}]_i) \leq \phi^{\text{LDC}}(\underline{x}; \underline{x}^n)$$

$$\phi^{\text{LDC}}(\underline{x}; \underline{x}^n) = \sum_{j=1}^p \phi_j(x_j; \underline{x}^n),$$

$$\phi_j(x_j; \underline{x}^n) = \sum_{i=1}^m \alpha_{ij} \psi_i \left( \frac{b_{ij}}{\alpha_{ij}}(x_j - x_j^n) + \mathbf{B}\underline{x}^n - \underline{c} \right)$$

$\phi^{\text{LDC}}$  satisfies the 3 conditions for monotonicity

# LANGE / DE PIERRO CONVEX ALGORITHM

$$\underline{x}^{n+1} = \arg \min_{\underline{x}} \phi^{\text{LDC}}(\underline{x}; \underline{x}^n)$$

$$\begin{aligned} x_j^{n+1} &= \arg \min_{x_j \geq 0} \phi_j(x_j; \underline{x}^n) \\ &= \arg \min_{x_j \geq 0} \sum_{i=1}^m \alpha_{ij} \psi_i \left( \frac{b_{ij}}{\alpha_{ij}} (x_j - x_j^n) + \mathbf{B} \underline{x}^n - \underline{c} \right) \end{aligned}$$

Advantages:

- Monotonically decreases  $\Phi$
- Converges globally to unique minimizer
- No matrix inversion required
- Can enforce nonnegativity constraint
- Parallelizable (all pixels updated simultaneously)

Disadvantages:

- Requires subiteration for minimization  
Solution: use 1-D Huber algorithm
- Very slow convergence (ala EM algorithm)  
Solution: update only a subset of the pixels simultaneously

# GROUPED COORDINATE DESCENT ALGORITHM

Construct surrogate function using Lange / De Pierro convexity method but for only a (large) subset of the pixels.

## Pixel Groups (2x3)

1	5	3	1	5	3	1	5
4	2	6	4	2	6	4	2
1	5	3	1	5	3	1	5
4	2	6	4	2	6	4	2
1	5	3	1	5	3	1	5
4	2	6	4	2	6	4	2

**Pixels separated => decoupled => fast convergence**  
**Many pixels per subiteration => parallelizable**

Retains advantages of Convex Algorithm, but converges faster.

Disadvantages:

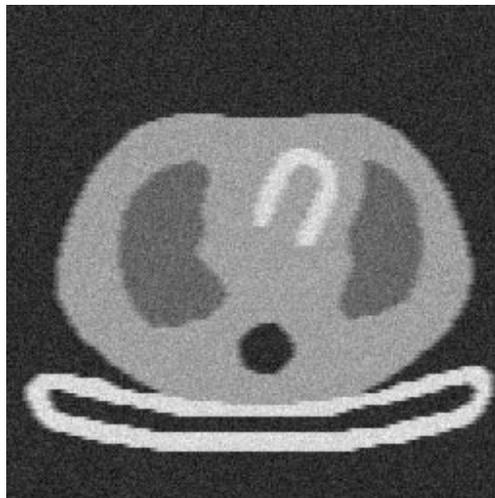
- Slightly less parallelizable.
- Slightly more complicated implementation
- Difficult to exploit structure of  $B$   
(e.g. FFTs for shift-invariant PSF, separable blur in PET)

# SIMULATION EXAMPLE

True object  $x$ :

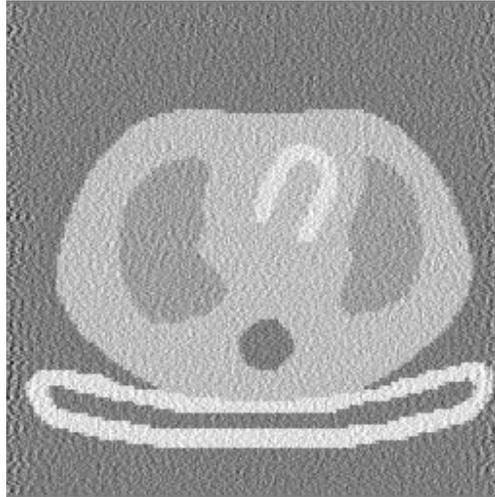


With 5 pixel horizontal motion blur and Gaussian noise,  $y$  is:



# RESTORED IMAGE

Wiener filter:

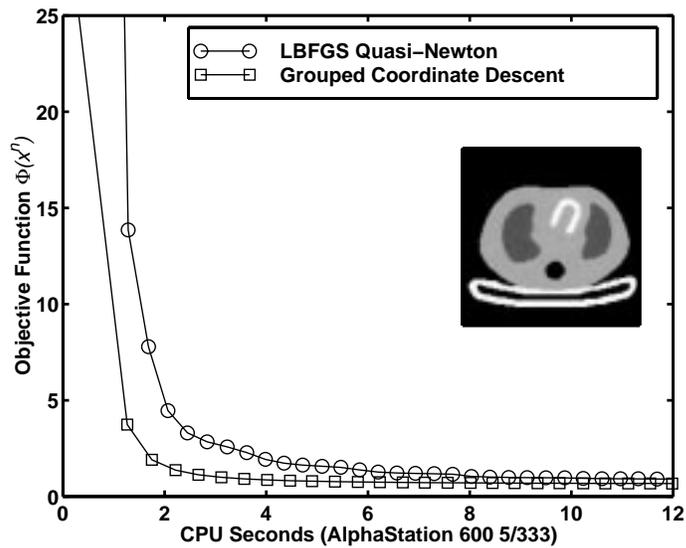
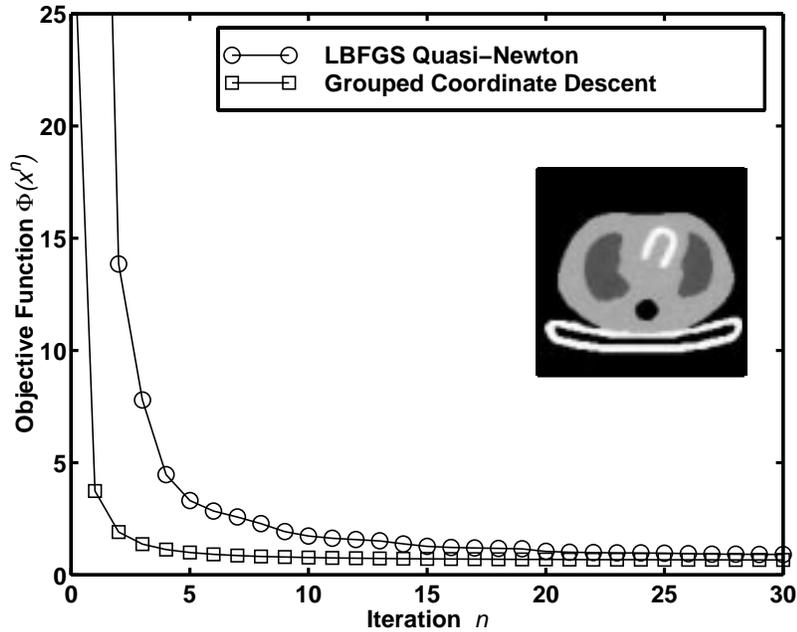


Edge-preserving restoration  $\hat{x}$ :



Huber function used for  $\psi_i$ 's for piece-wise smoothness.  
15 iterations of Grouped Coordinate Descent.

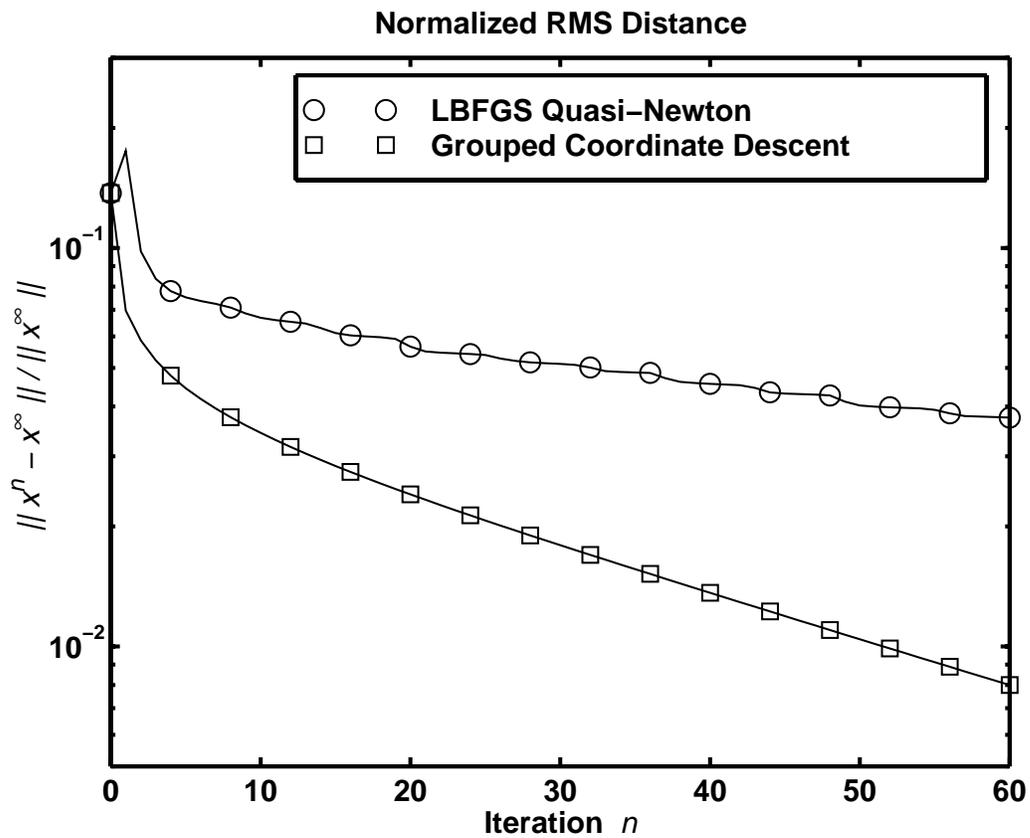
# CONVERGENCE RATES



LBFQS: Limited Memory Bound Constrained Quasi-Newton Method  
(R. Byrd, P. Lu, J. Nocedal, R. Schnabel, C. Zhu)

# NORMALIZED RMS DISTANCE

$$\frac{\|\underline{x}^n - \underline{x}^\infty\|}{\|\underline{x}^\infty\|}$$



$\underline{x}^\infty$ : 400 iterations of single-coordinate descent

(Thanks to Web Stayman for interfacing LBFGS with ASPIRE.)

# SUMMARY

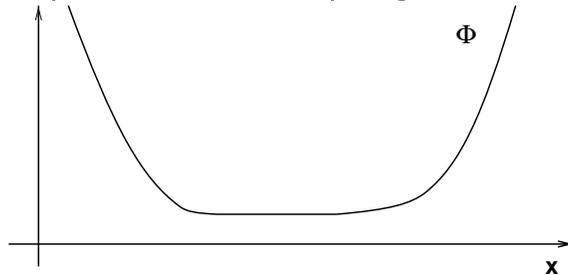
## Grouped Coordinate Descent Algorithm

- Accommodates non-quadratic cost function (for noise robustness and preserving edges)
- Monotonically decreases  $\Phi$
- Converges globally to unique minimizer
- Easily accommodates nonnegativity constraint
- Parallelizable
- Converges faster than a general-purpose optimization method

---

Future Work:

Extend convergence proofs for multiple global minimizers:



---

Slides and paper available from:

<http://www.eecs.umich.edu/~fessler/>