# Sistemi Intelligenti
## Stima MAP

Alberto Borghese

Università degli Studi di Milano
Laboratory of Applied Intelligent Systems (AIS-Lab)
Dipartimento di Scienze dell'Informazione
borghese@di.unimi.it

---

# Overview

Filtering images

MAP, Tikhonov and Poisson model of the noise

A-priori and Markov Random Fields

Cost function minimization

2 of 72

## Images are corrupted by noise...

i) When measurement of some physical parameter is performed, noise corruption cannot be avoided.

ii) Each pixel of a digital image measures a number of photons.

Therefore, from i) and ii)...

...Images are corrupted by noise!

http://ais-lab.dsi.unimi.it

3 / 46

---

## A general framework

▸ $f = \{f_1, f_2, f_M\}, \quad f_k \in R^M$      e.g. Pixel true luminance

▸ $g = \{g_1, g_2, g_M\} \quad g_k \in R^N$      e.g. Pixel measured luminance

▸ $g = A\,f + h + v$      -> determining f is a deblurring problem (the measuring mean transforms the image: scale + offset)

▸ $g = I\,f + v$      -> determining f is a denoising problem (the image is a copy of the real one with the addition of noise)

It is a general framework. It is a linear framework.

*h is the background radiation*
*v is the noise*

4 of 72

# Gaussian noise and likelihood

- Images are composed by a set of pixels, **f** (**f** is a vector!)
- Let us assume that the noise is Gaussian and that its mean and variance is equal for all pixels;
- Let $g_{n,i}$ be the measured value for the i-th pixel (n = noise);
- Let and $f_i$ be the true (noiseless) value for the i-th pixel;

- How can we quantify the probability to measure the image **f**, given the probability density function for each pixel?
- Likelihood function, $L(\mathbf{g_n} \mid \mathbf{f})$:

$$L(\mathbf{g_n} \mid \mathbf{f}; w) = \prod_{i=1}^{N} p(g_{n,i} \mid f_i) = \prod_{i=1}^{N} \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{g_{n,i} - f_i}{\sigma}\right)^2\right]$$

- $L(\mathbf{g_n} \mid \mathbf{f})$ describes the probability to measure the image $\mathbf{g_n}$, given the noise free value for each pixel, **f**. But we do not know these values….

# Statistical formulation of image restauration

Measuring an image **g** taken from an object, **f**, we want to determine **f**, when **g** is corrupted by noise: $\mathbf{g_n} = \mathbf{Af} + \mathbf{b} + \textit{noise} \rightarrow \mathbf{f}$? It is a typical **inverse problem**.

**A** is a linear operator that describes the transformation (mapping) from **f** to **g** (e.g. perspective projection, sensor transfer function, $\mathbf{A} = \mathbf{I}$ for denoising …).

**b** is the background radiation. It is the measure **g**, when no signal arrives to the sensor.

Each pixel is considered an independent process (white noise). **For each pixel** therefore, we want to find **f** that maximize: $p(\mathbf{g_n} ; \mathbf{f})$

Being the pixels independent, the total probability can be written in terms of product of independent probabilities (likelihood function):

$$L(g_n ; f) = \prod_{i=1}^{N} p(g_{n,i} ; f_i)$$

L is the likelihood function of $g_n$, given the object $f$.

## Do we get anywhere?

L is the likelihood function of $g_n$, given the object $f$.

$$L(g_n; f) = \prod_{i=1}^{N} p(g_i; f_i)$$

Determine $\{f_i\}$ such that L is maximized. Negative log-likelihood is usually considered to deal with sums:

$$-\log(L(.)) = -\sum_{i=1}^{N} \ln\left(p(g_{n,i}; f_i)\right)$$

$$f(g_{n,1}, g_{n,2}....g_{n,N}; f_{n,1}, f_{n,2}....f_{n,N}; 0, \sigma) = -\sum_{i=1}^{N} \ln\left\{\frac{1}{\sqrt{2\pi}\sigma} \cdot \exp\left[-\frac{1}{2}\left(\frac{g_{n,i} - f_i}{\sigma}\right)^2\right]\right\} \Rightarrow \qquad f = (A^T A)^{-1} A^T g_n$$

$$\min_{\{f_i\}}(f(.)) = \min\left\{-\sum_{i=1}^{N} \ln\left(\frac{1}{\sqrt{2\pi}\sigma}\right) + \frac{1}{2\sigma^2}\sum_{i=1}^{N}(g_{n,i} - f_i)^2\right\}_{\{f_i\}} \qquad \begin{array}{l} \text{if } A = I \\ f = g_n \end{array}$$

The system has a single solution, that is good. The solution is $f_i = g_{n,i}$, not a great result....

Can we do any better?

## Overview

Filtering images

MAP, Tikhonov and Poisson model of the noise

A-priori and Markov Random Fields

Cost function minimization

# The Bayesian framework

We assume that the object f is a realization of the "abstract" object F that can be characterized statistically as a density probability on F. f is extracted randomly from F.

The probability $p(g_n| f)$ becomes a conditional probability: $J_0 = p(g_n| f = f^*)$

Under this condition, the probability of observing f and $g_n$ (joint probability) can be written as the product of the conditional probability $p(g_n | f)$ by a-priori probability on f, $p_f$:

$$p(g_n, f) = p(g_n | f) p_f$$

As we are interested in determining f, we have to write the conditional probability of f given $g_n$ : $p(f | g_n)$. We apply Bayes theorem:

$$p(f | g_n) = \frac{p(g_n | f) p_f}{p_{g_n}} = L(g_n; f) \frac{p_f}{p_{g_n}}$$

# MAP estimate with logarithms

$$p(f | g_n) = \frac{p(g_n | f) p_f}{p_{g_n}} = L(g_n; f) \frac{p_f}{p_{g_n}}$$

Logarithms help:

$$-\ln(p(f | g_n)) = -\left\{ \frac{\ln(p(g_n | f) p_f)}{p_{g_n}} \right\} = -\left\{ \ln(p(g_n | f)) + \ln(p_f) - \ln(p_{g_n}) \right\}$$

We maximize the MAP of f | $g_n$, by minimizing:

$$\arg\min_f -\left\{ \ln\left( \frac{p(g_n | f) p_f}{p_{g_n}} \right) \right\} = \arg\min_f -\left\{ \ln(p(g_n | f)) + \ln(p_f) - \ln(p_{g_n}) \right\}$$

We explicitly observe that the marginal distribution of $p_{g,n}$ is not dependent on f. It does not affect the minimization and it can be neglected. It represents the statistical distribution of the measurements alone.

## MAP estimate with logarithms

We maximize the MAP of f | $g_n$, by minimizing:

$$\arg\min_{f} -\{\ln(p(g_n \mid f)p_f)\} = \arg\min_{f} -\{\ln(p(g_n \mid f)) + \ln(p_f)\}$$

$J_0(g_{n,i}; f)$  Likelihood = adherence to the data

A-priori $J_R(f)$

Depending on the shape of the noise (inside the likelihood) and the a-priori distribution of f(.), $J_R(f)$, we get different solutions.

## Gibb's priors

We often define the a-priori term, $J_R(f)$, as Gibb's prior:

$$p_f = \frac{1}{Z}\left\{e^{\left(-\frac{1}{\beta}U(f)\right)}\right\} \qquad Z = \int_{-\infty}^{+\infty} e^{-\frac{1}{\beta}U(f)} \, df = \cos t$$

$$J_R(f) = \ln(p_f) = -\ln(Z) - \frac{1}{\beta}U(f)$$

U(f) is also termed potential => $J_R(f)$ is a linear function of the potential U(f).

$\beta$ describes how fast the potential (the cost) decreases with U(f).

## Gaussian noise and a-priori term on the norm of the solution

$$f = \arg\min_{f} - \left\{ \ln\left(p(g_n \mid f) p_f\right) \right\} = \arg\min_{f} - \left\{ \ln(p(g_n \mid f)) + \ln(p_f) \right\} =$$

$$\arg\min_{f} \left\{ J_0(g_n; f) + J_R(f) \right\} =$$

Gaussian noise on the data: $\quad J_0(g_n; f) = \cos t. + \left( \sum_i \|g_{n,i} - Af_i\|^2 \right)$

We choose as a-priori term the squared norm of the function f, weighted by P.

$$p_f = \frac{1}{Z} \left\{ e^{\left(-\frac{1}{\beta}\|Pf\|^2\right)} \right\} \qquad \text{P = I} \qquad J_R(f) = \cos t + \frac{1}{\beta}\|f_i^2\|$$

$$f = \arg\min_{f} \left( \sum_i \|g_{n,i} - Af_i\|^2 + \frac{1}{\beta} \sum_i \|f_i\|^2 \right)$$

13 of 72

## Tikhonov regularization

$$f = \arg\min_{f} \left( \sum_i \|g_{n,i} - Af_i\|^2 + \frac{1}{\beta} \sum_i \|f_i\|^2 \right)$$

$$f = \arg\min_{f} \left( \sum_i \|g_{n,i} - Af_i\|^2 + \lambda \sum_i \|Pf_i\|^2 \right)$$

(cf. Ridge regression and Levemberg-Marquardt optimization)

It is a quadratic cost function. We find *f* minimizing with respect to f the cost function:

$$f : A^T g_n - A^T A f + \lambda P^T P f = 0 \quad => \quad A^T g_n = \left(A^T A + \lambda P^T P\right)$$

$$f : A^T g_n - A^T A f + \lambda P^T P f = 0 \quad => \quad A^T g_n = \left(A^T A + \lambda P^T P\right)f$$

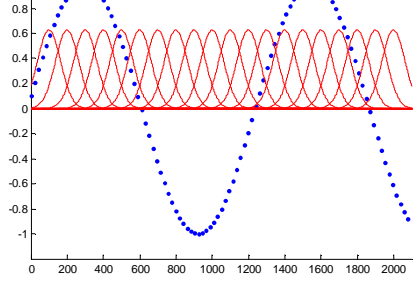$$\text{P = I} \qquad\qquad f : A^T g_n = \left(A^T A + \lambda I\right)f$$
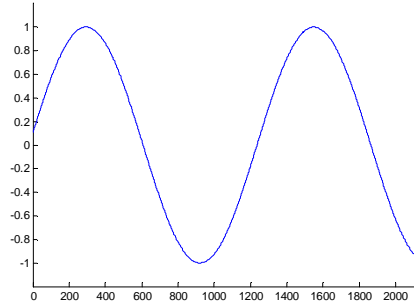
Poggio and Girosi, 1990

14 of 72

## Example



$$y(x) = \sum_k w_k h(x \mid x_k, a_k)$$

$$g = A\, f \quad con \ \{f_k\} = \{w_k\}$$

Good reconstruction when no noise is present.
Waved reconstruction with noise.

$$w = \left(A^T A\right)^{-1} A^T g$$

## Poisson case

$Noise_i = \|A\, f - g_{ni}\,\|$

We know the statistical distribution of the noise -> we now the statistical distribution of the second term. In case of Poisson noise we have:

For one pixel: $p(g_{ni}, f_i) = \left\{ \dfrac{e^{-Af_i}\left(Af_i\right)^{g_{n_i}}}{g_{n_i}!} \right\}$

$$-\ln\left(L(g_n; f)\right) = -\ln\left(\prod_{i=1}^{N} p\left(g_{n,i}; f_i\right)\right) = -\sum_{i=1}^{N} \left(-Af_i + g_{n,i}\ln(Af_i) - \ln\left(g_{n,i}!\right)\right)$$

To eliminate the factorial term, we normalize the likelihood by $L(g_n, g_n)$:

$$-\ln\left(\frac{L(g_n, f)}{L(g_{n,}g_n)}\right) = -\sum_{i=1}^{N}\left(g_n \ln(Af) - \ln(g_n) + g_n - Af\right) = KL\ divergence$$

$$= \sum_i g_n \ln\left(\frac{g_n}{Af} + Af - g_n\right)$$

It is not a distance!
It is not linear

## Gibbs priors and Regularization

$$\arg\min_f -\{\ln(p(g_n \mid f)p_f)\} = \arg\min_f -\{\ln(p(g_n \mid f)) + \ln(p_f)\}$$

Likelihood = adherence to the data

A-priori

Gaussian $\quad K(\sigma)\sum_i \left\| g_{n,i} - Af_i \right\|^2$

Poisson $\quad \sum_i g_{n,i} \ln\left(\dfrac{g_{n,i}}{Af} + Af_i - g_{n,i}\right)$

$$-\ln\left\{ \frac{1}{Z} e^{\left\{ -\frac{1}{\beta} U(\mathbf{f}) \right\}} \right\}$$

$$\boxed{J(f) = J_o(f) + \lambda J_R(f)}$$

$J_R(f) = U(f)$

## What happens if noise is Poisson?

$$f = \arg\min_f -\{\ln(p(g_n \mid f)p_f)\} = \arg\min_f -\{\ln(p(g_n \mid f)) + \ln(p_f)\}$$

**Poisson** noise model
Squared shape for the a-priori term

$$\sum_i g_{n,i} \ln\left(\frac{g_{n,i}}{Af} + Af_i - g_{n,i}\right)$$

$\| \lambda Pf \|^2$

$$f = \arg\min_f \sum_i g_{n,i} \ln\left(\frac{g_{n,i}}{Af} + Af_i - g_{n,i}\right) + \lambda\|Pf\|^2 \qquad \text{Regularization}$$

No analytical solution

## Overview

Filtering images

MAP, Tikhonov and Poisson model of the noise

A-priori and Markov Random Fields

Cost function minimization

## Which is the most adequate $p_f$ for images?

We usually ask to images to be smooth (we look at differential properties)

We look at the local gradient of the image: $\nabla f$.

One possibility is to use the square of the l-2 norm of the gradient: $\| \nabla f \|^2$

This is another form of Tikhonov regularization.

## Gibbs priors and Regularization

$$\arg\min_f -\{\ln(p(g_n \mid f)p_f)\} = \arg\min_f -\{\ln(p(g_n \mid f)) + \ln(p_f)\}$$

Likelihood =
adherence to the data

A-priori

Gaussian $\qquad K(\sigma)\sum_i \|g_{n,i} - Af_i\|^2$

$$-\ln\left\{\frac{1}{Z}e^{\left\{-\frac{1}{\beta}U(\mathbf{f})\right\}}\right\}$$

Poisson $\qquad \sum_i g_{n,i}\ln\left(\frac{g_{n,i}}{Af} + Af_i - g_{n,i}\right)$

$$\boxed{J(f) = J_o(f) + \lambda J_R(f)}$$

$J_R(f) = U(f)$

## A-priori on the derivatives

$$\arg\min_f -\{\ln(p(g_n \mid f)p_f)\} = \arg\min_f -\{\ln(p(g_n \mid f)) + \ln(p_f)\}$$

$$\boxed{J(f) = J_o(f) + \lambda J_R(f)} \qquad J_R(f) = funzione\left(\|\nabla f\|^2\right)$$

$$\arg\min_f \{(Af - g_n)^2 + \lambda\nabla f^2\} \arg\min_f \{\|(Af - g_n)^2\| + \lambda\|\nabla f\|^2\}$$

$$f: \{2A^T(Af - g_n) + 2\lambda\nabla f\} = 0$$

If we apporximate $\nabla f$ with the fiinite differences: $f_i - f_j$, we get a linear system.

## Non-quadratic a-priori: total variation

$$f = \arg\min_f -\{\ln(p(g_n \mid f) p_f)\} = \arg\min_f -\{\ln(p(g_n \mid f)) + \ln(p_f)\}$$

Poisson noise model

$$\sum_i g_{n,i} \ln\left(\frac{g_{n,i}}{Af} + Af_i - g_{n,i}\right)$$

**The a-priori term is a gradient and it is expressed in l$_2$ norm**

$$\sum_i \sqrt{\sum_p^P f_{p,i}^2} \qquad \sum_i \sqrt{\left(f_{x,i}^2 + f_{y,i}^2 + f_{z,i}^2\right)}$$

$$f = \arg\min_f \sum_i \left(\|g_n - Af\|^2 + \lambda \sqrt{\sum_p^P f_{p,i}^2}\right) \qquad \textbf{Total variation}$$

The derivative is not linear anymore because of the square root.

23 of 72

## Tikhonov regularization - simulations



Edge smoothing effect with Tikhonov-like regularization
Poisson noise model – $\lambda$ = 0.5
P is the gradient operator

24 of 72

## Total variation regularization - simulations



No appreciable edge smoothing with total variation
Poisson noise model - $\lambda = 0.5$
P is the gradient operator

## Tikhonov regularization – panoramic images



Edge smoothing effect with Tikhonov-like regularization
Poisson noise model - $\lambda = 0.5$
P is the gradient operator

# Total variation regularization – panoramic images

Original

Iter = 20 - Solution - F = 4386075.6946

Denoising effect - lambda = 0.5p = 1

No appreciable edge smoothing with total variation
Poisson noise model - $\lambda = 0.5$
P is the gradient operator

27 of 72

# Tikhonov regularization - endo-oral images

Original

Iter = 20 - Solution - F = 9759471.5548

Denoising effect - lambda = 0.1p = 2

Edge smoothing effect with Tikhonov-like regularization
Poisson noise model - $\lambda = 0.1$
P is the gradient operator

28 of 72

## Total variation – endo-oral images



No appreciable edge smoothing with total variation
Poisson noise model - $\lambda = 0.1$
P is the gradient operator

29 of 72

## A priori term – image gradients (no noise)



$p_x = p(i,j) - p(i-1,j)$
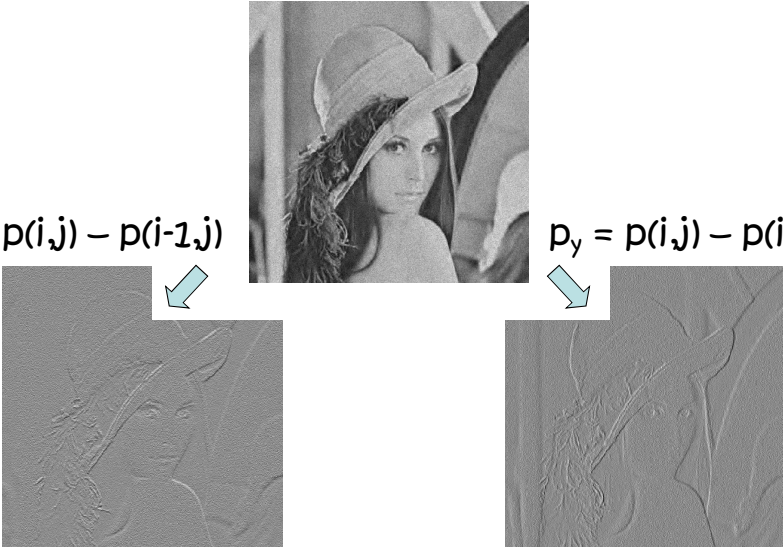
$p_y = p(i,j) - p(i,j-1)$

http://ais-lab.dsi.unimi.it

I. Frosio, M. Lucchese, N. A. Borghese

# A priori term – image gradients (noise)



$p_x = p(i,j) - p(i-1,j)$

$p_y = p(i,j) - p(i,j-1)$

http://ais-lab.dsi.unimi.it

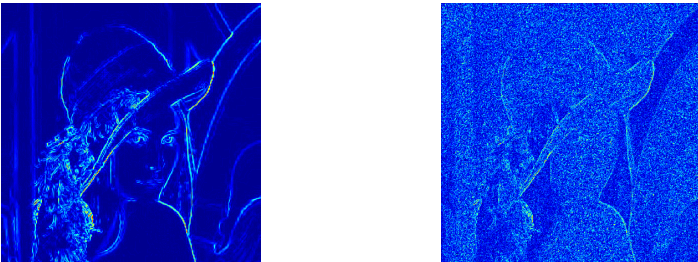I. Frosio, M. Lucchese, N. A. Borghese

# A priori term – norm of image gradient

No noise

Noise



In the real image, most of the areas are characterized by an (almost) null gradient norm;

**We can for instance suppose that $||\nabla p||$ is a random variable with Gaussian distribution, zero mean and variance equal to $\beta^2$.**

[Note that, in the noisy image, the norm of the gradient assume higher values → low $||\nabla p||$ means low noise!]
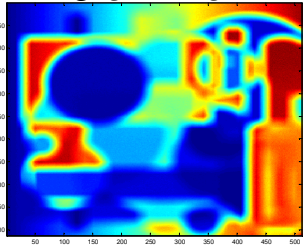
http://ais-lab.dsi.unimi.it

I. Frosio, M. Lucchese, N. A. Borghese

Tikhonov vs. TV (preview)

Original image

Tikhonov =>

TV =>

Filtered image

Difference

http://ais-lab.dsi.unimi.it

33 / 46

I. Frosio, M. Lucchese, N. A. Borghese



Cost introduced by the regularzation term

Cost increases quadratically with the local gradient in Tikhonov

34 of 72

17

# A-priori

We can insert in the a-priori term all the desirable characteristic of the image: local smoothness, edges, piece-wise constancy,....

The idea of defining a neighboring system is a natural one:

Neighbor region of $S_k$

Images have a natural neighboring system: the pixels structure. We want to consider the local properties of the image considering neighboring pixels (in particular differential properties - our vision system is particularly tuning to gradients both spatial and temporal). Ideas have been borrowed from physics.
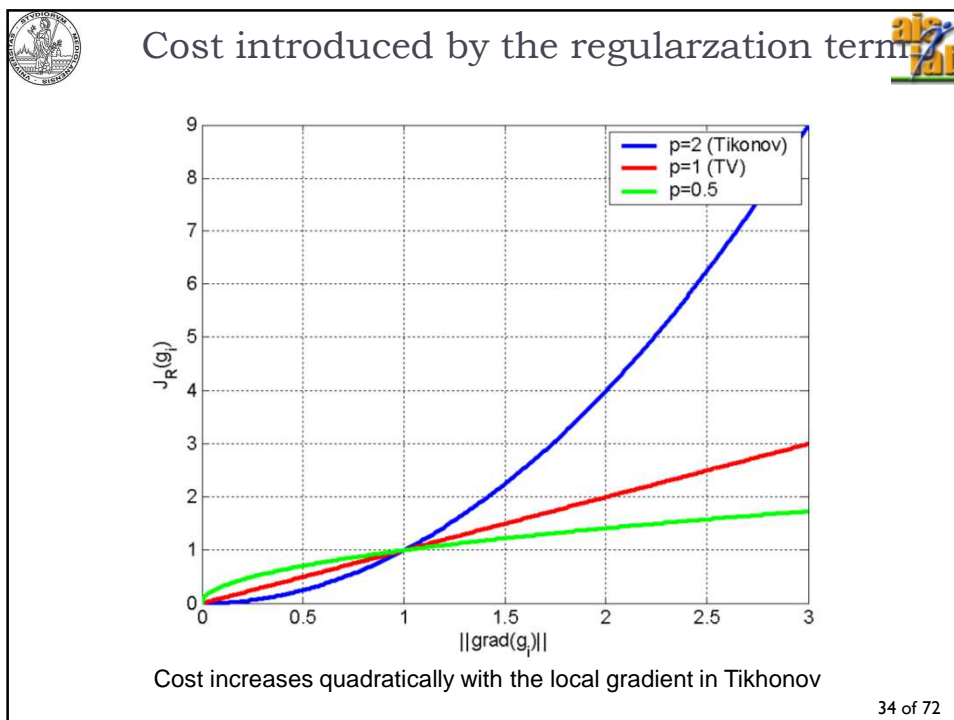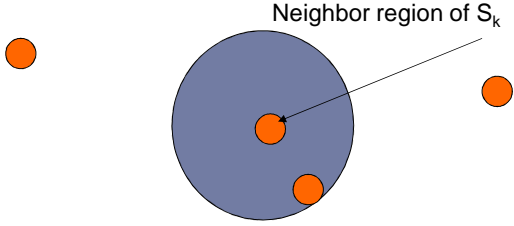
35 of 72

# Neighboring System

Let P be the set of pixels of the image: $P = \{p_1, p_2, \ldots p_P\}$

The neighboring system defined over P, S, is defined as $H = \{\mathcal{N}_p \mid p, \forall p \in P\}$, that has the following properties:

An element is not a neighbor of itself: $p_k \notin \mathcal{N}_{pk}$

Mutuality of the neighboring relationship: $p_k \in \mathcal{N}_{pj} \leftrightarrow p_j \in \mathcal{N}_{pk}$

(S, P) constitute a graph where P contains the nodes of the graph and S the links. An image can be seen also as a graph.

Depending on the distance from p, different neighboring systems can be defined:

|   | o |   |
|---|---|---|
| o | x | o |
|   | o |   |

First order neighboring System
4-neighboring System

| o | o | o |
|---|---|---|
| o | x | o |
| o | o | o |

Second order neighboring System
8-neighboring System

36 of 72

## Clique

Borrowed from phisics.



A clique *C*, for (S, P), is defined as a subset of vertices of S, an undirected graph, such that every two vertices in the subset are connected by an edge.

I can consider ordered sets of voxels, that are connected to p through S.
Types of cliques: single-site, pairs of neighboring sites, triples of neighboring sites,… up to the cardinality of $\mathcal{N}_p$

## Markov Random Field

Given (S, P) we can define a set of random values, $\{f_k(m)\}$ for each element defined by S, that is in $\mathcal{N}_p$. Therefore we define a **random field** , $\mathcal{F}$, over S:

$$\mathcal{F}(\mathcal{N}_p) = \{f_k(m) \mid m \in \mathcal{N}_p \} \; \forall p$$

Under the Markovian hypotheses:

$P(f(p)) \geq 0 \; \forall p$           Positivity

$P(f(p) \mid g(P-\{p\}) = P(f(p) \mid g(\mathcal{N}_p)\}$      Markovianity

2 expresses the fact that the probability of p assuming a certain value, f (e.g. a certain gradient), is the same considering in p all the pixel of P but p, or only the neighbor pixels, that is the value of f depends only on the value of the pixels in $\mathcal{N}_p$ and not in p.

the random field $\mathcal{F}$ is named **Markov Random Field**.

## Energy in a Markov Random Field

A "potential" function, $\phi(f)$, can be defined for a MRF. This is a scalar value that is a function of the random value associated to the pixels for all the possible elements of a clique:

$$\phi_c(f) = \sum_{j \in c} f(p_j)$$

If we consider all the possible cliques defined for each element p, we can define a potential energy function associated to the MRF:

$$U(f) = \sum_{c \in C} \phi_c(f)$$

The higher is the potential energy, the lower is the probability that the set of random values of the elements of the cliques is realized, that is the higher is the penalization for the associated configuration.

We want to go towards minimum energy.

## Gibbs prior

If we consider all the possible cliques defined for each element p, we can define a potential energy function associated to the MRF:

$$U(\mathbf{f}) = \sum_{c \in C} \phi_c(\mathbf{f})$$

The higher is the potential energy, the lower is the probability that the set of random values of the elements of the cliques is realized, that is the higher is the penalization for the associated configuration.

This is well captured by the Gibbs distribution, that describes the probability of a certain configuration to occur. It is a function exponentially decreasing of U:

$$P(\mathbf{f}) = \frac{1}{Z} e^{\left\{-\frac{1}{\beta} U(\mathbf{f})\right\}}$$

P(f) is a Gibbs random field, Hammersley-Clifford theorem (1971). $\beta$ regulates the decrease in probability and it is associated with temperature in physics. Z is a normalization constant. NB to define Gibbs random fields, P(f) > 0, P(f) → 0 U(f) → ∞: there are not configurations with 0 probability.

## Gibbs priors and Regularization

$$\arg\min_f -\{\ln(p(g_n\,|\,f)p_f)\} = \arg\min_f -\{\ln(p(g_n\,|\,f)) + \ln(p_f)\}$$

Likelihood =
adherence to the data

A-priori

Gaussian $\quad K(\sigma)\sum_i \left\|g_{n,i} - Af_i\right\|^2$

$$-\ln\left\{\frac{1}{Z}e^{\left\{-\frac{1}{\beta}U(\mathbf{f})\right\}}\right\}$$

Poisson $\quad \sum_i g_{n,i}\ln\left(\frac{g_{n,i}}{Af} + Af_i - g_{n,i}\right)$

$$\boxed{J(f) = J_o(f) + \lambda J_R(f)}$$

$J_R(f) = U(f)$

41 of 72

---

## Role of λ

$$K(\sigma)\sum_i \left\|g_{n,i} - Af_i\right\|^2$$

$$-\ln\left\{\frac{1}{Z}e^{\left\{-\frac{1}{\beta}U(\mathbf{f})\right\}}\right\}$$

$$\boxed{J(f) = J_o(f) + \lambda J_R(f)}$$

λ incorporates different elements here:
- the standard deviation of the noise in the likelihood
- the "temperature", that is the decrease in the energy of the configurations with their cost (β)
- the normalized constant Z.

λ has been investigated in the classical regularization theory (Engl et al., 1996), but not as deep in the Bayesian framework ➔ λ is set experimentally through cross-validation.

42 of 72

# How to set the regularization parameter

Analysis of the residual after the estimate = **Af** − **g**
- The residual should be equal to the noise distribution

**Gaussian case:**
- $\lambda$ is increased until $(r_i, r_j) = \Sigma^2$   ($||r||^2 = \sigma^2$)
- Sample covariance is equal to distribution covariance

**Poisson case:**
- $r_i$ tends to be larger, the larger is $g_i$.
- $\lambda$ is increased until $|r|^2 / g \rightarrow 1$

# Choice of the Gibbs priors

We choosed $|| \lambda Pf ||^2$ as a quadratic functional, but not specified P.

P is ofted chosen as a smoothing operator. The rationale is that the noise added to the image is often white (both Gaussian and Poisson) over the image as there is no correlation between adjacent pixels. Therefore its spatial content is unform and with a larger bandwidth that the signal.

As a smoothing operator P is often a differential operator, which penalizes edges.

$$J_R(\mathbf{f}) = \sum_{c \in C} \phi_c(d^k{}_c \mathbf{f})$$

 k is the order of the derivative
 $\phi_c$ can be $l_2$ norm (total variation), squared (Tikhonov)

k = 2 difference of gradients → piecewise linear areas.
k = 3 difference of Hessian → piecewise squared.
Neighbor of order higher than 2.

## Quadratic Priors with k = 0

k = 0 – No derivative, the same gray level – single site cliques.

$$J_R(\mathbf{f}) = \sum_{c \in C} \phi(d^k_c \mathbf{f}) = \sum_{c \in C} (d^0_c \mathbf{f})^2 = \sum_{p \in P} \mathbf{f}(p)^2$$

It has been applied to both Poisson and Gaussian noise models

Reduces bright spots and biases the solution to low intensity values.

## Quadratic Priors with k = 1

k = 1 – First order derivatives – pair-sites cliques.

$$J_R(\mathbf{f}) = \sum_{c \in C} \phi(d^1_c \mathbf{f}) = \sum_{p \in P} \sum_{m \in \mathcal{N}_p} \phi(d^0_c \mathbf{f})^2 = \sum_{p \in P} \sum_{m \in \mathcal{N}_p} \phi\left( \frac{f(p) - f(m)}{d(p,m)} \right)$$

d(p,m) takes into account anisotropies in computing the distance.

If we consider φ(.) a squared function, we have another form of Tikhonov regularization:

$$J_R(\mathbf{f}) = \sum_{p \in P} \sum_{m \in \mathcal{N}_p} \left( \frac{f(p) - f(m)}{d(p,m)} \right)^2$$

## Quadratic Priors with k = 1

k = 1 – First order derivatives – pair-sites cliques.

$$J_R(\mathbf{f}) = \sum_{p \in P} \sum_{m \in \mathcal{N}_p} \left( \frac{f(p) - f(m)}{d(p,m)} \right)^2$$

If we consider $\phi(.)$ a squared function, we have another form of Tikhonov regularization:

$\| Pf \|^2$

P is the convolution with the Laplacian operator:

$$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 0 \end{bmatrix}$$

First order neighboring System
4-neighboring System

$$\begin{bmatrix} -\frac{\sqrt{2}}{2} & -1 & -\frac{\sqrt{2}}{2} \\ -1 & 4+2\sqrt{2} & -1 \\ -\frac{\sqrt{2}}{2} & -1 & -\frac{\sqrt{2}}{2} \end{bmatrix}$$

Second order neighboring System
8-neighboring System

47 of 72

## Non-quadratic potential functions, k =

Quadratic functions priors imposes smoothness everywhere. Large true gradients of the solution are therefore penalized → smoothing sharp edges.

In imaging objects tend to be piecewise smooth, but different pieces of objects are separated by more or less sharp edges. We want to smooth inside the object but not the edge. A parallel worthwhile to be investigated is with anisotropic diffusion (Koenderink, 1987; Perona&Malik, 1990).

We search different potential functions (Geman&McClure, 85; Charbonnier et al., 1994, 1997; Hebert&Lehay, 1989).

48 of 72

## Non-quadratic potentials
## (Charbonier et al., 1997)

1. $\phi(t) \geq = 0 \quad \forall t \qquad \phi(0) = 0$ — Derives from the definition of potential

2. $\Phi'(t) \geq = 0 \quad \forall t$ — Semi-monotone derivatives

3. $\phi(t) = \phi(-t)$ — Positive and negative gradients are equally considered

4. $\phi(t) \in C^1$ — This is to avoid instability.

*Up to now quadratic potentials are OK*

5. $\dfrac{\varphi'(t)}{2t}$ — The potential increase rate should decrease with t.

6. $\lim\limits_{t \to \infty} \dfrac{\varphi'(t)}{2t} = 0$ — The potential increase rate should decrease for all t (at least for large values of t)

7. $\lim\limits_{t \to 0} \dfrac{\varphi'(t)}{2t} = \cos t > 0$ — The potential increases at least linearly for t = 0.

## Few non-quadratic functions
## (Vicedomini 2008)

| Regularization name | Potential function | Expression of $\varphi(t)$ | Expression of $\psi(t) = \varphi'(t)/2t$ | Convex |
|---|---|---|---|---|
| Quadratic-Potential | $\varphi_{QP}$ | $t^2$ | $1$ | yes |
| Geman-McClure | $\varphi_{GM}$ | $\dfrac{t^2}{1+t^2}$ | $\dfrac{1}{(1+t^2)^2}$ | no |
| Hebert-Leahy | $\varphi_{HL}$ | $\log(1+t^2)$ | $\dfrac{1}{1+t^2}$ | no |
| Huber | $\varphi_{HB}$ | $\begin{cases} t^2, & |t| \leq 1 \\ 2|t|-1, & |t| > 1 \end{cases}$ | $\begin{cases} 1, & |t| \leq 1 \\ 1/|t|, & |t| > 1 \end{cases}$ | yes |
| Hyper-Surface | $\varphi_{HS}$ | $2\sqrt{1+t^2} - 2$ | $\dfrac{1}{\sqrt{1+t^2}}$ | yes |

Asymptotic linear behavior

Asymptotic log-like behavior

Why not simply $\sqrt{t^2}$ ?

## Results



## Summary

MAP estimate can be seen as a statistical version of regularization.

**The regularization term** can be derived from the potential energy associated to an adequate neighbor system defined over the object (e.g. over the image).

Under this hypothesis the value assumed by the elements of the object to be reconstructed (e.g. restored or filtered image) represent a MRF.

Different neighbor systems and different potential functions allow defining different properties of the object.

For quadratic potential functions, Tikhonov regularizer are derived.

**The discrepancy term** for the data represents the likelihood and can accommodate different statistical models: Poison, Gaussian or even mixture models.

## Overview

Filtering images

MAP, Tikhonov and Poisson model of the noise

A-priori and Markov Random Fields

Cost function minimization

## Regularization term

$$J_{REG}(f) = \left\| \nabla f \right\|_2^q$$

For q = 1, it has a singularity in the origin for which its derivative cannot be computed. Solution is one of the potentials functions above, or a numerical solution:

$$J_{REG}(f_i) = \sqrt{\frac{df_i}{dx} + \frac{df_i}{dy} + .... + \varepsilon}$$

$\varepsilon = 2.22 \times 10^{-16}$

# Simulated images

Original, unnoisy

Original, noisy

# Gradient Descendent is slow

Algorithm

Set $\qquad u^{(0)} = \{g\}$

Compute $\quad \nabla J = \left[ \dfrac{\partial}{\partial u_1} J, ..., \dfrac{\partial}{\partial u_N} J \right]^T$

Update $\qquad u^{(k+1)} = u^{(k)} - \eta \nabla J$

η is a scalar parameter (damping factor), optimized at each iteration, such as it is guaranteed that J decreases (line search).

▸ Time expensive: ~ 210s (with Matlab) on 500x500 images

▸ We can improve the algorithm and / or the gradient computation

## One-step late EM (Green, 1990)

We derive it with fixed point optimization. Let us consider the cost function for Poisson noise:

$$J\left(g_{n,i} \mid g_i\right) = -\sum_{i=1}^{N}\left\{g_{n,i}\ln(g_i) - g_i\right\} + \lambda\sum_{i=1}^{N}\left\|\nabla g_i\right\|_2^2$$

We suppose all the pixel constant and the variation of each pixel are accumulated and applied to the next step (one-step late).

$$\frac{\partial J\left(g_{n,k} \mid g_k\right)}{\partial g_k} = \frac{\partial}{\partial g_k}\left\{-\left[g_{n,k}\ln(g_k) - g_k\right]\right\} + \lambda\cdot\frac{\partial}{\partial g_k}J_R(g_k) = -\frac{g_{n,k}}{g_k} + 1 + \lambda\cdot\frac{\partial}{\partial g_k}J_R(g_k) = 0$$

This cannot be solved directly, but it can be solved using fixed point iteration:

$$-\frac{g_{n,k}}{g_k} + 1 + \lambda\cdot\frac{\partial}{\partial g_k}J_R(g_k) = 0 \Rightarrow \frac{g_{n,k}}{g_k} = 1 + \lambda\cdot\frac{\partial}{\partial g_k}J_R(g_k) \Rightarrow g_k = \frac{g_{n,k}}{1 + \lambda\cdot\dfrac{\partial}{\partial g_k}J_R(g_k)}$$

## Expectation Maximization

From emission Tomography (Green, 1990; Panin et al., 1999)

$$u_i^{(new)} = \frac{u_i^{(old)}}{\sum\limits_j h_{i,j} + \lambda\dfrac{\partial}{\partial u_i}J_{REG}\left(u^{(old)}\right)}\sum_j\frac{h_{i,j}z_j}{\sum\limits_k h_{k,j}u_k^{(old)}}$$

In our case

$$H = \left[h_{i,j}\right] = I$$

The previous formula becomes

$$u_i^{(new)} = \frac{z_i}{1 + \lambda\dfrac{\partial}{\partial u_i}J_{REG}\left(u^{(old)}\right)}$$

## Observations

Semi-convergence properties.

Damping of the solution is required.

▸ Damped EM, $x^{k+1}=(1-t)x^k+t*EM(x^k)$ (damping, relaxation, reduction of the step length)

Solutions have been recently proposed for PET images (Mair&Zahnen, 2006).

Large increase in speed has been registered.

Sensitive to number of steps.



## Centered gradient is bad

$$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 0 & +1 \\ 0 & +1 & 0 \end{bmatrix}$$

If used centered gradient to computer the a-priori, we obtain a checkerboard effect



60 of 72

30

## Different gradient possibilities

We consider only two gradients: North-Center + West-Center

$$\left\|\nabla g(x_i, y_i)\right\|_2 = \sqrt{g_x(x_i, y_i)^2 + g_y(x_i, y_i)^2 + \varepsilon} =$$

$$= \sqrt{[g(x_i, y_i) - g(x_i - 1, y_i)]^2 + [g(x_i, y_i) - g(x_i, y_i - 1)]^2 + \varepsilon}$$

4 neighbors gradient

8 neighbors gradient

## Why not to change the norm?

We consider only two gradients: North-Center + West-Center

$$\left\|\nabla g(x_i, y_i)\right\|_1 = \left|g_x(x_i, y_i)\right| + \left|g_y(x_i, y_i)\right| = \left|g(x_i, y_i) - g(x_i - 1, y_i)\right| + \left|g(x_i, y_i) - g(x_i, y_i - 1)\right|$$
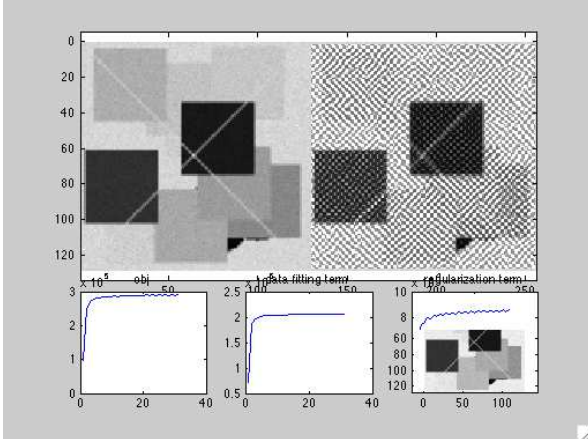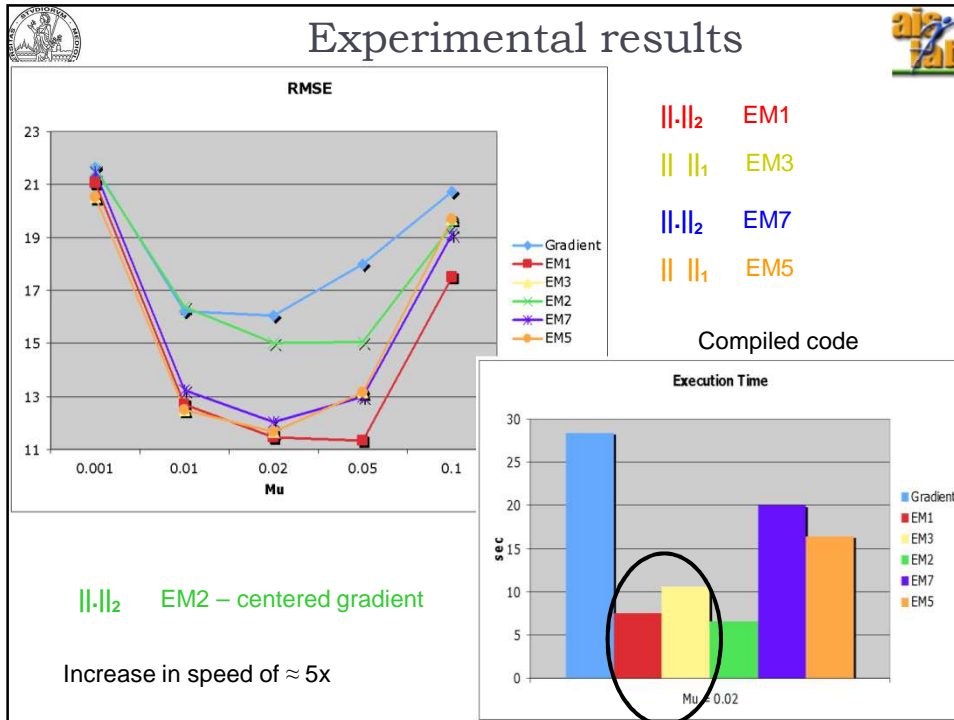
$$\frac{\partial J_R(\mathbf{g})}{\partial g_k} = \frac{\partial \sum_{i=1}^{N} \left\|\nabla g(x_i, y_i)\right\|_1}{\partial g_k} = \frac{\partial \left[\left\|\nabla g(x_k, y_k)\right\|_1 + \left\|\nabla g(x_k + 1, y_k)\right\|_1 + \left\|\nabla g(x_k, y_k + 1)\right\|_1\right]}{\partial g_k} =$$

$$= \frac{\partial}{\partial g_k}\left[\left|g(x_k, y_k) - g(x_k - 1, y_k)\right| + \left|g(x_k, y_k) - g(x_k, y_k - 1)\right|\right] +$$

$$\frac{\partial}{\partial g_k}\left[\left|g(x_k + 1, y_k) - g(x_k, y_k)\right| + \left|g(x_k + 1, y_k) - g(x_k + 1, y_k)\right|\right] +$$

$$\frac{\partial}{\partial g_k}\left[\left|g(x_k, y_k + 1) - g(x_k - 1, y_k + 1)\right| + \left|g(x_k, y_k + 1) - g(x_k, y_k)\right|\right] =$$

$$= \boxed{sign[g_x(x_k, y_k)] + sign[g_y(x_k, y_k)] - sign[g_x(x_k + 1, y_k)] - sign[g_y(x_k, y_k + 1)]}$$

We do not need ε anymore but we do not have continuity in the origin. May be we can relax Charbonnier et al. conditions….

## Experimental results



||·||₂ → $\|\cdot\|_2$   EM1

|| ||₁ → $\|\ \|_1$   EM3

||·||₂ → $\|\cdot\|_2$   EM7

|| ||₁ → $\|\ \|_1$   EM5

Compiled code

$\|\cdot\|_2$   EM2 – centered gradient

Increase in speed of ≈ 5x

---

## Beyond EM

$$J\left(g_{n,i} \mid g_i\right) = -\sum_{i=1}^{N}\left\{g_{n,i}\ln\left(g_i\right) - g_i\right\} + \lambda\sum_{i=1}^{N}\left\|\nabla g_i\right\|_2^{q}$$

is an optimization problem, in which g has two interesting properties:

$g(p) \geq 0$

$$\sum_{p} g(p) = \cos t$$

Flux conservation (preservation of the intensity of the image)

Moreover, J(.) is supposed convex. Under these hypotheses, the so
Called Kuhn-Tucker condition for the (unique) minimum should hold:

g*∇J(g*; gₙ) = 0

g* ≥ 0      ∇J(g*; gₙ) ≥ 0

## Split gradient (Lanteri, 2002)

$$J\left(g_{n,i} \mid g_i\right) = -\sum_{i=1}^{N}\left\{g_{n,i}\ln(g_i) - g_i\right\} + \lambda\sum_{i=1}^{N}\left\|\nabla g_i\right\|_2^q$$

Singularity when gradient is 0 and q < 2.
The idea is to obtain a term > 0 strictly at the denominator.

$\nabla J(g; gn) = U(g; gn) + V(g; gn)$   with  $U(g; gn) \geq 0$; $V(g; gn) > 0$

Kuhn-Tucker condition becomes:

$g^*\nabla J(g^*; gn) = 0$   ➔ $g^*(U(g; gn) + V(g; gn)) = 0$

We can write fixed point iteration and obtain:

$g^{(t+1)} = g(t)\ U(g; gn)\ /\ V(g; gn))$

## Split-gradient Algorithm

**Inizialization**. Choose $g^{(0)}$, that can be coincident with $g_n$ and compute the flux, that is the $c = \Sigma g_{n,i}$.

**Iteration** in two steps: update + normalization.

Update:
$$\hat{g}^{(t+1)} = g^{(t)} + a^{(t)}g^{(t)}\left(\frac{U(g; g_n) - V(g; g_n)}{V(g; g_n)}\right)$$

$$c^{(t+1)} = \sum_p g^{(t+1)}(p)$$

Normalization through flux conservation:

$$g^{(t+1)}(p) = \frac{c}{c^{(t+1)}}\hat{g}^{(t+1)}(p)$$

## Relaxed Split-gradient Algorithm
## ($\alpha = 1$)

**Inizialization**. Choose $g^{(0)}$, that can be coincident with $g_n$ and compute the flux, that is the $c = \Sigma g_{n,i}$ .

**Iteration** in two steps: update + normalization.

Update: $\quad \hat{g}^{(t+1)} = g^{(t)} + a^{(t)} g^{(t)} \left( \dfrac{U(g; g_n) - V(g; g_n)}{V(g; g_n)} \right) = g^{(t)} \left( \dfrac{U(g; g_n)}{V(g; g_n)} \right)$

$$c^{(t+1)} = \sum_p g^{(t+1)}(p)$$

Normalization through flux conservation:

$$g^{(t+1)}(p) = \frac{c}{c^{(t+1)}} \hat{g}^{(t+1)}(p)$$

that has a very attractive multiplicative factor. This is also a Scaled gradient algorithm (Bertero et al., 2008)

## Determination of U(.) and V(.)

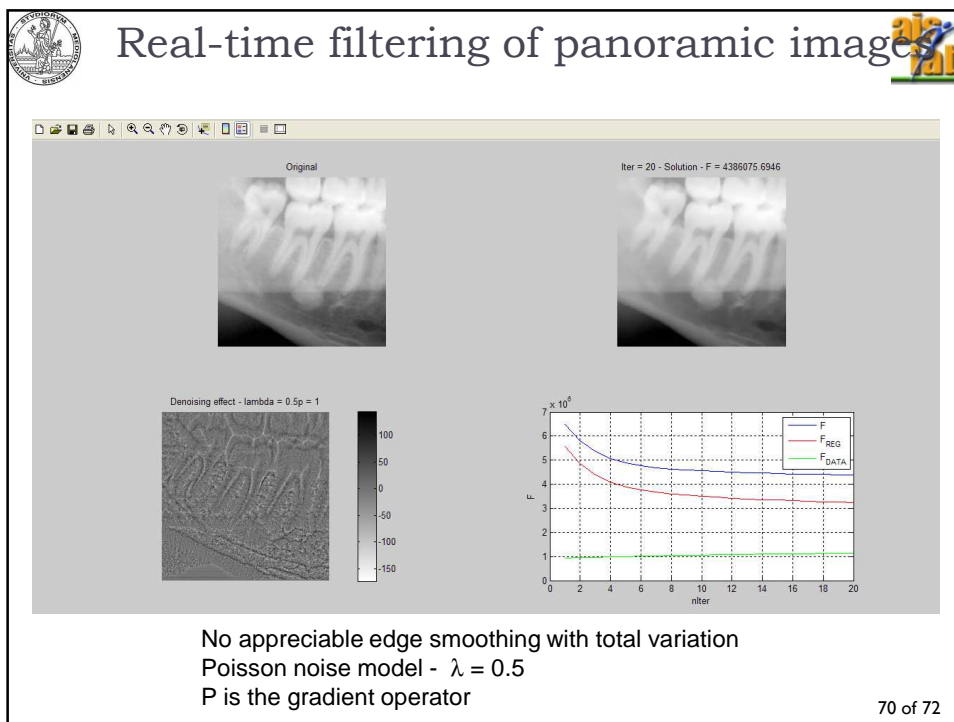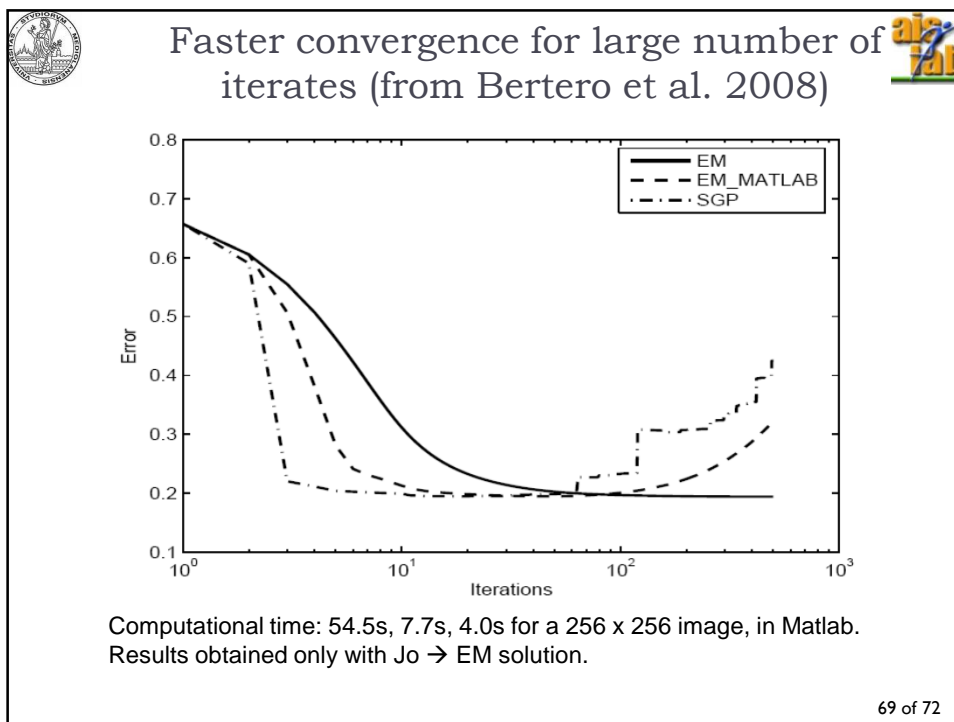$$J\left(g_{n,i} \mid g_i\right) = -\sum_{i=1}^{N} \left\{g_{n,i} \ln(g_i) - g_i\right\} + \lambda \sum_{i=1}^{N} \left\| \nabla g_i \right\|_2^q = J_o + \lambda J_R$$

For the likelihood term: $\nabla J_0$

|  | U | V |
|---|---|---|
| Gaussian case | $2g_n$ | $2g$ |
|  | $2A^T g_n$ | $2(A^T Ag + b)$ |
| Poisson case | $g_n / g$ | 1 |
|  | $A^T g_n / (Ag + b)$ | |

For the regularization term: $\nabla J_R$ the derivatives of the potential function have to be considered (Bertero et al., in preparation) and grouped into positive and strictly positive values.

## Faster convergence for large number of iterates (from Bertero et al. 2008)



Computational time: 54.5s, 7.7s, 4.0s for a 256 x 256 image, in Matlab.
Results obtained only with Jo → EM solution.

## Real-time filtering of panoramic images



No appreciable edge smoothing with total variation
Poisson noise model - $\lambda = 0.5$
P is the gradient operator

## Application for intensive algebraic methods

Denoising – Bayesian filtering

Deconvolution (tomosynthesis, volumetric reconstruction from limited angle of view)

Deconvolution (CB-CT, FanBeam CT)

….

Amenable to be implemented on CUDA architectures → Real-time volumetric reconstruction.

## Overview

Filtering images

MAP, Tikhonov and Poisson model of the noise

A-priori and Markov Random Fields

Cost function minimization