

MOCCA:
a primal/dual algorithm
for nonconvex composite functions
with applications to CT imaging

Rina Foygel Barber

Dept. of Statistics, University of Chicago

<http://www.stat.uchicago.edu/~rina/mocca.html>

Collaborators

- Algorithm & optimization work:
collaboration with Emil Sidky
- Application to CT:
collaboration with Emil Sidky, Taly Gilat-Schmidt, & Xiaochuan Pan



Emil Sidky
Dept. of Radiology
U. Chicago

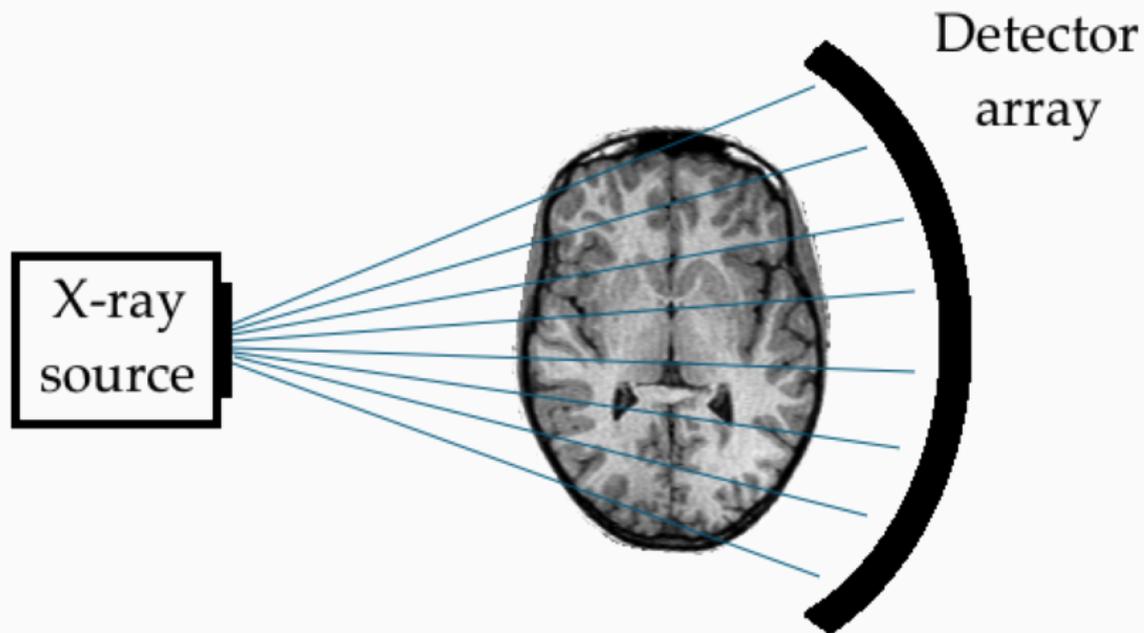


Xiaochuan Pan

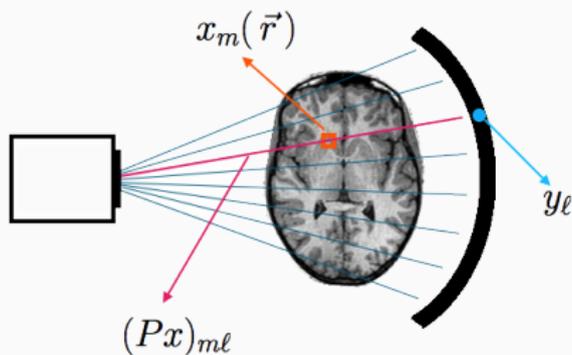


Taly Gilat-Schmidt
Dept. Biomedical Eng.
Marquette U.

Computed tomography (CT) imaging



Computed tomography (CT) imaging



- Measure: $y_\ell =$ number of photons detected along ray ℓ
- Want to estimate the materials at each point inside the object:

$$x_m(\vec{r}) = \text{density of material } m \text{ at location } \vec{r}$$

- Distribution of y is \approx determined by projections of x :

$$(Px)_{m\ell} = \text{amount of material } m \text{ along ray } \ell$$

Computed tomography (CT) imaging

If the X-ray beam is monochromatic,

for each ray ℓ the number of photons detected is

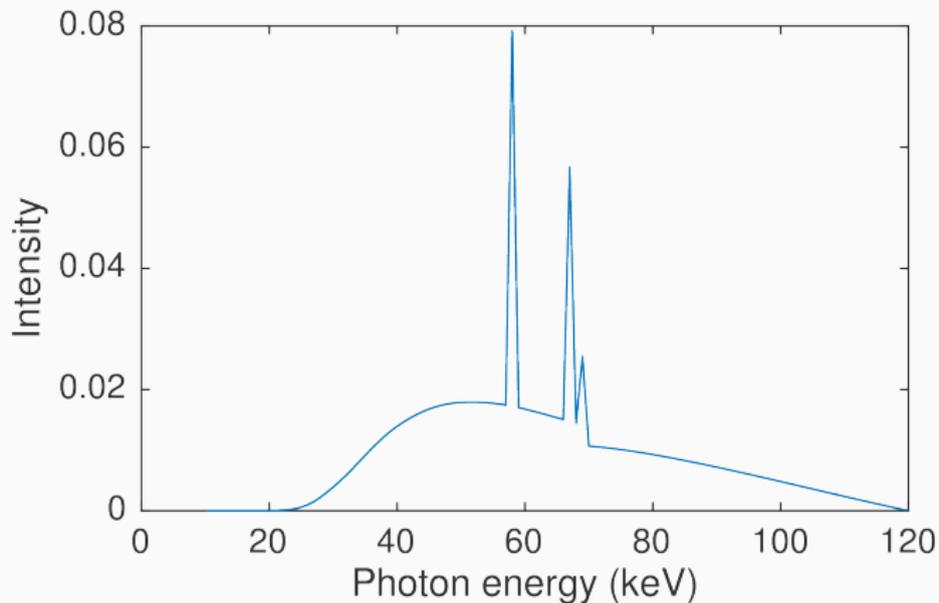
$$y_\ell \approx \text{Poisson} \left(I_{\text{total}} \cdot \exp \left\{ - \sum_m \mu_m \cdot \underbrace{(Px)_{m\ell}}_{\substack{\text{amount of material } m \\ \text{along ray } \ell}} \right\} \right)$$

μ_m = attenuation coefficient for material m

I_{total} = total intensity of X-ray spectrum / detector sensitivity

Computed tomography (CT) imaging

X-ray beam used in CT is polychromatic:



Computed tomography (CT) imaging

For polychromatic X-ray beam:

$$y_{\ell} \approx \text{Poisson} \left(I_{\text{total}} \int_E S(E) \cdot \exp \left\{ - \sum_m \underbrace{\mu_m(E)}_{\substack{\text{attenuation coefficient for} \\ \text{material } m \text{ at energy } E}} \cdot (Px)_{m\ell} \right\} dE \right)$$

$S(E)$ = distribution of X-ray spectrum intensity /
detector sensitivity across energies E

Computed tomography (CT) imaging

Existing algorithms for CT treat the measurements as a log linear function of the image:

$$\log(\mathbb{E}[y]) \approx \text{Linear function of } Px$$

- Filtered back projection (FBP) — used in clinical CT

Computed tomography (CT) imaging

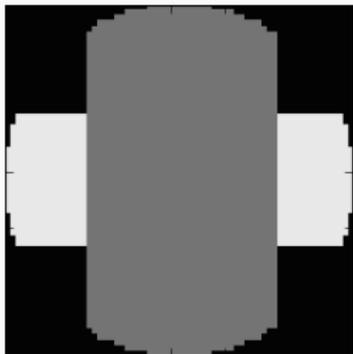
$$\begin{aligned} & \log \left(\mathbb{E} \left[\frac{y_\ell}{I_{\text{total}}} \right] \right) \\ &= \log \left(\int_E S(E) \cdot \exp \left\{ - \sum_m \mu_m(E) \cdot (Px)_{m\ell} \right\} dE \right) \end{aligned}$$

If we swap $\log(\cdot)$ with averaging:

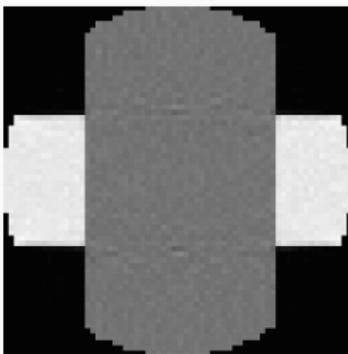
$$\approx - \sum_m \left[\int_E S(E) \cdot \mu_m(E) dE \right] \cdot (Px)_{m\ell}$$

Computed tomography (CT) imaging

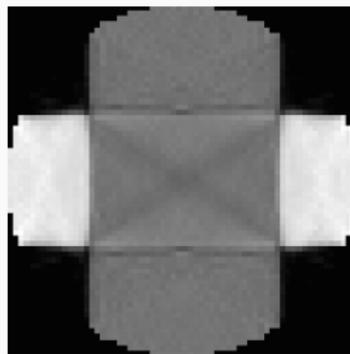
Ignoring the X-ray spectrum leads to beam hardening:



true object



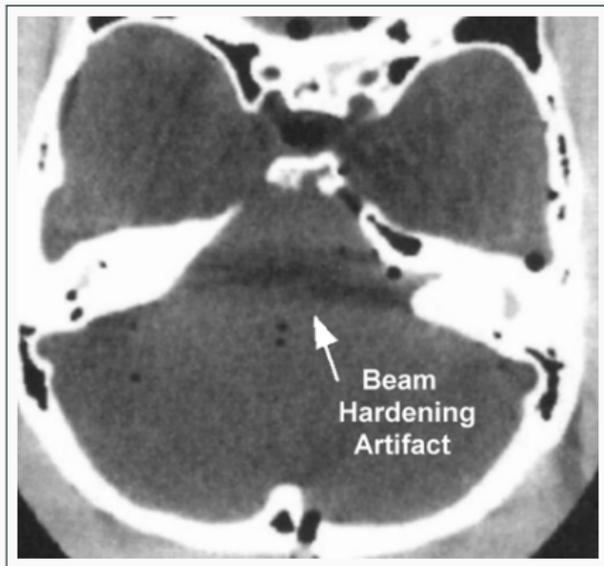
full Poisson model



log-linear
Poisson model

Computed tomography (CT) imaging

Beam hardening in practice:



Goldman, J. Nucl. Med. Technol., 2007

CT optimization problem

After discretization into pixels, want to minimize

$$\sum_{\text{rays } \ell} \mathcal{L}\left(y_{\ell}; \underbrace{\sum_{\text{energy } i} s_{li} \cdot \exp\left\{-\left(\mu^{\top} P x\right)_{li}\right\}}_{\text{Poisson negative log-likelihood}}\right) + \left(\text{Total variation constraints, etc}\right)$$

Vector x = discretized materials map

CT optimization problem

Spectral CT: photon detection is split
into multiple energy “windows” (bands):

$$\sum_{\substack{\text{windows } w \\ \text{rays } \ell}} \mathcal{L} \left(y_{w\ell}; \underbrace{\sum_{\text{energy } i} s_{wli} \cdot \exp \left\{ -(\mu^\top P x)_{li} \right\}}_{\text{Poisson negative log-likelihood}} \right) + \left(\text{Total variation constraints, etc} \right)$$

Vector x = discretized materials map

Optimization problem

General problem:

Want to minimize

$$F(Kx) + G(x)$$

where F and G might be nonconvex and/or nondifferentiable

Optimization problem: differentiable case

If F is differentiable & G has an easy proximal map:

- Proximal gradient descent:

$$\begin{cases} \tilde{x}_{t+1} = x_t - \frac{1}{\eta} K^\top \nabla F(Kx_t), \\ x_{t+1} = \arg \min \left\{ \frac{1}{2} \|x - \tilde{x}_{t+1}\|_2^2 + \frac{1}{\eta} G(x) \right\} \end{cases}$$

- Accelerated version: add an extrapolation step,

$$x_{t+1} \leftarrow x_{t+1} + \theta(x_{t+1} - x_t)$$

Convex: Beck & Teboulle 2009

Nonconvex: Loh & Wainwright 2013; Ochs et al 2014

Optimization problem: convex case

If F, G are convex:

ADMM (alternating direction method of multipliers)

- Rewrite optimization:

$$\min_{x,w} \max_u \left\{ F(w) + G(x) + \langle u, Kx - w \rangle + \frac{\sigma}{2} \|Kx - w\|_2^2 \right\}$$

- Alternate between minimizing over x and w , and updating u

Optimization problem: convex case

CP (Chambolle-Pock algorithm)

$$\text{Saddle point problem } \min_x \max_y \left\{ \underbrace{\langle Kx, y \rangle - \overbrace{F^*(y)}^{\text{Fenchel conjugate of } F}}_{F(Kx) = \max \text{ over } y} + G(x) \right\}$$

Iterate:

$$x_{t+1} = \arg \min \left\{ \langle Kx, y_t \rangle + G(x) + \frac{1}{2\tau} \|x - x_t\|_2^2 \right\}$$
$$y_{t+1} = \arg \max \left\{ \langle K\bar{x}_{t+1}, y \rangle - F^*(y) - \frac{1}{2\sigma} \|y - y_t\|_2^2 \right\}$$

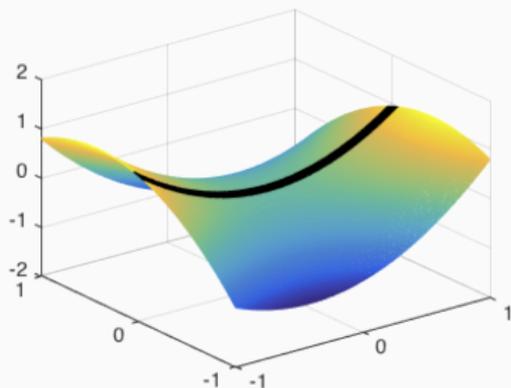
extrapolation $x_{t+1} + \theta(x_{t+1} - x_t)$

- Equivalent to ADMM with an added preconditioning step

Optimization problem: convex case

Can we run CP or ADMM if F & G are nonconvex?

- Example: $x \mapsto F(Kx) + G(x)$ is convex,
but F is strongly concave in some directions



- ADMM / CP may diverge immediately
- CP may converge to the wrong solution because $F^{**} \neq F$

MOCCA algorithm

Main idea:

1. Take local convex approximations to F and G
2. Take one step (or a few steps) of the CP algorithm
3. Repeat until convergence

MOCCA \approx majorization/minimization + primal/dual updates

Main question:

How should we construct the local convex approximations?

MOCCA algorithm

- Split F & G into convex + differentiable components:

$$F = F_{\text{cvx}} + F_{\text{diff}}, \quad G = G_{\text{cvx}} + G_{\text{diff}}$$

- Convex approximations at step t :

$$F_t(w) = F_{\text{cvx}}(w) + \left[F_{\text{diff}}(z_F^t) + \langle w - z_F^t, \nabla F_{\text{diff}}(z_F^t) \rangle \right]$$

$$G_t(x) = G_{\text{cvx}}(x) + \left[G_{\text{diff}}(z_G^t) + \langle x - z_G^t, \nabla G_{\text{diff}}(z_G^t) \rangle \right]$$

MOCCA algorithm

- How do we pick expansion points z_F^t and z_G^t ?

$$\underbrace{F(Kx)}_{\text{dual variable } y} + \underbrace{G(x)}_{\text{primal variable } x}$$

- $z_G^t =$ primal variable x_t
- $z_F^t =$ primal point that mirrors the dual variable y_t

Iterate:

$$x_{t+1} = \arg \min \left\{ \langle Kx, y_t \rangle + G_t(x) + \frac{1}{2\tau} \|x - x_t\|_2^2 \right\}$$

$$y_{t+1} = \arg \max \left\{ \langle K\bar{x}_{t+1}, y \rangle - F_t^*(y) - \frac{1}{2\sigma} \|y - y_t\|_2^2 \right\}$$

$$z_F^{t+1} = \frac{1}{\sigma}(y_t - y_{t+1}) + K\bar{x}_{t+1}, \quad z_G^{t+1} = x_{t+1}$$

- Step sizes σ, τ should satisfy $\sigma\tau \|K\|^2 < 1$.

(As in Chambolle & Pock 2011)

- Can use a preconditioning step to avoid computing $\|K\|$

(Pock & Chambolle 2011)

Case study: nonconvex total variation

The problem:

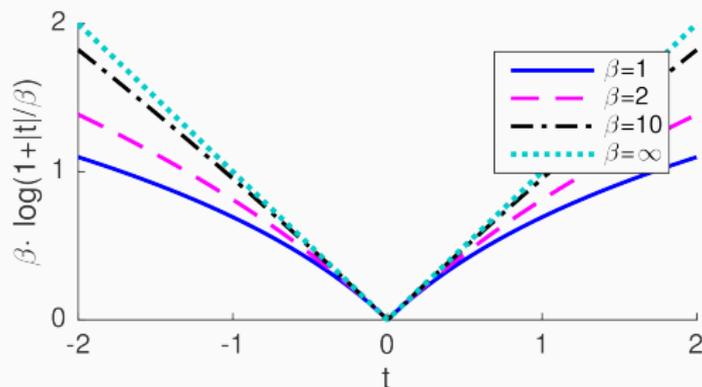
- True signal $x^* \in \mathbb{R}^d$ has total-variation sparsity
(nearby pixels often have identical values)
- Problem: minimize loss $\mathcal{L}(x)$ subject to sparsity in $\nabla_{2d}x$

2-dim. gradient operator
- Common approach: penalize $\|\nabla_{2d}x\|_1$
 \rightsquigarrow bias due to shrinkage on large gradient values

Case study: nonconvex total variation

Use a nonconvex TV penalty to reduce bias from shrinkage:

$$\log \text{TV}_\beta(x) = \sum_i \beta \cdot \log(1 + |(\nabla_{2d} x)_i|/\beta)$$



Equivalent to $\|x\|_{\text{TV}} = \|\nabla_{2d} x\|_1$ when $\beta = \infty$.

Parekh & Selesnick (2015)

Related to reweighted ℓ_1 sparsity, Candès et al (2008)

Case study: nonconvex total variation

Optimization problem for least squares loss:

$$\text{minimize } \frac{1}{2} \|b - Ax\|_2^2 + \nu \cdot \log \text{TV}_\beta(x)$$

$$\log \text{TV}_\beta(x) = \underbrace{\|\nabla_{2d} x\|_1}_{\text{convex}} + \underbrace{\left[\beta \log(1 + |\nabla_{2d} x|/\beta) - \|\nabla_{2d} x\|_1 \right]}_{h(\nabla_{2d} x) = \text{differentiable}}$$

Define:

$$\begin{aligned} F_{\text{cvx}}(w) &= \nu \cdot \|w\|_1 & G_{\text{cvx}}(x) &= \frac{1}{2} \|b - Ax\|_2^2 \\ F_{\text{diff}}(w) &= \nu \cdot h(w) & G_{\text{diff}}(x) &\equiv 0 \end{aligned}$$

Case study: nonconvex total variation

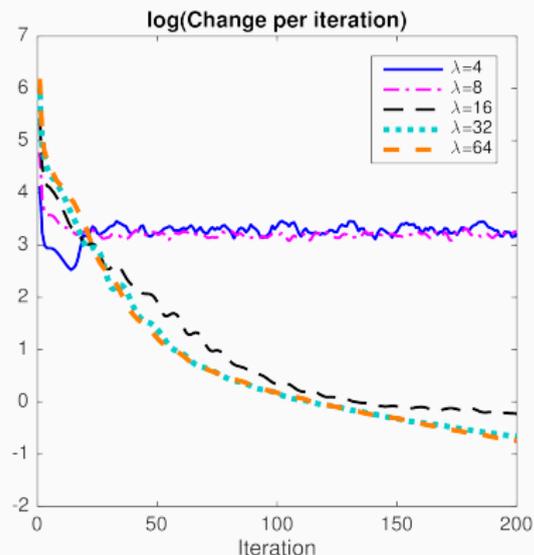
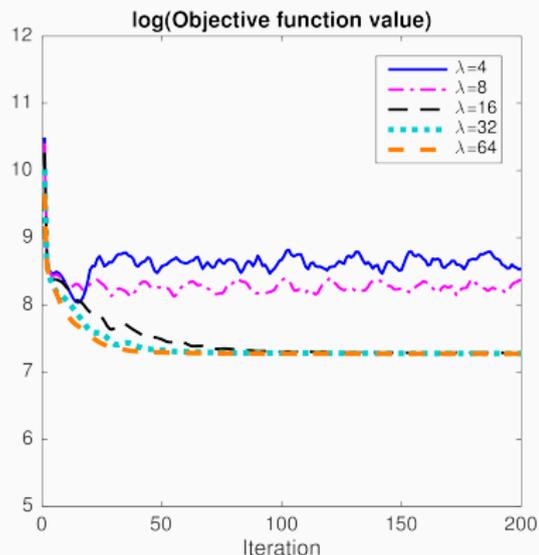
MOCCA for least squares + nonconvex TV

$$x_{t+1} = (\mathbf{I} + \tau A^\top A)^{-1} (x_t + \tau A^\top b - \tau \nabla_{2d}^\top y_t)$$

$$y_{t+1} = \text{Truncate}_\nu (y_t + \sigma \nabla_{2d} \bar{x}_{t+1} - \lambda \nabla h(z_F^t)) + \lambda \nabla h(z_F^t)$$

$$z_F^{t+1} = \frac{1}{\sigma} (y_t - y_{t+1}) + K \bar{x}_{t+1}$$

Case study: nonconvex total variation

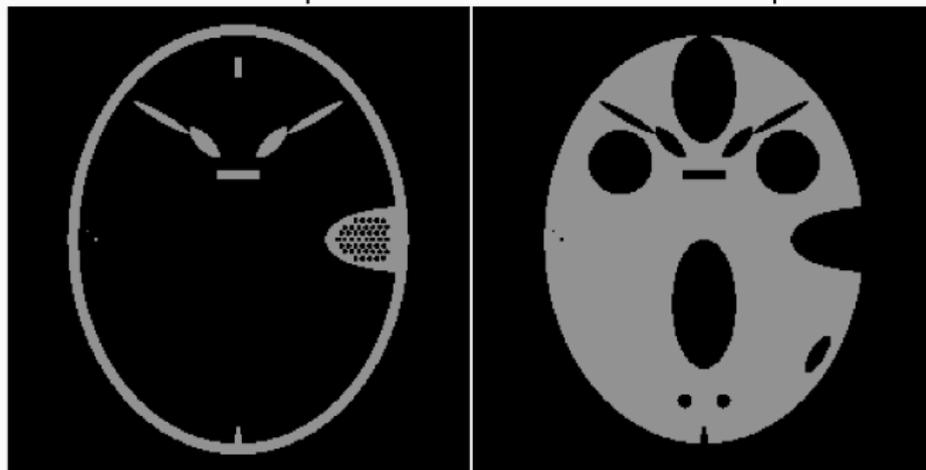


Problem size: $x \in \mathbb{R}^{25 \times 25}$ with block structure; 200 measurements

Tuning parameter λ : $\sigma = \frac{\lambda}{2}$, $\tau = \frac{1}{2\lambda}$

Application to spectral CT

Simulated CT measurements from object with 2 materials:



Bone

Brain

FORBILD head phantom (Lauritsch & Bruder 2001)

Application to spectral CT

Minimize:

$$\sum_{\substack{\text{windows } w \\ \text{rays } \ell}} \mathcal{L} \left(y_{w\ell}; \underbrace{\sum_{\text{energy } i} s_{w\ell i} \cdot \exp \left\{ -(\mu^\top P x)_{\ell i} \right\}}_{\text{Poisson negative log-likelihood}} \right) + \left(\text{Total variation constraints, etc} \right)$$

Application to spectral CT

Minimize:

$$\underbrace{\mathcal{L}(\mu^\top P \cdot x)}_{\text{Poisson negative log-likelihood}} + \underbrace{\delta \left(\begin{array}{l} \|x_{\text{bone}}\|_{\text{TV}} \leq \gamma_{\text{bone}} \\ \& \\ \|x_{\text{brain}}\|_{\text{TV}} \leq \gamma_{\text{brain}} \end{array} \right)}_{\text{convex indicator function}}$$

Application to spectral CT

MOCCA setup: minimize $F(Kx) + G(x)$

$$w = \begin{pmatrix} w_1 \\ w_2 \\ w_3 \end{pmatrix} = \begin{pmatrix} \mu^\top P \\ \nabla_{\text{bone}} \\ \nabla_{\text{brain}} \end{pmatrix} \cdot x = Kx$$

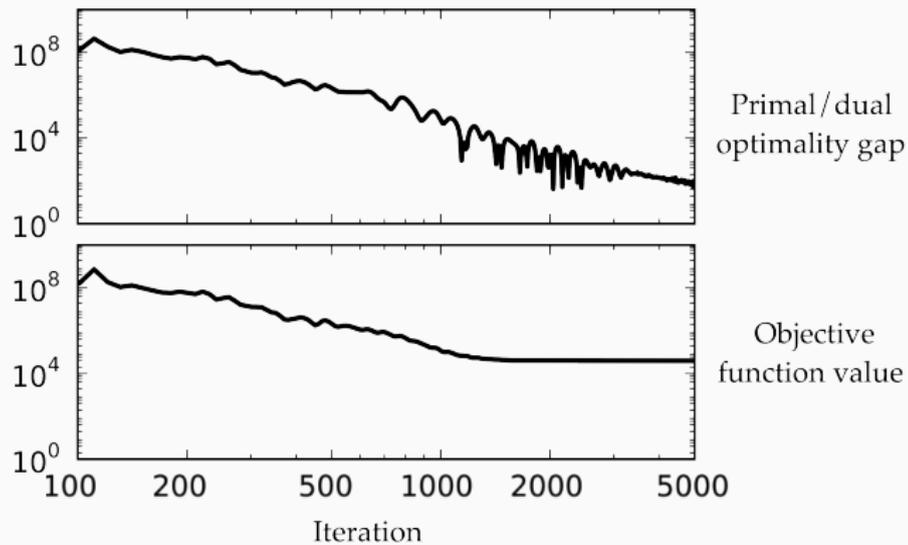
$$\left\{ \begin{array}{l} F(w) = \begin{pmatrix} \text{local convex/concave} \\ \text{quadratic approx. to } \mathcal{L}(w_1) \end{pmatrix} + \delta \begin{pmatrix} \|w_2\|_1 \leq \gamma_{\text{bone}} \\ \& \\ \|w_3\|_1 \leq \gamma_{\text{brain}} \end{pmatrix} \\ G(x) \equiv 0 \end{array} \right.$$

Application to spectral CT

Algorithm:

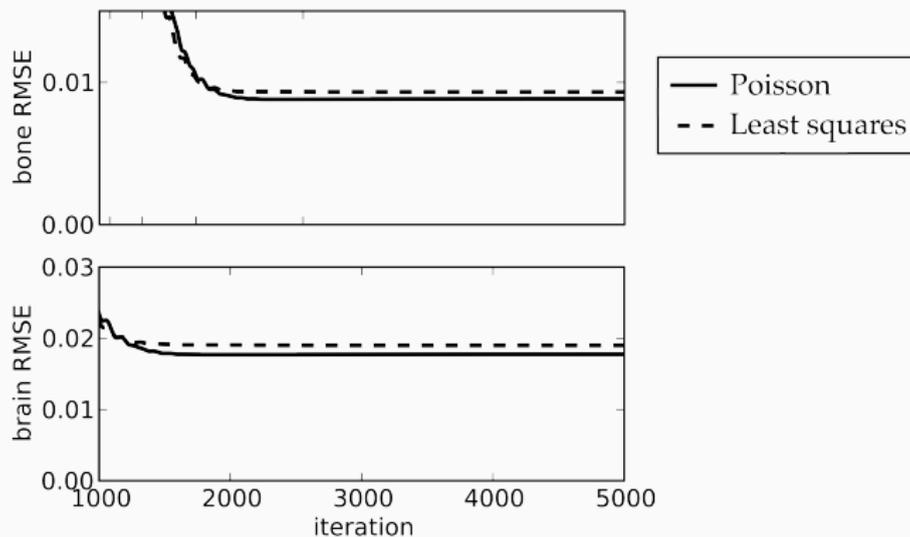
1. Take one step of the MOCCA algorithm
2. Update local convex/concave quadratic approximation to $\mathcal{L}(\cdot)$
3. Update step sizes
4. Repeat until convergence

Application to spectral CT



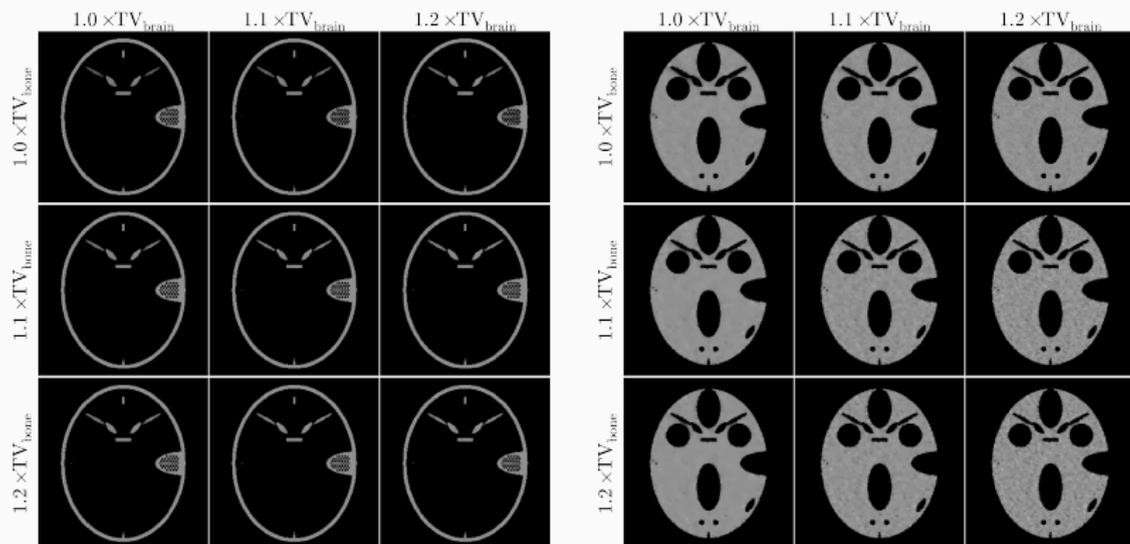
Application to spectral CT

Using the Poisson likelihood vs. a least squares loss:



Application to spectral CT

How critical is the choice of TV constraints γ_{bone} & γ_{brain} ?



Theoretical results

Question 1:

If MOCCA converges, has it converged to the right solution?

Theorem 1: critical points

If the MOCCA algorithm converges with

$$(x_t, y_t, z_t) \rightarrow (\hat{x}, \hat{y}, \hat{z})$$

then \hat{x} is a critical point of the optimization problem,

$$0 \in K^\top \partial F_{\text{cvx}}(K\hat{x}) + K^\top \nabla F_{\text{diff}}(K\hat{x}) + \partial G_{\text{cvx}}(\hat{x}) + \nabla G_{\text{diff}}(\hat{x})$$

Theoretical results

Question 2:

Is MOCCA guaranteed to converge (& at what rate)?

Theoretical results

Stable MOCCA algorithm (with “inner loop”)

At stage t ,

1. Run the “inner loop”: fixing expansion points (z_F^t, z_G^t) ,
update (x, y) variables L_{t+1} times
2. Update (x, y) variables by averaging over stage t :

$$(x_{t+1}, y_{t+1}) = \frac{1}{L_{t+1}} \sum_{\ell=1}^{L_{t+1}} (x_{t+1;\ell}, y_{t+1;\ell})$$

3. Update expansion points by averaging over stage t :

$$\begin{cases} z_F^{t+1} = \frac{1}{L_{t+1}} \sum_{\ell=1}^{L_{t+1}} \frac{1}{\sigma} (y_{t+1;\ell-1} - y_{t+1;\ell}) + K \bar{x}_{t+1;\ell} \\ z_G^{t+1} = \frac{1}{L_{t+1}} \sum_{\ell=1}^{L_{t+1}} x_{t+1;\ell} \end{cases}$$

Theoretical results

Background—restricted strong convexity:

- Definition: a loss function $\mathcal{L}(x)$ satisfies RSC if

$$\langle x - x', \partial\mathcal{L}(x) - \partial\mathcal{L}(x') \rangle \gtrsim \|x - x'\|_2^2 - \frac{\log(d)}{n} \|x - x'\|_1^2$$

- Convex: accurate recovery of sparse/structured signals in high-dimensional statistics

Negahban et al 2009

- Nonconvex: local minima guaranteed to be near global min for (differentiable loss) + (sparsity penalty)

Loh & Wainwright 2013

Theoretical results

Restricted convexity/smoothness assumptions for MOCCA:

- F_{cvx} is Λ_F -convex and F_{diff} is Θ_F -smooth
- G_{cvx} is Λ_G -convex and G_{diff} is Θ_G -smooth
- The overall optimization problem is nearly convex:

$$\underbrace{(Kx)^\top (\Lambda_F - \Theta_F)(Kx)}_{\text{Convexity of F}} + \underbrace{x^\top (\Lambda_G - \Theta_G)x}_{\text{Convexity of G}} \succeq C_{\text{cvx}} \|x\|_2^2 - \tau^2 \|x\|_{\text{restrict}}^2 \cdot$$

ℓ_1 norm / any
structured norm



- Optimization is over bounded region $\{x : \|x\|_{\text{restrict}} \leq R\}$

Theoretical results

Theorem 2: convergence guarantee

For the stable form of the MOCCA algorithm with $L_t \sim C^t$,

$$\|x_t - x^*\|_2 \lesssim C^{-t/2} + \tau R,$$

for any critical point x^* with $\|x^*\|_{\text{restrict}} \leq R$.

Number of iterations to calculate x_t is $L_1 + \dots + L_t \sim C^t$

$$\rightsquigarrow \|x_t - x^*\|_2 \sim \frac{1}{\sqrt{(\text{computational cost})}} + \tau R$$

Theoretical results

Main ingredient: contraction property

Consider two convex approximations:

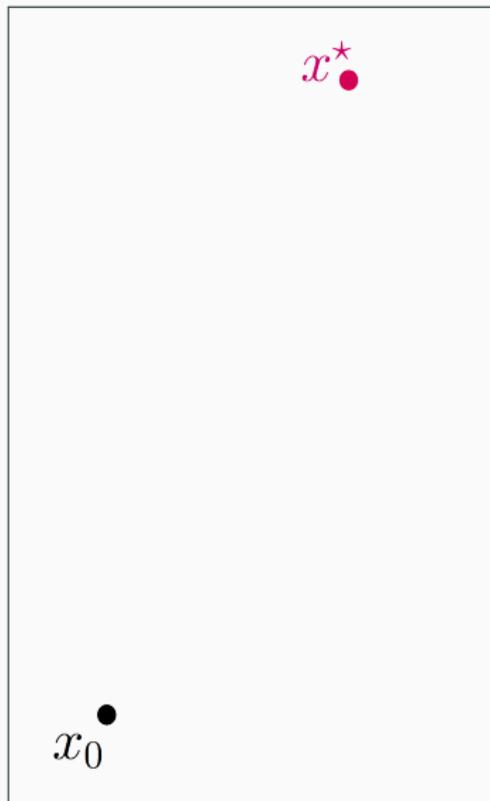
$$\left\{ \begin{array}{l} F_z(Kx) + G_z(x) \\ F_{z'}(Kx) + G_{z'}(x) \end{array} \right. \quad \text{with minimizers} \quad \left. \begin{array}{l} x_z^* \\ x_{z'}^* \end{array} \right\}$$

Then

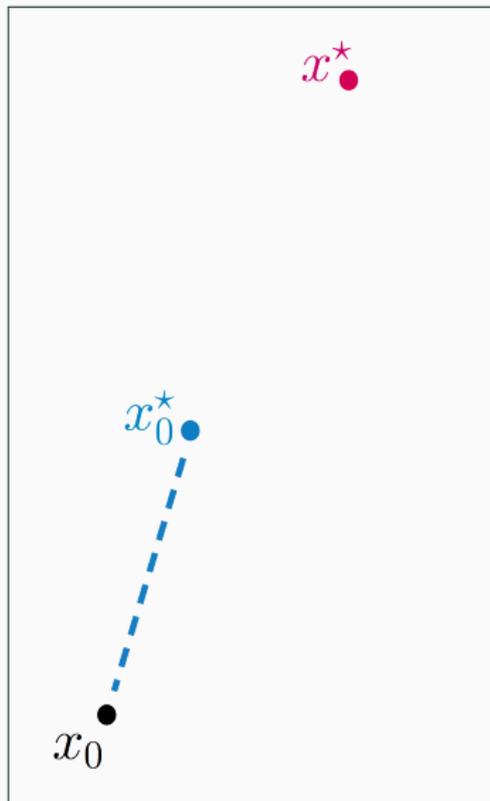
$$\left\| \begin{pmatrix} x_z^* - x_{z'}^* \\ Kx_z^* - Kx_{z'}^* \end{pmatrix} \right\| \leq (1 - \epsilon) \left\| \begin{pmatrix} z_G - z'_G \\ z_F - z'_F \end{pmatrix} \right\| + C \cdot \tau R$$

for some $\epsilon > 0$ and $C < \infty$.

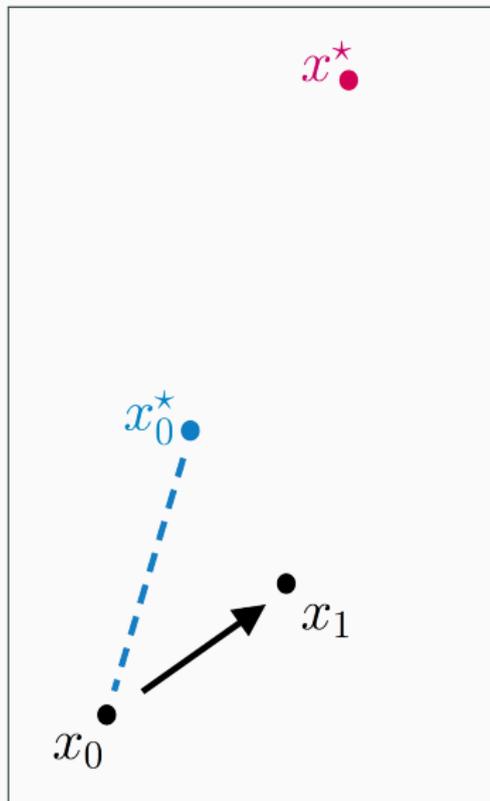
Theoretical results



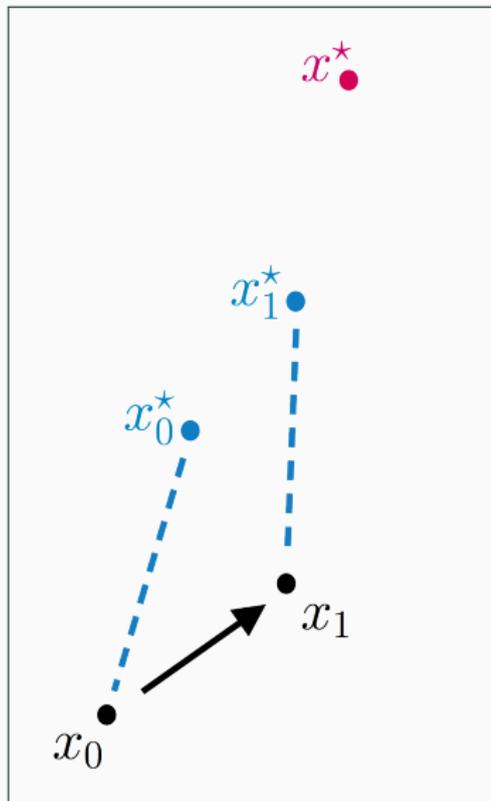
Theoretical results



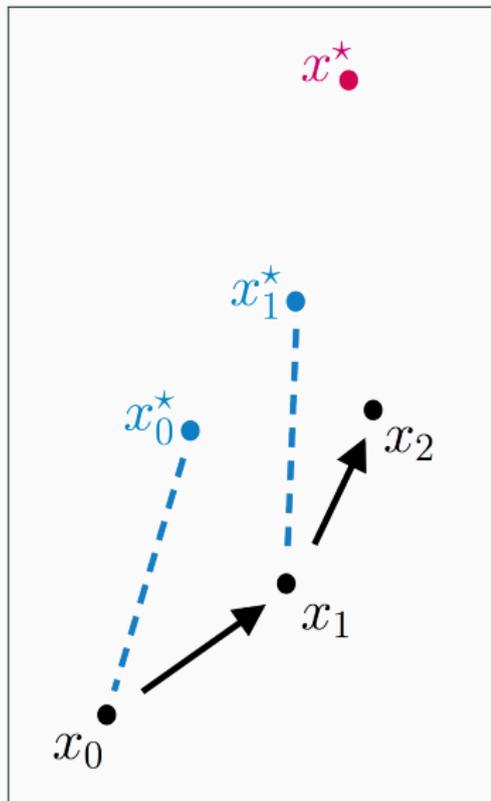
Theoretical results



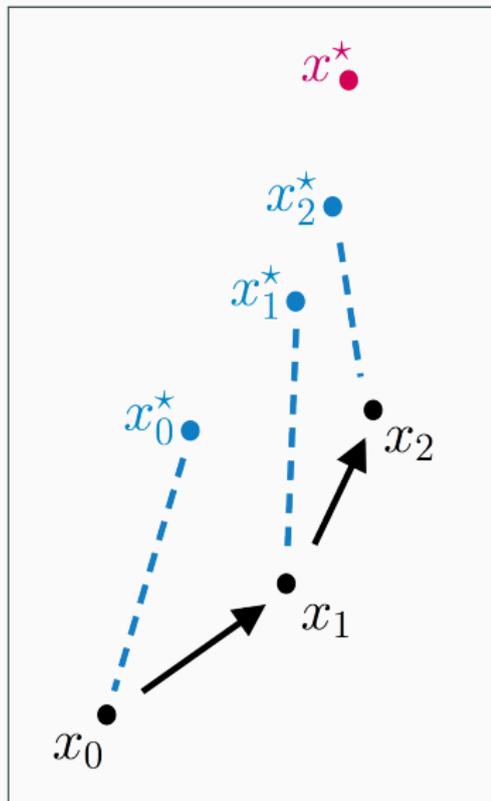
Theoretical results



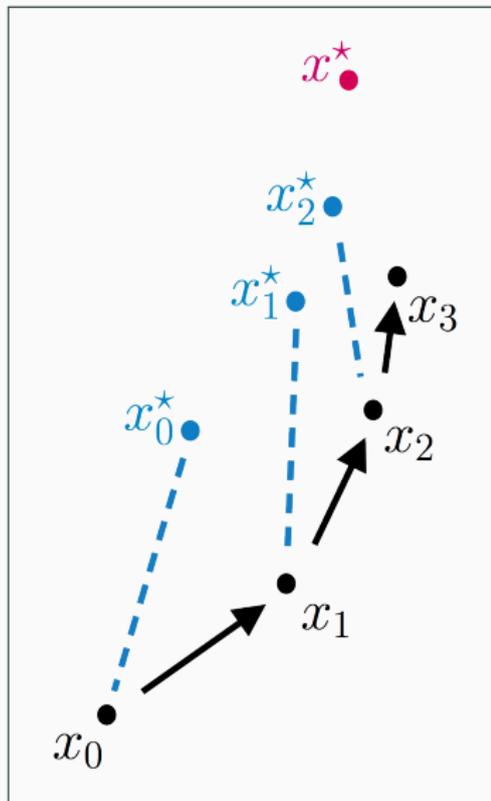
Theoretical results



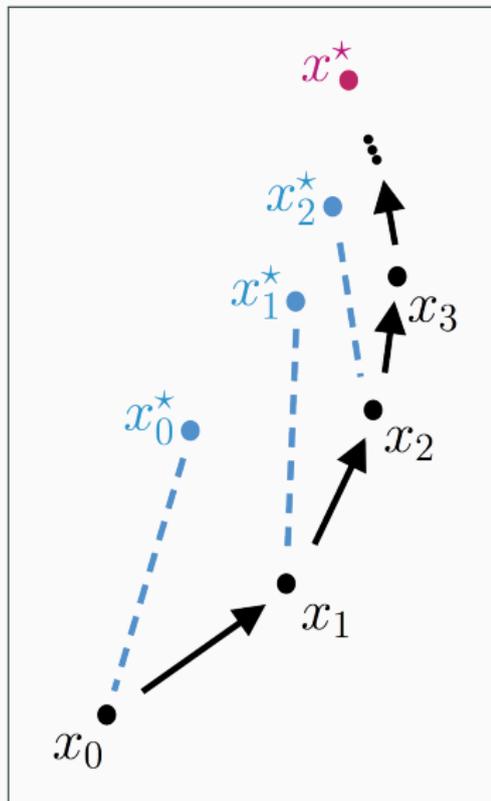
Theoretical results



Theoretical results



Theoretical results



Future directions

Optimization & theory:

- Is the stable “inner loop” version of MOCCA necessary?
 - Without RSC, guarantee convergence to stationary point?
 - Adaptive step sizes for faster convergence?
-

CT imaging:

- Detector sensitivity is not known exactly & may vary over detector cells \rightsquigarrow data-adaptive calibration?
- Apply MOCCA directly to Poisson likelihood, without quadratic approximation?

Thank you!

Website: <http://www.stat.uchicago.edu/~rina/mocca.html>

Funding: partially supported by NIH Grants R21EB015094, CA158446, CA182264, and EB018102. The contents of this presentation are solely the responsibility of the authors and do not necessarily represent the official views of the National Institutes of Health.