# Deconvolution of interferometric images

## Marco Prato

Department of Physics, Computer Science and Mathematics
University of Modena and Reggio Emilia, Italy

Have you ever done one of those spot-the-difference newspaper puzzles where you have to find the missing details using two very similar cartoons? The quick way to solve them is to cut out the two images, place one on top of the other, and shine a light through the paper.

It might sound like cheating but it's actually science: you're using the light pattern from one image to show up differences in the other.

Scientists use a very similar process called interferometry to measure small things with incredibly high accuracy by comparing light or radio beams.

To understand interferometry, you need to understand interference. In everyday life, interference simply means getting in the way or meddling, but in physics it has a much more specific meaning.
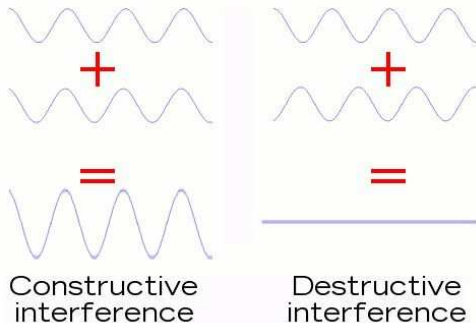
Interference is what happens when two waves carrying energy meet up and overlap. The energy they carry gets mixed up together so, instead of two waves, you get a third wave whose shape and size depends on the patterns of the original two waves. When waves combine like this, the process is called superposition.



*Wherever a crest coincides with a trough, the water surface is flattened.*

**Double crest-** A crest coincide with a crest

**Double trough-** A trough coincide with a trough

If the two waves are in step each other, they add themselves and increase the size of their peaks (amplitude). When waves add together to make bigger waves, scientists call it constructive interference.
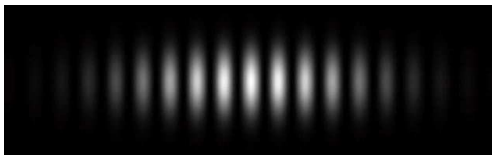
On the contrary, if the two waves are out of step, they subtract energy one to the other and make them smaller. This is what scientists call destructive interference.



Constructive interference          Destructive interference

The extent to which one wave is in step with another is known as its phase. If two identical waves are "in phase", it means their peaks align so, if we add them together, we get a new wave that's twice as big but otherwise exactly the same as the original waves.

Similarly, if two waves are completely out of phase, the peaks of one exactly coincide with the troughs of the other so adding the waves together gives you nothing at all.

In between these two extremes are all sorts of other possibilities where one wave is partly in phase with the other. Adding two waves like this creates a third wave that has an unusual, rising and falling pattern of peaks and troughs. Shine a wave like this onto a screen and you get a characteristic pattern of light and dark areas called interference fringes. This pattern is what you study and measure with an interferometer.

The basic idea of interferometry involves taking a beam of light and splitting it into two equal halves using what's called a beam-splitter.

If you shine light at it, half the light passes straight through and half of it reflects back.

One of the beams shines onto a mirror and from there to a screen, camera, or other detector. The other beam shines at or through something you want to measure, onto a second mirror, back through the beam splitter, and onto the same screen.

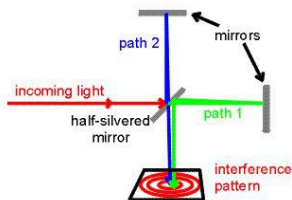This second beam travels an extra distance to the first beam, so it gets slightly out of step (out of phase).

When the two light beams meet up at the screen, they overlap and interfere, and the phase difference between them creates a pattern of light and dark areas.



The light areas are places where the two beams have added together (constructively) and become brighter; the dark areas are places where the beams have subtracted from one another (destructively).

The exact pattern of interference depends on the different way or the extra distance that one of the beams has traveled. By inspecting and measuring the fringes, you can calculate this with great accuracy and that gives you an exact measurement of whatever it is you're trying to find.

Instead of the interference fringes falling on a simple screen, often they're directed into a camera to produce a permanent image called an interferogram. In another arrangement, the interferogram is made by a detector (like the CCD image sensor used in older digital cameras) that converts the pattern of fluctuating optical interference fringes into an electrical signal that can be very easily analyzed with a computer.

Interferometers are widely used in all kinds of scientific and engineering applications for making accurate measurements. By scanning interferometers over objects, you can also make very precise maps of surfaces.
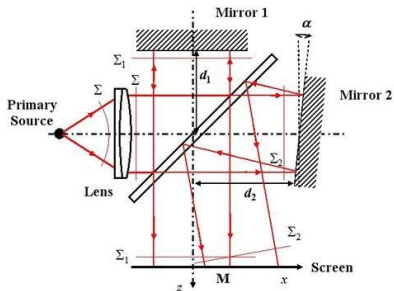
Astronomers also use interferometers to combine signals from telescopes so they work in the same way as larger and much more powerful instruments that can penetrate deeper into space. Some of these interferometers work with light waves; others use radio waves (similar to light waves but with much longer wavelengths and lower frequencies).



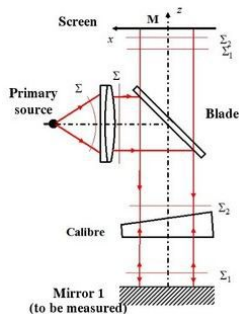Figure: The Keck interferometer (Photo courtesy of NASA Jet Propulsion Laboratory).

## Michelson interferometer

Probably best known for the part it played in the famous Michelson-Morley experiment in 1881.
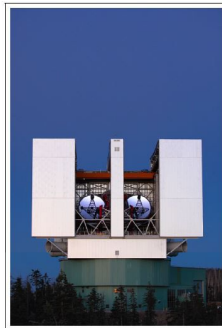
## Fizeau interferometer

It makes clearer and sharper fringes that are easier to see and measure. It's widely used for making optical and engineering measurements.

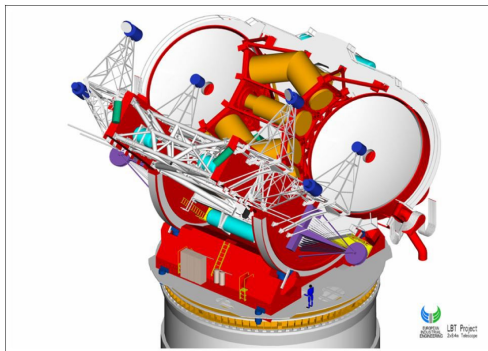The Large Binocular Telescope (LBT) is one of the most advanced telescopes operating in the world.

It is located near Safford, in Arizona (USA), within the Mt. Graham International Observatory that includes other two telescopes: the Vatican Advanced Technology Telescope and the Heinrich Hertz Submillimeter Telescope.

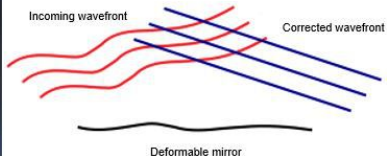The LBT project is a collaboration between the Italian astronomical community, represented by the National Institute for Astrophysics (INAF), the German astronomical community (mainly represented by the Max-Planck Institute for Astronomy at Heidelberg), the University of Arizona and other USA institutions.

The LBT consists of two 8.4 meters primary mirrors separated by a center-to-center distance of 14.4 meters.

To improve the performance of the telescope by reducing the effects of the wavefront distortions.

A deformable mirror can change its shape in real time.

First Light Adaptive Optics (FLAO) system @ LBT (2011)

Very close to (theoretical) Airy pattern (i.e., diffraction limit): Strehl ratio[1] values up to 0.9 in K-band



_____

[1] The Strehl ratio is the ratio of peak diffraction intensity of an aberrated versus perfect waveform. In the case of AO images this parameter can be estimated by the astronomers during the observation and provided with an error of few percent (about 4-5%).

Two interferometers are planned for LBT:

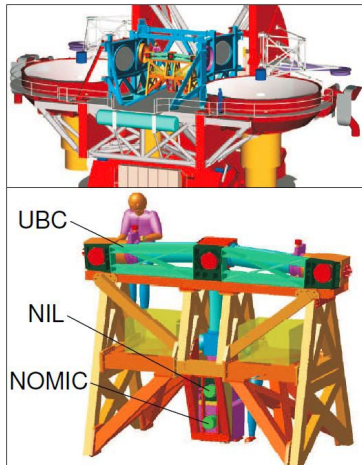- the LBT Interferometer (LBTI), designed for high spatial resolution, high dynamic range imaging in the thermal infrared and already operating on Mount Graham. The instrument consists of three main subsystems.

  - The Universal Beam Combiner (UBC) brings the radiation from the two optical trains to a common axis on the midline of the telescope
  - The Nulling Interferometer for the LBT (NIL) creates and overlays the two pupil images with the necessary $\pi$ phase shift
  - the Nulling-Optimized Mid-IR Camera (NOMIC) serves as the detector system

- the forthcoming Lbt INterferometric Camera - Near InfraRed Visible Adaptive INterferometer for Astronomy (LINC–NIRVANA), in advanced realization stage by a German–Italian consortium led by MPIA, Heidelberg. Two planned stages:
  - LINC–mode: interferometric near-infrared imaging with classical AO correction (with a single on–axis natural guide star).
  - NIRVANA–mode: layer–oriented (two–layers) multi–conjugated adaptive optics, whose detailed performance will depend not only on the atmospheric conditions (as classical AO does), but strongly also on the number of stars, their spatial distribution, and their magnitudes.

Any image acquired with a telescope has a limited band, i.e., the modulation transfer function (MTF)[2] of the telescope is zero outside a certain circular domain. The radius of this domain can be analytically calculated.

The monochromatic[3] PSF of LINC–NIRVANA, in the case of perfect optics and with no atmosphere, with the baseline aligned along $\xi$, is:

$$h(\xi, \eta) = \frac{2\Omega_1^2}{\pi} \left( \frac{J_1(\Omega_1|\theta|)}{\Omega_1|\theta|} \right)^2 \cos^2(\Omega_2\xi),$$

where $\theta^2 = \xi^2 + \eta^2$, $\Omega_1 = \pi D/\lambda$ ($D = 8.4$m, the primary mirror of LBT), $\Omega_2 = \pi B/\lambda$ ($B = 14.4$m, the baseline length) and $J_1$ is the Bessel function of the first kind defined e.g. as

$$J_1(x) = \frac{1}{\pi} \int_0^\pi \cos(\tau - x\sin(\tau))d\tau, \qquad \forall x \in \mathbb{R}.$$

---

[2]The MTF is formally defined as the magnitude (absolute value) of the Fourier transform of the point spread function (PSF, that is, the impulse response of the optics, the image of a point source).

[3]For simplicity, the bands are supposed monochromatics, i.e., only the central wavelength is considered.

The MTF of such a PSF is formed by a central cone of diameter proportional to $\Omega_1$ (corresponding to the primary $8.4$m mirror) and two side lobs separated by a center-to-center distance proportional to $\Omega_2$ (corresponding to the interferometric part of the PSF).

Fizeau
Interferometry

Fizeau
Interferometry

Fizeau Interferometry

With a few observations at different parallactic angles, and **the right algorithm**, the image can be reconstructed.

Figure: MTF of three interferometric PSFs (angles $0°$, $60°$, $120°$) combined



Figure: MTF of the PSF of a single $22.8$m telescope

Under suitable physical approximations, the mapping that transforms the object into the radiation incoming on the detector can be defined by

$$\bar{y}^{(k)} = H^{(k)}x + b^{(k)} \ ,$$

where

- $\bar{y}^{(k)} = \{\bar{y}_i^{(k)}\}_{i=1}^m \in \mathbb{R}^m$ is a vector of sampled values of the radiation before detection (also called exact data values);
- $x = \{x_j\}_{j=1}^n \in \mathbb{R}^n$ is the unknown object;
- $b^{(k)} = \{b_i^{(k)}\}_{i=1}^m$ is the background radiation affecting the $k$–th measured image (due to sky emission, dark current, etc.);
- $H^{(k)}$ is the $m \times n$ imaging matrix such that

$$H_{i,j}^{(k)} \geq 0 \ ; \ \sum_{i=1}^m H_{i,j}^{(k)} > 0 \ , \ \forall j \ ; \ \sum_{j=1}^n H_{i,j}^{(k)} > 0 \ , \ \forall i \ .$$

We assume periodic boundary conditions, so that $H^{(k)}$ is the block–circulant with circulant blocks matrix corresponding to the $k$–th PSF $h^{(k)}$.

We model the images according to the model proposed by Snyder et al.[4] for images acquired with a CCD camera, i.e., each pixel is affected by:

- photon counting noise (described by a Poisson distribution);
- additive read–out noise (RON), described by a Gaussian distribution.

The $i$–th pixel $y_i^{(k)}$ of the $k$–th measured image is then modeled as a realization of the random variable

$$Y_i^{(k)} = \mathcal{P}((H^{(k)}x + b^{(k)})_i) + \mathcal{N}(0, \sigma^2),$$

being

- $\mathcal{P}(\lambda)$ a Poisson random variable with mean and variance equal to $\lambda$;
- $\mathcal{N}(0, \sigma^2)$ is a Gaussian random variable with zero mean and variance equal to $\sigma^2$.

[4]Snyder DL, Hammoud AM and White RL 1993 Image recovery from data acquired with a charge-coupled-device camera, *J Opt Soc Am A* **10** 1014–1023

Three possibilities:

- address the mixed Poisson–Gaussian model[5];

- use a variance stabilizing transformation (e.g., generalized Anscombe) based approach to remove signal–dependency by rendering the noise approximately Gaussian[6];

- add $\sigma^2$ both to the detected images and the corresponding backgrounds and view all the pixel values of the detected images as realizations of suitable Poisson random variables[7].

---

[5]Chouzenoux E, Jezierska A, Pesquet J-C, Talbot H 2015, A convex approach for image restoration with exact Poisson–Gaussian likelihood, *SIAM J Imaging Sci* **8** 662–2682

[6]Starck J-L, Murtagh F, Bijaoui A 1998, *Image processing and data analysis*, Cambridge University Press, Cambridge

[7]Snyder DL, Helstrom CW, Lanterman AD, Faisal M, White RL 1995, Compensation for readout noise in CCD images, *J Opt Soc Am A* **12** 272–283

For $K$ given detected images $y = (y^{(1)}, \ldots, y^{(K)})$, let us introduce the likelihood function defined by

$$L_y^Y(x) = \sum_{k=1}^{K} p_{Y^{(k)}}(y^{(k)}; x) \ .$$

The ML-estimate of the unknown object is any object $x^*$ that maximizes the likelihood function

$$x^* = \underset{x \in \mathbb{R}^n}{\operatorname{argmax}} \ L_y^Y(x) \ .$$

Sums are better than products $\longrightarrow$ we consider the logarithm of the likelihood.

Minimization is more standard that maximization $\longrightarrow$ we consider the negative logarithm.

Therefore we introduce the functional

$$J_0(x; y) = -A \ln L_y^Y(x) + B \ ,$$

where $A, B$ are suitable constants, and we solve

$$x^* = \underset{x \in \mathbb{R}^n}{\operatorname{argmin}} \ J_0(x; y) \ .$$

### Example

In the case of additive white Gaussian noise, by a suitable choice of the constants $A, B$, we obtain

$$J_0(x; y) = \sum_{k=1}^{K} \|H^{(k)}x + b^{(k)} - y^{(k)}\|^2 \ ,$$

and therefore the ML approach coincides with the well–known least–squares (LS) approach.

$J_0$ is convex, strictly convex if and only if the equations $H^{(k)}x = 0$ have only the solution $x = 0$, and it has always global minima.

Problem: the condition number of $H^{(k)}$ can be very large.

## Example

In the case of Poisson noise, by exploiting the Stirling approximation $\ln(n!) \approx n \ln(n) - n$, the functional $J_0(x; y)$ becomes the sum of the so-called Kullback–Leibler (KL) divergences of $H^{(k)}x + b^{(k)}$ from $y^{(k)}$:

$$J_0(x; y) = \sum_{k=1}^{K} D_{KL}(y^{(k)}; H^{(k)}x + b^{(k)})$$

$$= \sum_{k=1}^{K} \sum_{i=1}^{m} \left\{ y_i^{(k)} \ln \frac{y_i^{(k)}}{(H^{(k)}x + b^{(k)})_i} + (H^{(k)}x + b^{(k)})_i - y_i^{(k)} \right\} .$$

Domain: non-negative orthant.

$J_0$ is convex, strictly convex if the equations $H^{(k)}x = 0$ have only the solution $x = 0$, non-negative and locally bounded $\rightarrow$ it has global minima.

The continuous version of $J_0$ and its minimization is an ill-posed problem $\rightarrow$ noise strongly affects the minima of the discrete problem (checkerboard effect).

## Example

In the case of Gauss+Poisson noise, the functional $J_0(x)$ is given by

$$J_0(x; y) = - \sum_{k=1}^{K} \sum_{i=1}^{m} \ln \sum_{l=0}^{+\infty} \frac{\exp^{-(H^{(k)}x + b^{(k)})_i} (H^{(k)}x + b^{(k)})_i^l}{l!} e^{-\frac{1}{2\sigma^2}(l - y_i^{(k)})^2}.$$

$J_0$ is convex, strictly convex if the equations $H^{(k)}x = 0$ have the unique solution $x = 0$, non–negative and locally bounded $\longrightarrow$ it has global minima on the non-negative orthant.

No result about the ill-posedness of this minimization problem (numerical experience shows that also in this case the minimum points are affected by the checkerboard effect).

The previous examples demonstrate that, in the case of image reconstruction, ML problems are ill-posed or ill-conditioned $\longrightarrow$ the minimum points $x^*$ do not provide sensible estimates $\bar{x}$ of the unknown object.

Therefore, one must be very careful in applying to these problems methods derived from optimization theory (in particular, second order methods).

Numerical experience (but not only) demonstrates that first order methods (Landweber, steepest descent, conjugate gradient, etc.) can provide acceptable (regularized) solutions by early stopping.

Since objects are non-negative, the non-negativity constraint must always be introduced in the formulation of the minimization problems.

The previous remark is not surprising in the framework of inverse problems theory. Indeed it is generally accepted that, if the formulation of the problem does not use some additional information on the object, then the resulting problem is ill-posed. This is what happens in the ML approach because we only use information about the noise with, possibly, the addition of the non-negativity constraint.

The additional information may consist, for instance, in prescribed bounds on the solution and/or its derivatives up to a certain order (in general not greater than two), and can be introduced by means of suitable functionals $J_R$ in the objective function:

$$\text{minimize} \quad J(x; y) = J_0(x; y) + \mu J_R(x)$$
$$\text{subject to} \quad x \geq 0 \ .$$

As follows from the examples discussed above, we can assume that both are convex so that we have a convex minimization problem.

### Theorem (Karush-Kuhn-Tucker)

*Consider the minimization problem*

$$\text{minimize} \quad f(x)$$
$$\text{subject to} \quad g_i(x) \geq 0 \ , \ \ i = 1, \ldots, p$$

*If $x^*$ is a regular point for the constraints[a] and a relative minimum point for the problem, then $\exists \mu \in \mathbb{R}^p$ such that $\mu \geq 0$ and*

$$\nabla f(x^*) - \sum_{i=1}^{p} \mu_i \nabla g_i(x^*) = 0 \ \ ,$$
$$\mu_i g_i(x^*) = 0 \ , \ \ i = 1, \ldots, p \ \ .$$

[a] If $x^*$ is a point satisfying the constraints $g(x^*) \geq 0$ and $J$ is the set of indices $j$ for which $g_j(x^*) = 0$, then $x^*$ is said to be a regular point for the constraints if the vectors $\nabla g_j(x^*)$, $j \in J$, are linearly independent.

Since $J(x; y)$ is convex all its minima are global. Then the KKT conditions are necessary and sufficient conditions for a point $x^*$ to be a minimum of $J(x; y)$

$$x^* \nabla J(x^*; y) = 0 \quad ,$$
$$x^* \geq 0 \quad , \quad \nabla J(x^*; y) \geq 0 \quad .$$

Let us consider now the following decomposition of the gradient[8,9]

$$-\nabla J(x; y) = U(x; y) - V(x; y) \quad ; \quad U(x; y) \geq 0 \quad , \quad V(x; y) > 0 \quad .$$

Then we can write the first KKT condition as a fixed point equation

$$x^* = T_y(x^*) \quad ,$$

with

$$T_y(x) = x \frac{U(x; y)}{V(x; y)} \quad .$$

[8]Lantéri A, Roche M, Cuevas O, Aime C 2001, A general method to devise maximum likelihood signal restoration multiplicative algorithms with non-negativity constraints, *Signal Process* **81**, 945–974

[9]Lantéri A, Roche M, Aime C 2002, Penalized maximum likelihood image restoration with positivity constraints: multiplicative algorithms, *Inverse Probl* **18**, 1397–1419

The operator $T_y(\cdot)$ is:

- well defined, since $V(x; y) > 0$;
- continuous if the functional $J(x; y)$ is continuously differentiable.

By applying the method of successive approximations we get the following iterative algorithm

$a)$ choose $x^{(0)} > 0$;

$b)$ for $\ell = 0, 1, \ldots$, compute

$$x^{(\ell+1)} = x^{(\ell)} \frac{U(x^{(\ell)}; y)}{V(x^{(\ell)}; y)} \ .$$

### Theorem

*If the sequence of the iterates $\{x^{(\ell)}\}_{\ell \in \mathbb{N}}$ is convergent to $x^*$ and if $U(x; y) > 0$ for any $x > 0$, then $x^*$ is a minimum point of $J(x; y)$.*

### Proof.

It is sufficient to prove that $x^*$ satisfies the KKT conditions. □

## Remark

- about the convergence of the algorithm nothing can be said at this stage of the analysis since the operator $T_y(\cdot)$, in general, is not a contraction;
- all the iterates are automatically non-negative;
- the algorithm is a <span style="color:red">scaled gradient method</span>, with step-size 1:

$$x^{(\ell+1)} = x^{(\ell)} - S_k \nabla J(x^{(\ell)}; y)$$

where

$$S_\ell = \operatorname{diag}\left\{ \frac{x_j^{(\ell)}}{V_j(x^{(\ell)}; y)} \right\} \quad . \tag{1}$$

### Remark

In the papers of Lantéri et al. the algorithm is presented as a descent method with a step-size selection:

$$x^{(\ell+1)} = x^{(\ell)} + \lambda_\ell \frac{x^{(\ell)}}{V(x^{(\ell)};y)} \left\{ U(x^{(\ell)};y) - V(x^{(\ell)};y) \right\} \qquad (2)$$

The step-size $\lambda_\ell > 0$ is chosen in the following way:

1) an upper bound $\lambda_\ell^{(0)}$ is determined in order to ensure that $x^{(\ell+1)} \geq 0$. This is obtained by looking at the values of $j$ such that $x_j^{(\ell)} > 0$ and $[\nabla J(x^{(\ell)};y)]_j > 0$. If we denote by $I_+$ the set of these index values, then

$$\lambda_\ell^{(0)} = \min_{j \in I_+} \left\{ \frac{V_j(x^{(\ell)};y)}{V_j(x^{(\ell)};y) - U_j(x^{(\ell)};y)} \right\} \geq 1 \ .$$

2) the step-size $\lambda_\ell$ is optimized by a line search in the interval $(0, \lambda_\ell^{(0)}]$ using, for instance, the Armijo rule.

In such a way, the convergence of the method is ensured.

If we particularize the general algorithm to the noise models, we obtain two well-known algorithms proposed for image reconstruction.

Indeed, in the case of Gauss noise we obtain

$$x^{(\ell+1)} = x^{(\ell)} \frac{\sum_{k=1}^{K} (H^{(k)})^T y^{(k)}}{\sum_{k=1}^{K} (H^{(k)})^T H^{(k)} x^{(\ell)} + b^{(k)}}$$

$$\left( U_0(x;y) = 2 \sum_{k=1}^{K} (H^{(k)})^T y^{(k)} \quad , \quad V_0(x;y) = 2 \sum_{k=1}^{K} (H^{(k)})^T H^{(k)} x + b^{(k)} \right)$$

and this is the image iterative space reconstruction algorithm (ISRA), introduced by Daube-Witherspoon and Muehllehner[10], whose asymptotic convergence has been proved by De Pierro[11].

More precisely the original algorithm is with $b^{(k)} = 0$, but the proof of convergence can be easily extended to the case $b^{(k)} \neq 0$.

---

[10] Daube-Witherspoon ME, Muehllehner G 1986, An iterative image space reconstruction algorithm suitable for volume ECT, *IEEE T Med Imaging* **5**, 61–66

[11] De Pierro AR 1987, On the convergence of the iterative image space reconstruction algorithm for volume ECT, *IEEE T Med Imaging* **6**, 174–175

In the case of Poisson noise we obtain

$$x^{(\ell+1)} = \frac{x^{(\ell)}}{\sum_{k=1}^{K}(H^{(k)})^T \mathbb{1}} \sum_{k=1}^{K} (H^{(k)})^T \frac{y^{(k)}}{H^{(k)}x^{(\ell)} + b^{(k)}}$$

$$\left( U_0(x;y) = \sum_{k=1}^{K}(H^{(k)})^T \frac{y^{(k)}}{H^{(k)}x + b^{(k)}} \quad , \quad V_0(x;y) = \sum_{k=1}^{K}(H^{(k)})^T \mathbb{1} \right)$$

and this is the generalization to the multi–image case of the expectation maximization (EM) algorithm proposed by Shepp and Vardi[12] and known as Richardson-Lucy (RL) algorithm.

The convergence proof in the case $b^{(k)} = 0$ is based on the following property: if the matrix $H^{(k)}$ is normalized in such a way that $(H^{(k)})^T \mathbb{1} = \mathbb{1}$, then

$$\sum_{j=1}^{n} x_j^* = \sum_{j=1}^{n} x_j^{(\ell)} = \frac{1}{K} \sum_{k=1}^{K} \sum_{i=1}^{m} y_i^{(k)} \qquad \text{(flux conservation)}.$$

This property is not satisfied in the case $b^{(k)} \neq 0$. Therefore the convergence of the algorithm seems not to be proved in such a case.

---

[12] Shepp LA, Vardi Y 1982, Maximum likelihood reconstruction for emission tomography, *IEEE T Med Imaging* **1**, 113–122

In the case of a regularized functional, the general algorithm takes the form

$$x^{(\ell+1)} = x^{(\ell)} \frac{U_0(x^{(\ell)}; y) + \mu U_R(x^{(\ell)})}{V_0(x^{(\ell)}; y) + \mu V_R(x^{(\ell)})} \ ,$$

where $U_0(x; y), V_0(x; y)$ come from the likelihood while $U_R(x), V_R(x)$ come from the prior.

Table: $D$ is a matrix with non-negative entries and this example includes regularization in terms of the discrete Laplacian

| $J_R(x)$ | $U_R(x)$ | $V_R(x)$ |
|---|---|---|
| $\frac{1}{2}\|\|(I - D)x\|\|_2^2$ | $(D + D^T)x$ | $(I + D^T D)x$ |
| $\|\|x\|\|_1$ | 0 | 1 |

## Flux conservation

If the reconstructed image has to be used for quantitative analysis it is important to guarantee flux conservation.

For simplicity we assume that each matrix $H^{(k)}$ is normalized in such a way that $(H^{(k)})^T \mathbb{1} = \mathbb{1}$, so that

$$\sum_{i=1}^{m} \left( \sum_{j=1}^{n} H_{i,j}^{(k)} x_j \right) = \sum_{j=1}^{n} x_j \ .$$

Then, the flux condition or flux constraint is defined by

$$\sum_{j=1}^{n} x_j = \frac{1}{K} \sum_{k=1}^{K} \sum_{i=1}^{m} \{ y_i^{(k)} - b_i^{(k)} \} \doteq c \ .$$

If we introduce the flux constraint, the problem is modified as follows

$$\text{minimize} \qquad J(x; y) = J_0(x; y) + \mu J_R(x) \qquad (3)$$

$$\text{subject to} \qquad x \geq 0 \ , \ \sum_{j=1}^{n} x_j = c \ .$$

We denote by $\mathcal{C}$ the closed and convex set that is the intersection of the nonnegative orthant with the affine subspace defined by the flux condition. We remark that $\mathcal{C}$ is compact so that any sequence contained in $\mathcal{C}$ will contain convergent subsequences.

In order to solve problem (3), a scaled gradient projection (SGP) method has been proposed[13,14], which can be considered a generalization of the scaled gradient method (2).

Notations:

- for a given vector $x \in \mathbb{R}^n$, $\|x\|_D$ is the norm induced by the $n \times n$ symmetric positive definite matrix $D$ (i.e., $\|x\|_D = \sqrt{x^T D x}$);

- for some given positive scalars $c_1$ and $c_2$, $\mathcal{D}$ is the set of the $n \times n$ symmetric positive definite matrices $D$ such that

$$c_1\|x\|^2 \leq x^T D x \leq c_2\|x\|^2, \quad \forall\, x \in \mathbb{R}^n; \tag{4}$$

- $P_{\mathcal{C},D}(x)$ is the projection of $x \in \mathbb{R}^n$ over $\mathcal{C}$ in the norm $\|\cdot\|_D$, that is

$$P_{\mathcal{C},D}(x) = \operatorname*{argmin}_{z \in \mathcal{C}}\|z - x\|_D = \operatorname*{argmin}_{z \in \mathcal{C}} \left( \frac{1}{2} z^T D z - z^T D x \right).$$

[13]Bonettini S, Zanella R, Zanni L 2009, A scaled gradient projection method for constrained image deblurring, *Inverse Probl* **25**, 015002

[14]Bonettini S, Prato M 2015, New convergence results for the scaled gradient projection method, *Inverse Probl* **31**, 095008

| 1. *Initialization.* | Let $\alpha_{\min}, \alpha_{\max} \in \mathbb{R}$ be such that $0 < \alpha_{\min} < \alpha_{\max}$, $\beta, \gamma \in (0,1)$ and let $M$ be a positive integer. Set $x^{(0)} \in \mathcal{C}$, $\quad D_0 \in \mathcal{D}$, $\quad \alpha_0 \in [\alpha_{\min}, \alpha_{\max}]$. |

FOR $\ell = 0, 1, 2, \ldots$

| 2. *Projection.* | Compute the descent direction $d^{(\ell)} = P_{\mathcal{C}, D_\ell^{-1}}(x^{(\ell)} - \alpha_\ell D_\ell \nabla J(x^{(\ell)}; y)) - x^{(\ell)}$. |
| 3. *Line-search.* | Set $\lambda_\ell = 1$ and $\bar{J} = \max_{0 \leq j \leq \min\{\ell, M-1\}} J(x^{(\ell-j)}; y)$. |

WHILE $J(x^{(\ell)} + \lambda_\ell d^{(\ell)}; y) > \bar{J} + \gamma \lambda_\ell \nabla J(x^{(\ell)}; y)^T d^{(\ell)}$
$\quad\quad \lambda_\ell = \beta \lambda_\ell$
END

Set $x^{(\ell+1)} = x^{(\ell)} + \lambda_\ell d^{(\ell)}$.

| 4. *Update.* | Define $D_{\ell+1} \in \mathcal{D}$ and $\alpha_{\ell+1} \in [\alpha_{\min}, \alpha_{\max}]$. |

END

Several reasons make this approach appealing for problem (3):

- it is very simple: it belongs to the class of standard scaled gradient methods with variable step-length $\alpha_\ell$ and non-monotone line-search strategy;
- due to the special constraints of the problem and to appropriate choices of $D_\ell$, the projection operation in step 2 can be non-expensive;
- the iterative scheme can achieve good convergence rate by exploiting the effective step-length selection rules recently proposed in literature (Barzilai-Borwein rules[15,16], adaptive alternating strategies[17,18], Ritz values[19,20],...).

[15] Barzilai J, Borwein JM 1988, Two-point step size gradient methods, *IMA J Numer Anal* **8**, 141–148

[16] Birgin EG, Martinez JM, Raydan M 2000, Nonmonotone spectral projected gradient methods on convex sets, *SIAM J Optim* **10**, 1196–1211

[17] Dai YH 2003, Alternate stepsize gradient method, *Optimization* **52**, 395–415

[18] Zhou B, Gao L, Dai YH 2006, Gradient methods with adaptive step-sizes, *Comput Optim Appl* **35**, 69–86

[19] Fletcher R 2012, A limited memory steepest descent method, *Math Program* **135**, 413–436

[20] Porta F, Prato M, Zanni L 2015, A new steplength selection for scaled gradient methods with application to image deblurring, *J Sci Comput* **65**, 895–919

The choice of the scaling matrix $D_\ell$ must avoid to introduce significant computational costs and, in particular, it must keep the projection $P_{\mathcal{C}, D_\ell^{-1}}(\cdot)$ in step 2 non-expensive.

Diagonal scaling $\longrightarrow$ the projection is obtained by solving a separable quadratic program.

Modification of the scaling matrix defined in (1):

$$D_\ell = \mathsf{diag}\left\{ \max\left\{ c_1, \frac{x_j^{(\ell)}}{V_j(x^{(\ell)}; y)} \right\} \right\},$$

where $c_1 > 0$ is a prefixed threshold.

Bounds (4) are satisfied by choosing $c_2 = c/\nu$, with

$$\nu = \min_j \left\{ \min_{x \in \mathcal{C}} \{V_j(x; y)\} \right\}.$$

---

3. *Line-search.* Set $\lambda_\ell = 1$ and $\bar{J} = \max\limits_{0 \leq j \leq \min\{\ell, M-1\}} J(x^{(\ell-j)}; y)$.

WHILE $J(x^{(\ell)} + \lambda_\ell d^{(\ell)}; y) > \bar{J} + \gamma \lambda_\ell \nabla J(x^{(\ell)}; y)^T d^{(\ell)}$
   $\lambda_\ell = \beta \lambda_\ell$
END

---

The line-search step of the SGP consists in a non-monotone strategy that uses successive reductions of $\lambda_\ell$ to make $J(x^{(\ell+1)}; y)$ lower than the maximum of the objective function on the last $M$ iterations.

Of course, if $M = 1$ then the strategy reduces to the standard Armijo rule.

The updating rule for the step-length $\alpha_\ell$ is crucial for improving the convergence rate of the scheme.

Barzilai-Borwein (BB) rules: regard the matrix $B(\alpha_\ell) = (\alpha_\ell D_\ell)^{-1}$ as an approximation of the Hessian $\nabla^2 J(x^{(\ell)}; y)$ and force a quasi-Newton property on $B(\alpha_\ell)$:

$$\alpha_\ell^{\text{BB1}} = \underset{\alpha \in \mathbb{R}}{\text{argmin}} \| B(\alpha) s^{(\ell-1)} - z^{(\ell-1)} \|$$

or

$$\alpha_\ell^{\text{BB2}} = \underset{\alpha \in \mathbb{R}}{\text{argmin}} \| s^{(\ell-1)} - B(\alpha)^{-1} z^{(\ell-1)} \|,$$

where $s^{(\ell-1)} = x^{(\ell)} - x^{(\ell-1)}$ and $z^{(\ell-1)} = \nabla J(x^{(\ell)}) - \nabla J(x^{(\ell-1)})$.

In this way, the following step-lengths are obtained

$$\alpha_\ell^{\text{BB1}} = \frac{s^{(\ell-1)^T} D_\ell^{-1} D_\ell^{-1} s^{(\ell-1)}}{s^{(\ell-1)^T} D_\ell^{-1} z^{(\ell-1)}} \quad , \quad \alpha_\ell^{\text{BB2}} = \frac{s^{(\ell-1)^T} D_\ell z^{(\ell-1)}}{z^{(\ell-1)^T} D_\ell D_\ell z^{(\ell-1)}}.$$

The step-length selection rule implemented within SGP is the ABB$_{min1}$ strategy[21], which consists in the following adaptive alternation scheme:

IF $\ell \leq 20$ THEN
$\qquad \alpha_\ell = \min_{j=\max\{1, \ell+1-M_\alpha\}, \ldots, \ell} \alpha_j^{(2)};$ $\qquad (\odot)$
ELSE IF $\alpha_\ell^{(2)}/\alpha_\ell^{(1)} \leq \tau_\ell$ THEN
$\qquad$ Set $\alpha_\ell$ as in $(\odot)$
$\qquad \tau_{\ell+1} = 0.9 \cdot \tau_\ell;$
ELSE
$\qquad \alpha_\ell = \alpha_\ell^{(1)};$ $\qquad \tau_{\ell+1} = 1.1 \cdot \tau_\ell;$
ENDIF

where $M_\alpha$ is a prefixed positive integer and $\tau_1 \in (0, 1)$.

---

[21] Frassoldati G, Zanni L, Zanghirati G 2008, New adaptive stepsize selections in gradient methods, *J Ind Manag Optim* **4**, 299–312

Convergence results of the sequence generated by SGP have been proved:

- stationarity of any limit point for nonconvex $J$[22];

- convergence to a minimum point (if it exists) for convex $J$ under conditions on the eigenvalues of the scaling matrices sequence $\{D_\ell\}_{\ell \in \mathbb{N}}$:

$$\mu_\ell^2 = 1 + \xi_\ell, \quad \xi_\ell \geq 0, \quad \sum_{\ell=0}^{\infty} \xi_\ell < \infty.$$

  If $\nabla J$ is either globally Lipschitz continuous or locally Lipschitz continuous and $J$ is level bounded on $\mathrm{dom}(J)$, then a $\mathcal{O}(1/\ell)$ convergence rate with respect to the objective function holds[23];

- convergence to a limit point (if it exists) for Kurdyka–Łojasiewicz functions $J$, plus $J$–dependent results on the convergence rate of both sequences $\{x^{(\ell)}\}_{\ell \in \mathbb{N}}$ and $\{J(x^{(\ell)})\}_{\ell \in \mathbb{N}}$[24].

---

[22]Bonettini S, Zanella R, Zanni L 2009, A scaled gradient projection method for constrained image deblurring, *Inverse Probl* **25**, 015002

[23]Bonettini S, Prato M 2015, New convergence results for the scaled gradient projection method, *Inverse Probl* **31**, 095008

[24]Bonettini S, Loris I, Porta F, Prato M, Rebegoldi S 2017, On the convergence of a linesearch based proximal-gradient method for nonconvex optimization, *Inverse Probl* **33**, 055005

SGP, together with other powerful tools for simulating observations and reconstructing real data, is included in the Software Package AIRY 7.0[25], a freely downloadable (http://lagrange.oca.eu/caos) IDL–based package of the Code for Adaptive Optics Systems Problem-Solving Environment (CAOS PSE).



[25]La Camera A et al 2016, The software package AIRY 7.0: new efficient deconvolution methods for post-adaptive optics data, *Proceedings of SPIE* **9909**, 99097T

- Tikhonov regularizations

$$J_R(x) = \frac{1}{2} \sum_{\boldsymbol{n}} |x(\boldsymbol{n})|^2 \ , \quad J_R(x) = \frac{1}{2} \sum_{\boldsymbol{n}} \boldsymbol{D}^2(\boldsymbol{n}) \ , \quad J_R(x) = \frac{1}{2} \sum_{\boldsymbol{n}} (\Delta x)(\boldsymbol{n})^2$$

- Cross-Entropy regularization [Byrne CL 1993, *IEEE T Image Proc* **2**, 96–103]

$$J_R(x) = \sum_{\boldsymbol{n}} \left\{ x(\boldsymbol{n}) \ln\left(\frac{x(\boldsymbol{n})}{\bar{x}(\boldsymbol{n})}\right) + \bar{x}(\boldsymbol{n}) - x(\boldsymbol{n}) \right\}$$

- $\ell_1$ regularization

$$J_R(\boldsymbol{f}) = \sum_{\boldsymbol{n}} x(\boldsymbol{n})$$

- Hypersurface potential [Charbonnier P et al 1997, *IEEE T Image Proc* **6**, 298–311]

$$J_R(\boldsymbol{f}) = \sum_{\boldsymbol{n}} \sqrt{\delta^2 + \boldsymbol{D}^2(\boldsymbol{n})} \ , \quad \delta > 0$$

---

$$\boldsymbol{n} = (n_1, n_2) \quad , \quad x(\boldsymbol{n}) = x(n_1, n_2) \quad , \quad \boldsymbol{n}_{1\pm} = (n_1 \pm 1, n_2) \quad , \quad \boldsymbol{n}_{2\pm} = (n_1, n_2 \pm 1)$$

$$\boldsymbol{D}^2(\boldsymbol{n}) = \left[x(\boldsymbol{n}_{1+}) - x(\boldsymbol{n})\right]^2 + \left[x(\boldsymbol{n}_{2+}) - x(\boldsymbol{n})\right]^2$$

- Markov random field regularization [Geman S, Geman D 1984, *IEEE T Pattern Anal Mach Intell* **6**, 721–741]

$$J_R(x) = \frac{1}{2} \sum_{\boldsymbol{n}} \sum_{\boldsymbol{n}' \in \mathcal{N}(\boldsymbol{n})} \sqrt{\delta^2 + \left( \frac{x(\boldsymbol{n}) - x(\boldsymbol{n}')}{\epsilon(\boldsymbol{n}')} \right)^2} \ ,$$

where $\delta > 0$, $\mathcal{N}(\boldsymbol{n})$ is a symmetric neighborhood made up of the eight first neighbors of $\boldsymbol{n}$ and $\epsilon(\boldsymbol{n}')$ is equal to 1 for the horizontal and vertical neighbors and equal to $\sqrt{2}$ for the diagonal ones.

- MISTRAL regularization [Mugnier LM et al 2004, *J Opt Soc Am A* **21**, 1841–1854]

$$J_R(\boldsymbol{f}) = \sum_{\boldsymbol{n}} \left\{ |\boldsymbol{D}(\boldsymbol{n})| - \delta \ln \left( 1 + \frac{|\boldsymbol{D}(\boldsymbol{n})|}{\delta} \right) \right\} \ , \ \delta > 0$$

---

$$\boldsymbol{n} = (n_1, n_2) \ , \quad x(\boldsymbol{n}) = x(n_1, n_2) \ , \quad \boldsymbol{n}_{1\pm} = (n_1 \pm 1, n_2) \ , \quad \boldsymbol{n}_{2\pm} = (n_1, n_2 \pm 1)$$

$$|\boldsymbol{D}(\boldsymbol{n})| = \sqrt{\left[ x(\boldsymbol{n}_{1+}) - x(\boldsymbol{n}) \right]^2 + \left[ x(\boldsymbol{n}_{2+}) - x(\boldsymbol{n}) \right]^2}$$

- PSFs corresponding to three equispaced orientations of the baseline ($0°$, $60°$, $120°$) generated with the code LOST[26](SR $\sim$ 70%)
- two test objects: a $512 \times 512$ HST image of the planetary nebula NGC7027 (with two magnitudes, 10 and 15) and a model of an open star cluster based on an image of the Pleiades, consisting of 9 stars with magnitudes ranging from 12.86 to 15.64
- blurred images perturbed with a constant background of about 13.5 mag arcsec$^{-2}$, corresponding to observations in K-band, and with both Poisson and Gaussian noises ($\sigma = 10 \; e^-$/px)
- added $\sigma^2$ to both images and backgrounds and chosen $J_0$ as the Kullback–Leibler divergence (no explicit regularization)

---

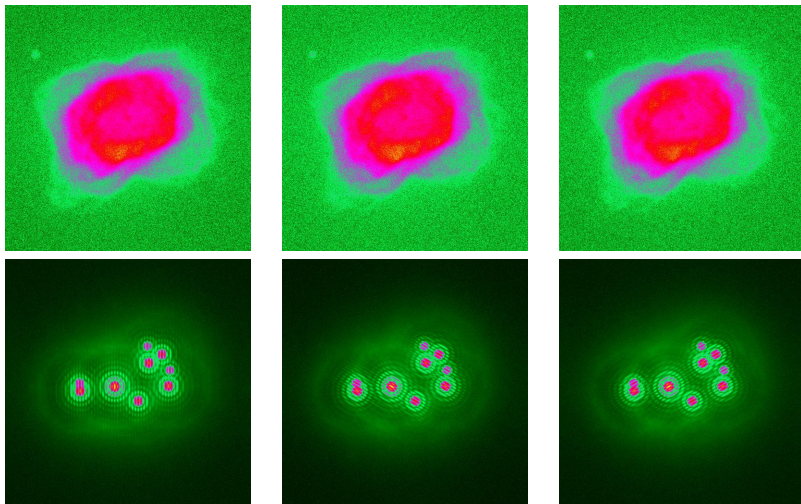[26]Arcidiacono C et al 2004, Layer-oriented simulation tool, *Appl Optics* **43**, 4288–4302

Figure: Blurred and noisy images for the nebula (magnitude 15 - top row) and the star cluster (bottom row)

Table: Reconstruction of the nebula using three equispaced $512 \times 512$ images (CUDA version of the algorithms based on GPUlib). Number of iterations tuned in order to minimize the relative RMSE w.r.t. the true object.

| | | | | |
|---|---|---|---|---|
| $m = 10$ | | | | |
| Algorithm | It | Err | Sec | SpUp |
| RL | 3401 | 0.032 | 4364 | - |
| RL_CUDA | 3401 | 0.032 | 48.00 | 90.9 |
| SGP | 144 | 0.033 | 220.7 | - |
| SGP_CUDA | 144 | 0.033 | 3.563 | 61.9 |
| $m = 15$ | | | | |
| Algorithm | It | Err | Sec | SpUp |
| RL | 353 | 0.091 | 441.5 | - |
| RL_CUDA | 353 | 0.091 | 4.937 | 89.4 |
| SGP | 16 | 0.087 | 26.14 | - |
| SGP_CUDA | 16 | 0.087 | 0.546 | 47.9 |

Table: Reconstruction of the star cluster with three $512 \times 512$ equispaced images. The error is the average relative error in the magnitudes. Number of iterations tuned in order to satisfy a stopping criterion on the difference between two successive values of $J_0$.

| Algorithm | It | Err | Sec | SpUp |
|---|---|---|---|---|
| $\mu$ = 1e-3 | | | | |
| RL | 319 | 2.39e-4 | 393.4 | - |
| RL_CUDA | 319 | 2.38e-4 | 4.641 | 84.8 |
| SGP | 71 | 1.35e-3 | 97.80 | - |
| SGP_CUDA | 71 | 1.29e-3 | 1.641 | 59.6 |
| $\mu$ = 1e-5 | | | | |
| RL | 1385 | 6.65e-5 | 1703 | - |
| RL_CUDA | 1385 | 6.64e-5 | 19.38 | 87.9 |
| SGP | 337 | 5.89e-4 | 455.2 | - |
| SGP_CUDA | 337 | 1.79e-4 | 7.187 | 63.3 |
| $\mu$ = 1e-7 | | | | |
| RL | 7472 | 5.64e-5 | 9180 | - |
| RL_CUDA | 7472 | 5.98e-5 | 104.8 | 87.6 |
| SGP | 572 | 7.37e-5 | 772.6 | - |
| SGP_CUDA | 572 | 7.05e-5 | 12.20 | 63.3 |

$$\text{av\_rel\_er} = \frac{1}{q} \sum_{j=1}^{q} \frac{|m_j - \widetilde{m}_j|}{\widetilde{m}_j} \qquad , \qquad |J_0(x^{(k+1)}; y) - J_0(x^{(k)}; y)| \leq \mu \, J_0(x^{(k)}; y)$$
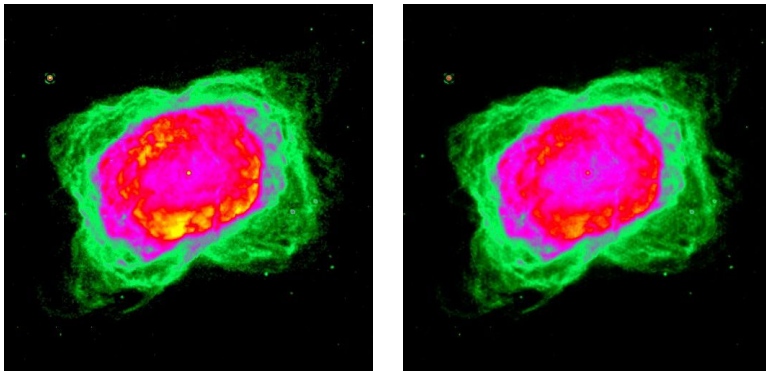
Figure: True image of the nebula (magnitude 10 - left panel) and SGP reconstruction (right panel)

- Seven interferometric images of Io observed with LBTI during UT 2013 December 24.
- PSF derived from the image of the star HD-78141.
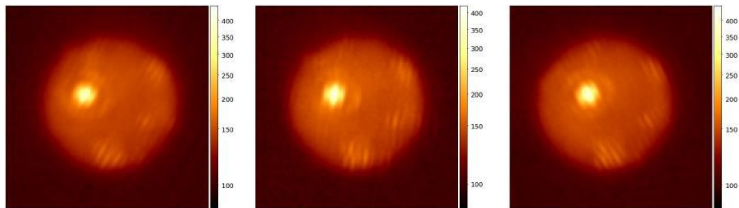- During the observation time of about 1 hour the Io relative orbital rotation is of $7.7°$.



Figure: Three (of seven) interferometric images, showing the variation of the parallactic angle of about $60°$

Idea[27,28]: considering the image to be reconstructed as the sum of two components $x_E$ (extended) and $x_P$ (pointwise)

New formulation:

$$\min_{(x_E, x_P) \in \overline{\Omega}} J(x_E, x_P; y) \equiv J_0(x_E + x_P; y) + \mu J_R(x_E),$$

where

$$\overline{\Omega} = \left\{ (x_E, x_P) \in \left( \mathbb{R}^{N \times N}_{\geq 0} \right)^2 \mid x_P(\boldsymbol{n}) = 0 \,\forall \boldsymbol{n} \in S \setminus P \,, \sum_{\boldsymbol{n} \in S} (x_E + x_P)(\boldsymbol{n}) = c \right\}$$

and $P$ is a given prefixed sub-region of the $N \times N$ region $S$ where the bright sources are located.

---

[27] De Mol C, Defrise M 2004, Inverse imaging with mixed penalties, *Proc Int Symp on Electromagnetic Theory*, 798–800
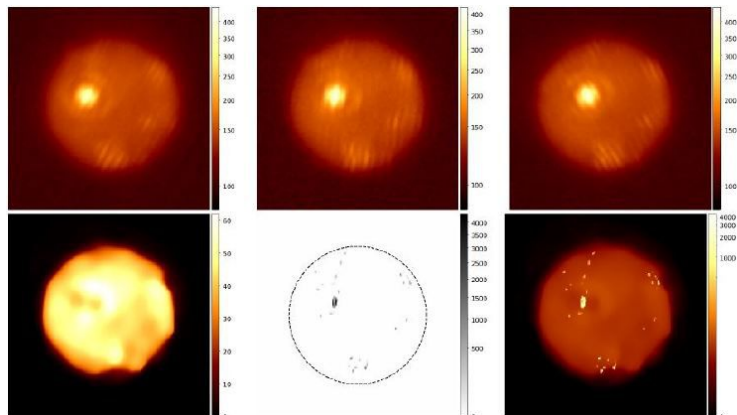
[28] Giovannelli J-F, Coulais A 2005, Positive deconvolution for superimposed extended source and point sources, *Astron Astrophys* **439**, 401–412

- Step 1 - Deconvolve the observed images with SGP with an edge–preserving regularization and a small value of $\mu$ (sharpening of the image + removal of the interferometric fringes)
- Step 2 - Determine the centroids of the bright regions and produce a mask which is one over small regions around the centroids, and zero elsewhere
- Step 3 - Apply MCM to the observed images and a regularizer which looks appropriate to the underlying structure
- Step 4 - If $x_E$ is the reconstruction of the structure in Step 3, then we write the unknown object as $x = x' + x_E$ and we can recover $x'$ by applying SGP, without regularization, to the observed images

The final result is the sum of the results of Step 3 and Step 4.

- **Initialization**: the fluxes of the point-wise objects in $x_P^{(0)}$ are chosen as those of the corresponding pixels of the background-subtracted observed image $y^{(1)}$. The remaining flux of the measured images (i.e., the value obtained by subtracting the flux of $x_P^{(0)}$ from the total flux $c$) is then spread on a constant $N \times N$ matrix $x_E^{(0)}$, which represents the starting point for the extended object.

- **$\mu$ parameter**: since $\mu$ provides a balance between the two terms of $J$, then one can estimate the value of $J_0$ and the order of magnitude of $J_R$, do a search around the value of $\mu$ provided by the quotient $J_0/J_R$ and look for a solution which could be the best for his purposes.

- **$\delta$ parameter**: one can compute the mean value $\delta_{\text{mean}}$ of the gradient on the observed image. Then a search of $\delta$ around the value of $\delta_{\text{mean}}$ is desirable, in order to find the best value for the user.

All the other SGP parameters have been optimized according to a huge amount of tests carried out in several applications and have been left unchanged.

First row: Three (of seven) interferometric images, showing the variation of the parallactic angle of about $60°$. Second row: Reconstructed surface of Io as obtained at Step 3 with MRF regularization, $\delta = 1$ and $\mu = 0.05$ (left), reconstruction of the hot spots (middle) and complete reconstruction (right).

Generalizations of the approach have been developed to address:

- **Boundary effect corrections**. If the target $x$ is not completely contained in the image domain, then the previous deconvolution methods produce annoying boundary artifacts. The proposed approach can be generalized in order to reconstruct $x$ over a domain broader than that of the detected images[29].

- **Blind deconvolution**. If the PSF is not available, then it might be included in the unknowns of the minimization problem, where the corresponding feasible set results in the set of non–negative arrays normalized to 1 and with maximum value deduced by the Strehl ratio of the telescope. The proposed approach can be exploited within an alternating minimization scheme to provide both the target $x$ and the PSFs $h^{(1)}, \ldots, h^{(K)}$[30,31].

---

[29] Prato M, Cavicchioli R, Zanni L, Boccacci P, Bertero M 2012, Efficient deconvolution methods for astronomical imaging: algorithms and IDL-GPU codes, *Astron Astrophys* **539**, A133

[30] Prato M, La Camera A, Bonettini S, Bertero M 2013, A convergent blind deconvolution method for post-adaptive-optics astronomical imaging, *Inverse Probl* **29**, 065017

[31] Prato M, La Camera A, Bonettini S, Rebegoldi S, Bertero M, Boccacci P 2015, A blind deconvolution method for ground based telescopes and Fizeau interferometers, *New Astron* **40**, 1–13