

Reformulation of the model using curvelets

F. Parzer

October 29, 2020

Abstract

OVERVIEW

Section 1 contains a short summary of the relevant informations of the problem as we understood it.

In Section 2, a formulation of the model is presented where the distribution function of the galaxy is approximated using curvelets.

Finally, in Section 3, the optimization problem is formulated. We also present the computational aspects and discuss how approximate uncertainty quantification in form of covariance estimates and confidence intervals could be performed.

CONTENTS

1	Formulation of the model	2
1.1	Problem summary	2
1.2	The observation operator	3
1.3	The measurement error	3
1.4	The prior distribution	3
2	Discretization of the distribution function	4
2.1	Curvelet preliminaries	4
2.2	The curvelet model	7
2.3	The finite-dimensional approximate model	7
3	The computational approach	8
3.1	Transform into nonlinear least-squares problem	8
3.2	The constrained Gauss-Newton method	10
3.3	Uncertainty quantification	10
4	What remains to be done	12
4.1	Open questions	12
4.2	Wishlist	12
	Bibliography	13

1. FORMULATION OF THE MODEL

1.1. Problem summary. The problem at hand is to estimate the line-of-sight velocity distribution of a galaxy.

A galaxy is modelled as a collection G of *stellar populations*. Each stellar population consists of a number of stars and is characterized by two parameters. The first parameter is the *metallicity* $m_p \in \mathbb{R}$. As the logarithm of a ratio, it is a real number. The second parameter is the age $a_p > 0$ of the stellar population in Gyr (Gigayears). In the *MILES simple stellar population model*, these two numbers determine the wavelength spectrum of the stellar population. A wavelength spectrum is a function $s : I \times \Theta \rightarrow \mathbb{R}^+$, where $I \subset \mathbb{R}^+$ is the range of relevant wavelengths and $\Theta \subset \mathbb{R} \times \mathbb{R}^+$ is the space of possible values of the parameter $\theta_p = (m_p, a_p)$.

By the *distribution function of the galaxy* we mean the histogram of the parameter θ with respect to stellar mass, i.e. a function $f : \Theta \rightarrow \mathbb{R}^+$ that maps to each parameter value θ the total stellar mass of all stellar populations with the corresponding value. The function f necessarily satisfies

$$f(\theta) \geq 0 \text{ and } \int_{\Theta} f(\theta) d\theta = M_{\text{total}}, \quad (1.1)$$

where M_{total} is the total stellar mass of the galaxy. From practical experience, f is typically a multimodal function which is characterized by a few sickle-shaped "hills" and otherwise close to zero. This corresponds to the physical fact that the distribution of the stellar mass is concentrated at a few clusters with similar age and metallicity. Thus, we expect that in a localized basis or frame of functions (see [Section 2.1](#)), f can be well-approximated by only a few basis elements. That is, we expect *sparsity* in the discretization.

The *composite spectrum* of the stellar population of a galaxy is then given through a superposition of the spectra of all stellar populations:

$$S(\lambda; f) = \int_{\Theta} f(\theta) s(\lambda; \theta) d\theta. \quad (1.2)$$

The Doppler effect shifts the observed spectrum. The size of this shift is determined by the component of the velocity of the stellar population along the line-of-sight from that population to earth. Therefore, in order to know the unshifted spectrum, we have to estimate the *line-of-sight velocity distribution* (LOSVD) of the galaxy. This is a collection $(v_p)_{p \in G} \subset \mathbb{R}$ of real numbers, where each v_p represents the line-of-sight velocity component of the stellar population p . One can normalize the LOSVD by estimating first the line-of-sight velocity of the galaxy (is this the average velocity of all stellar populations, or something else?), and then shift the LOSVD by this amount. Thus, the LOSVD is centered around zero.

Observed LOSVDs are typically modelled by *Gauss-Hermite expansions*. We define

$$L(v, \theta_v) := \mathcal{N}(\hat{v}; 0, 1) \left[\sum_{m=0}^M h_m H_m(\hat{v}) \right],$$

where $\mathcal{N}(\cdot; 0, 1)$ denotes the density function of the standard normal distribution, and where furthermore

$$\begin{aligned} \hat{v} &:= \frac{v - V}{\sigma}, \\ H_m(x) &:= \frac{H_m^{\text{phys}}(x)}{\sqrt{m! 2^m}}, \\ H_m^{\text{phys}}(x) &:= (-1)^m e^{x^2} \frac{d^m}{dx^m} [e^{-x^2}], \end{aligned}$$

and θ_v is the collection of parameters given by

$$\theta_v = (V, \sigma, h_0, \dots, h_M)$$

Thus, taking into account the Doppler effect and the LOSVD of the galaxy, the shifted spectrum is given by

$$\bar{y}(\lambda) = \int_{-\infty}^{\infty} \frac{1}{1 + v/c} S\left(\frac{\lambda}{1 + v/c}; f\right) L(v; \theta) dv. \quad (1.3)$$

1.2. The observation operator. We condense the above model in an observation operator

$$\mathcal{G} : \mathbb{F} \times \mathbb{R}^{M+2} \longrightarrow \mathbb{R}^{N_y}, \quad \mathcal{G}(f, \boldsymbol{\theta}_v) = \bar{\mathbf{y}}(f, \boldsymbol{\theta}_v),$$

which maps the parameters f and $\boldsymbol{\theta}_v$ to the discrete, shifted composite spectrum $\bar{\mathbf{y}} \in \mathbb{R}^{N_y}$. Here, N_y is the number of wavelet bins, and the vector $\bar{\mathbf{y}}$ is obtained from the function \bar{y} through a discretization operation \mathbf{P}_{N_y} , i.e.

$$\bar{\mathbf{y}} = \mathbf{P}_{N_y}(\bar{y}).$$

We assume that the distribution function f lies in function space \mathbb{F} , given by

$$\mathbb{F} = \left\{ f \in L^2(\Theta) : f \geq 0, \int_{\Theta} f(\theta_s) d\theta_s = M_{\text{total}} \right\},$$

where $\Theta \subset \mathbb{R} \times \mathbb{R}^+$ is the range of possible values of the metallicity and age.

The continuous spectrum \bar{y} is computed from

$$\bar{y}(\lambda; f, \boldsymbol{\theta}_v) = \int_{-\infty}^{\infty} \frac{1}{1+v/c} S\left(\frac{\lambda}{1+v/c}; f\right) L(v; \boldsymbol{\theta}) dv,$$

where

$$S(\lambda; f) = \int_{\Theta} f(\theta) s(\lambda; \boldsymbol{\theta}) d\boldsymbol{\theta},$$

and

$$L(v, [V, \sigma, h_0, \dots, h_M]) := \mathcal{N}\left(\frac{v-V}{\sigma}; 0, 1\right) \left[\sum_{m=0}^M h_m H_m\left(\frac{v-V}{\sigma}\right) \right].$$

If we insert these expressions in the formula for \bar{y} , we obtain an explicit expression for $\mathcal{G}(f, \boldsymbol{\theta}_v)$:

$$\mathcal{G}(f, \boldsymbol{\theta}_v) = \int_{-\infty}^{\infty} \frac{1}{1+v/c} \left(\int_{\Theta} f(\boldsymbol{\theta}) s\left(\frac{\lambda}{1+v/c}; \boldsymbol{\theta}\right) d\boldsymbol{\theta} \right) \mathcal{N}\left(\frac{v-V}{\sigma}; 0, 1\right) \left(\sum_{m=0}^M h_m H_m\left(\frac{v-V}{\sigma}\right) \right) dv. \quad (1.4)$$

1.3. The measurement error. We will assume that there is only additive measurement error, i.e. that the actual measurement is given by

$$\mathbf{y} = \mathcal{G}(f, \boldsymbol{\theta}_v) + \boldsymbol{\epsilon}.$$

We will furthermore assume that

- $\boldsymbol{\epsilon}$ is multivariate normally distributed (on \mathbb{R}^{N_y}),
- uncorrelated per wavelength bin,
- and that the expected noise level is equal for each wavelength bin.

That is, we assign the noise distribution

$$\boldsymbol{\epsilon} \sim \mathcal{N}(0, \delta^2 \mathbf{I}_{N_y}), \quad (1.5)$$

where $\mathbf{I}_{N_y} \in \mathbb{R}^{N_y \times N_y}$ is the N_y -dimensional identity matrix and $\delta > 0$ is the noise level. Typically, the signal-to-noise ratio lies in the range $[30, 300]$.

1.4. The prior distribution. The expected sparsity of the distribution function f is represented in our model by assigning a sparsity-promoting prior. We choose a Laplace-like prior that also accounts for the additional constraints (1.1) that f has to satisfy. Formally, we would like to write this in terms of a probability density

$$f \sim \frac{1}{C} \exp\left(-\sqrt{2} \left\| \Sigma_1^{-1/2} f \right\|_1\right),$$

and then restrict this function to only those values that satisfy the necessary conditions (1.1) (recall that the factor $\sqrt{2}$ is just a scaling factor that ensures that the distribution has covariance Σ_1 and not $\frac{1}{2}\Sigma_1$).

However, since f is an infinite-dimensional object and probability density functions on infinite-dimensional spaces are not well-defined, the only way to make rigorous sense of a prior probability on the distribution function $f : \Theta \rightarrow \mathbb{R}^+$ is as a probability measure on the infinite-dimensional function space $L^2(\Theta)$. In our particular case, this probability measure is given by

$$\mu_0(df \mid \Sigma_1) = \frac{1}{C_0} \text{Lap}(df \mid 0, \Sigma_1) \mathbb{1}_{\mathbb{F}}(f), \quad (1.6)$$

where $\text{Lap}(df \mid 0, \Sigma_1)$ is the infinite-dimensional Laplace-measure on $L^2(\Theta)$ with mean 0 (in order to promote sparsity) and covariance operator Σ_1 (see `IEKF11.pdf` for the exact construction of the infinite-dimensional Laplace measure). The indicator function $\mathbb{1}_{\mathbb{F}}(f)$ is equal to 1 if $f \in \mathbb{F}$, and 0 else. Consequently, it projects the prior probability on the constraints (1.1).

This measure-theoretic formulation can be treated as a formal detail and does not affect our actual implementation, since we will only deal with finite-dimensional discretizations of f . However, it might be relevant for rigorous interpretation of some results.

For the other parameter θ_v , which is given by the vector $(V, \sigma, h_0, \dots, h_M)$, we will assume a simple uncorrelated Gaussian prior, i.e.

$$\theta_v \sim \mathcal{N}(\bar{\theta}_v, \Sigma_2), \quad (1.7)$$

where $\bar{\theta}_v = (\bar{V}, \bar{\sigma}, \bar{h}_0, \dots, \bar{h}_M)$ is the prior mean and $\Sigma_2 = \text{diag}(\delta_v, \delta_{h_0}, \dots, \delta_{h_M}) \in \mathbb{R}^{(M+2) \times (M+2)}$ is a diagonal covariance matrix. The parameters $\bar{\theta}_v$ and Σ_2 have to be chosen by us, and we should evaluate in numerical experiments how this choice affects our inference.

If we assign the priors (1.6) and (1.7) and use the assumption (1.5), then the statistical formulation of our inverse problem is

$$\mathbf{y} = \mathcal{G}(f, \theta_v) + \epsilon, \quad (1.8)$$

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \delta^2 \mathbb{I}_{N_y}), \quad (1.9)$$

$$f \sim \mu_0(df \mid \Sigma_1), \quad (1.10)$$

$$\theta_v \sim \mathcal{N}(\bar{\theta}_v, \Sigma_2), \quad (1.11)$$

$$f \perp\!\!\!\perp \theta_v, \quad (1.12)$$

where $f \perp\!\!\!\perp \theta_v$ means that we assume that the unknown parameters f and θ_v are independent.

2. DISCRETIZATION OF THE DISTRIBUTION FUNCTION

As mentioned above, the distribution function f is expected to have a particular shape: It is negligibly small on most of its domain, except for a few curved hills in which most of its mass is concentrated. Because of this fact, in our first meeting we came up with the idea of finding a sparse approximation of f using *curvelets*.

2.1. Curvelet preliminaries. Curvelets are a useful modification of wavelets for sparse geometrical image representation. Their advantage over wavelets is that they allow for a better representation of directional components of images, i.e. curved edges. For some exemplary applications of curvelets to astronomy, see [15] and [9].

Mathematical definition

Mathematically, curvelets yield a discretization of the function space $L^2(\mathbb{R}^2)$. This discretization is particularly suited for the approximation of images that are smooth except for discontinuities along C^2 curves. For a very good review on curvelets in signal processing, see [11].

We define the curvelet dictionary as follows: Consider the so-called *scaled Meyer windows*

$$V(t) := \begin{cases} 1, & |t| \leq \frac{1}{3}, \\ \cos \left[\frac{\pi}{2} \nu(3|t| - 1) \right], & \frac{1}{3} \leq |t| \leq \frac{2}{3}, \\ 0, & \text{else, } t \in \mathbb{R}, \end{cases} \quad (2.1)$$

and

$$W(r) := \begin{cases} \cos \left[\frac{\pi}{2} \nu(5 - 6r) \right], & \frac{2}{3} \leq r \leq \frac{5}{6}, \\ 1, & \frac{5}{6} \leq r \leq \frac{4}{3}, \\ \cos \left[\frac{\pi}{2} \nu(3r - 4) \right], & \frac{4}{3} \leq r \leq \frac{5}{3}, \\ 0, & \text{else,} \end{cases} \quad (2.2)$$

where $\nu \in C^\infty(\mathbb{R})$ can be any smooth function that satisfies

- $\nu \equiv 0$ on $(-\infty, 0]$ and $\nu \equiv 1$ on $[1, \infty)$,
- $\nu(x) + \nu(1 - x) = 1$, for all $x \in \mathbb{R}$.

Next, we define the *a-scaled window*

$$U_a(r\omega) := a^{3/4} W(ar) V\left(\frac{\omega}{\sqrt{a}}\right), \quad a \in [0, 1], \quad r \geq 0, \quad \omega \in [0, 2\pi).$$

This window is used as a function in polar coordinates on the frequency domain, i.e. for a frequency $\xi = [\xi^1, \xi^2]^\top \in \mathbb{R}^2$ we define

$$r := \sqrt{|\xi^1|^2 + |\xi^2|^2}, \quad \text{and} \quad \omega := \arctan\left(\frac{\xi^1}{\xi^2}\right).$$

Then, we define the base curvelet in frequency domain as

$$\hat{\varphi}_{a,0,0}(\xi) := U_a(\xi), \quad \xi \in \mathbb{R}^2. \quad (2.3)$$

Via rotation and translation, this generates the curvelet family

$$\varphi_{a,b,\theta}(x) := \varphi_{a,0,0}(R_\theta(x - b)), \quad (2.4)$$

where

$$R_\theta := \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}, \quad \theta \in [0, 2\pi), \quad (2.5)$$

is the 2-dimensional rotation around the angle θ .

The discrete curvelet frame

The curvelet family $(\varphi_{a,b,\theta})$ defines a frame of $L^2(\mathbb{R}^2)$. Recall that frames generalize the concept of a basis to not necessarily linearly independent sequences. Thus, they allow for *redundancy* in the representation of an object. Frames are defined as follows (see, for example, [12]):

Definition 2.1 Let \mathbb{H} be a Hilbert space. A sequence $(\phi_i)_{i \in I}$ is a *frame* of \mathbb{H} if there are two constants $A, B > 0$, with $A \leq B$, such that

$$A \|f\|^2 \leq \sum_{i \in I} |\langle f, \phi_i \rangle|^2 \leq B \|f\|^2, \quad \forall f \in \mathbb{H}. \quad (2.6)$$

If $A = B$, then the frame is called *tight*.

Since the frame constant A must be strictly positive, condition (2.6) implies that the cardinality of the index set I cannot be smaller than the dimension of the space \mathbb{H} . Frames generalize the concept of a basis in the further sense that they allow for coefficient expansions of elements of \mathbb{H} . The next result can be found as theorem 5.5 in [12].

Theorem 2.2 Let $(\phi_n)_{n \in I}$ be a frame of a Hilbert space \mathbb{H} . Then, for every $f \in \mathbb{H}$, we have

$$f = \sum_{i \in I} \langle f, \phi_i \rangle \tilde{\phi}_i = \sum_{i \in I} \langle f, \tilde{\phi}_i \rangle \phi_i,$$

where $\tilde{\phi}_i := (\Phi^* \Phi)^{-1} \phi_i$ are the elements of the so-called dual frame, and $\Phi : \mathbb{H} \rightarrow \mathbb{R}^I$ is the frame analysis operator, given by

$$(\Phi f)_i := \langle f, \phi_i \rangle.$$

If the frame is tight with constant A , then $\tilde{\phi}_i = \frac{1}{A} \phi_i$, for all $i \in I$.

In particular, this shows that any function f can be reconstructed from its frame coefficients.

Finally, the next theorem demonstrates the usual construction of a curvelet frame of $L^2(\mathbb{R}^2)$.

Theorem 2.3 Let

$$\begin{aligned} a_j &:= 2^{-j}, \\ \theta_{j,l} &:= \frac{\pi l 2^{-\lceil j/2 \rceil}}{2}, \\ b_k^{j,l} &:= R_{\theta_{j,l}}^{-1} \left[\frac{k^1/2^j}{k^2/2^{j/2}} \right], \quad k \in \mathbb{Z}^2, \quad l \in \{0, 1, \dots, 4 \cdot 2^{-\lceil j/2 \rceil} - 1\}, \quad j \in \mathbb{N}. \end{aligned}$$

We define the scaled window function U_j by

$$U_j(r, \omega) := 2^{-3j/4} W(2^{-j} r) V\left(\frac{2 - 2^{-\lceil j/2 \rceil} \omega}{\pi}\right).$$

and the discrete family of curvelets by

$$\begin{aligned} \hat{\phi}_{j,0,0}(\xi) &:= U_j(\xi), \\ \phi_{j,k,l}(x) &:= \phi_{j,0,0}(R_{\theta_{j,l}}(x - b_k^{j,l})), \\ \phi_{-1,k,0} &:= \phi_{-1}(x - k). \end{aligned}$$

Let $I_c = \{ \mathbf{i} = (j, k, l) : k \in \mathbb{Z}^2, l \in \{0, 1, \dots, 4 \cdot 2^{-\lceil j/2 \rceil} - 1\}, j \in \mathbb{N} \}$. Then $(\phi_{\mathbf{i}})_{\mathbf{i} \in I_c}$ is a tight frame of $L^2(\mathbb{R}^2)$ with frame bound $A = 1$.

For a proof of this statement and more properties of curvelets, see [8].

Since the frame bound of the curvelet frame is 1, a particularly simple coefficient expansion holds for all $f \in L^2(\mathbb{R}^2)$:

$$f = \sum_{\mathbf{i} \in I_c} \langle f, \phi_{\mathbf{i}} \rangle_2 \phi_{\mathbf{i}}, \quad \forall f \in L^2(\mathbb{R}^2), \quad (2.7)$$

where $\langle \cdot, \cdot \rangle_2$ denotes the L^2 -inner product. Furthermore, the curvelet expansion satisfies Parseval's identity,

$$\sum_{\mathbf{i} \in I_c} |\langle f, \phi_{\mathbf{i}} \rangle_2|^2 = \|f\|_2^2.$$

Similar to the fast Fourier transform, there exists a very efficient algorithm for computing the curvelet coefficients $\langle f, \phi_{\mathbf{i}} \rangle$ of a given function f , called the *fast curvelet transform*. If f is given as a discrete image $(f(n_1, n_2))$ of $N \cdot N$ pixels, then the computational complexity of the discrete curvelet transform is $O(N^2 \log N)$ [7].

2.2. The curvelet model. We can now reformulate our model from [Section 1.3](#) using curvelets. Suppose the distribution f has the curvelet expansion

$$f = \sum_{\mathbf{i} \in I_c} f_{\mathbf{i}} \phi_{\mathbf{i}}, \quad (2.8)$$

where $f_{\mathbf{i}} = \langle f, \phi_{\mathbf{i}} \rangle$. Let $\tilde{f} := \Phi f = (f_{\mathbf{i}})_{\mathbf{i} \in I_c}$. If we insert (2.8) in (1.2), we have

$$\begin{aligned} S(\lambda, f) &= \sum_{\mathbf{i} \in I_c} f_{\mathbf{i}} \underbrace{\int \phi(\theta_s) s(\lambda; \theta_s) d\theta_s}_{=: s_{\mathbf{i}}(\lambda)} \\ &=: S(\lambda, \tilde{f}). \end{aligned}$$

Then, we can rewrite the observation operator as

$$\mathcal{G}(\tilde{f}, \theta_v) = \sum_{\mathbf{i} \in I_c} f_{\mathbf{i}} \int_{-\infty}^{\infty} \frac{1}{1+v/c} s_{\mathbf{i}} \left(\frac{1}{1+v/c} \right) L(v; \theta_v) dv. \quad (2.9)$$

2.3. The finite-dimensional approximate model. Suppose next we treat f as a discrete image $\mathbf{f}[n_1, n_2]$. Correspondingly, we let Σ_1 be the corresponding projection of the (infinite-dimensional) covariance matrix on the discretization. We can compute its curvelet coefficients $\tilde{\mathbf{f}} = \Phi \mathbf{f}$ using the fast curvelet transform. Since curvelets are not necessarily nonnegative, the nonnegativity condition in (1.1) does not directly translate into a nonnegativity condition on the curvelet coefficients. Instead, we have the condition

$$\Phi^{-1} \tilde{\mathbf{f}} \geq 0,$$

where Φ^{-1} is the finite-dimensional curvelet reconstruction operator.

The integral condition in (1.1) translates to

$$\sum_{\mathbf{i}} w_{\mathbf{i}} \tilde{\mathbf{f}}_{\mathbf{i}} = M_{\text{total}}, \quad (2.10)$$

where $w_{\mathbf{i}} = \int \phi_{\mathbf{i}}(\theta_s) d\theta_s$.

Note that we have only a finite number of wavelength bins so that our measurement is actually m -dimensional, where m is the number of bins. Thus, we need no additional discretization of the measurements.

Our discretized model then looks like this:

$$\mathbf{y} = \mathcal{G}(\tilde{\mathbf{f}}, \theta_v) + \epsilon, \quad (2.11)$$

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \delta^2 \mathbb{I}_m), \quad (2.12)$$

$$\tilde{\mathbf{f}} \sim \text{Lap}(\mathbf{0}, \tilde{\Sigma}_1), \quad (2.13)$$

$$\theta_v \sim \mathcal{N}(\bar{\theta}_v, \Sigma_2), \quad (2.14)$$

$$\tilde{\mathbf{f}} \perp\!\!\!\perp \theta_v, \quad (2.15)$$

where $\tilde{\Sigma}_1 = \Phi \Sigma_1 \Phi^{\top}$ is the covariance matrix for the curvelet coefficients and $\text{Lap}(\mathbf{0}, \tilde{\Sigma}_1)$ is the Laplace measure (see (1.6)). Its probability density function is given by

$$p_0(\tilde{\mathbf{f}}) = \frac{1}{C_0} \exp \left(-\sqrt{2} \left\| \tilde{\Sigma}_1^{-1/2} \tilde{\mathbf{f}} \right\|_1 \right) \mathbb{1}_{\mathbb{F}}(\tilde{\mathbf{f}}).$$

where the set $\mathbb{F} = \{ \tilde{\mathbf{g}} = [g_1, \dots, g_M] : \Phi^{-1} \tilde{\mathbf{g}} \geq \mathbf{0}, \sum_{\mathbf{i}} w_{\mathbf{i}} \tilde{g}_{\mathbf{i}} = M_{\text{total}} \}$ is the set of all vectors $\tilde{\mathbf{g}}$ of curvelet coefficients such that $\mathbf{g} = \Phi^{-1} \tilde{\mathbf{g}}$ satisfies the constraints (1.1), and C_0 is a normalization constant.

3. THE COMPUTATIONAL APPROACH

Since any sampling-based approach to uncertainty quantification is likely to be computationally infeasible, we have to make due with mode-based inference. That means we compute point estimates for $\tilde{\mathbf{f}}$ and $\boldsymbol{\theta}_v$ through maximum-a-posteriori estimation, and perform rudimentary uncertainty quantification using mode-based estimates of the posterior covariance and confidence regions.

By Bayes' rule, the posterior probability density p_{post} for the finite-dimensional model (2.11)-(2.15) satisfies

$$p_{\text{post}}(\tilde{\mathbf{f}}, \boldsymbol{\theta}_v | y) \propto \exp \left(-\frac{1}{2\delta^2} \left\| \mathbf{y} - \mathcal{G}(\tilde{\mathbf{f}}, \boldsymbol{\theta}_v) \right\|_2^2 - \sqrt{2} \left\| \boldsymbol{\Sigma}_1^{-1/2} \tilde{\mathbf{f}} \right\|_1 - \frac{1}{2} \left\| \boldsymbol{\Sigma}_2^{-1/2} (\boldsymbol{\theta}_v - \bar{\boldsymbol{\theta}}_v) \right\|_2^2 \right) \mathbb{1}_{\mathbb{F}}(\tilde{\mathbf{f}}),$$

where \mathbb{F} is the set of feasible points defined above.

Any mode of p_{post} is a minimizer of the following constrained minimization problem

$$\min \quad \frac{1}{2\delta^2} \left\| \mathbf{y} - \mathcal{G}(\tilde{\mathbf{f}}, \boldsymbol{\theta}_v) \right\|_2^2 + \sqrt{2} \left\| \boldsymbol{\Sigma}_1^{-1/2} \tilde{\mathbf{f}} \right\|_1 + \frac{1}{2} \left\| \boldsymbol{\Sigma}_2^{-1/2} (\boldsymbol{\theta}_v - \bar{\boldsymbol{\theta}}_v) \right\|_2^2, \quad (3.1)$$

$$\text{subject to} \quad \boldsymbol{\Phi}^{-1} \tilde{\mathbf{f}} \geq 0, \quad (3.2)$$

$$\sum_i w_i \tilde{f}_i = M_{\text{total}}. \quad (3.3)$$

If we want point estimates, we have to solve this optimization problem.

3.1. Transform into nonlinear least-squares problem. In general, l^2 -regularized least-squares problems are easier to solve than l^1 -regularized least-squares problems. However, we can transform the latter into the former.

Consider a generic l^1 -regularized nonlinear least-squares problem

$$\min_{\mathbf{x}} \frac{1}{2} \left\| \mathbf{y} - \mathbf{F}(\mathbf{x}) \right\|_2^2 + \sqrt{2} \left\| \boldsymbol{\Sigma}^{-1/2} \mathbf{x} \right\|_1. \quad (3.4)$$

(In statistics, this is called "LASSO"). The problem (3.4) can be reparametrized into an l^2 -regularized problem of the form

$$\min_{\mathbf{z}} \frac{1}{2} \left\| \mathbf{y} - \mathbf{F} \circ \mathbf{T}(\mathbf{z}) \right\|_2^2 + \frac{1}{2} \left\| \boldsymbol{\Sigma}^{-1/2} \mathbf{z} \right\|_2^2, \quad (3.5)$$

(in statistics, this is called "ridge regression") where \mathbf{T} is a suitable nonlinear transformation.

Consider the univariate transform $\tau : \mathbb{R} \rightarrow \mathbb{R}$ given by

$$\tau(x) = \frac{1}{2\sqrt{2}} \text{sign}(x) |x|^2. \quad (3.6)$$

This function is continuous and bijective (see [16]).

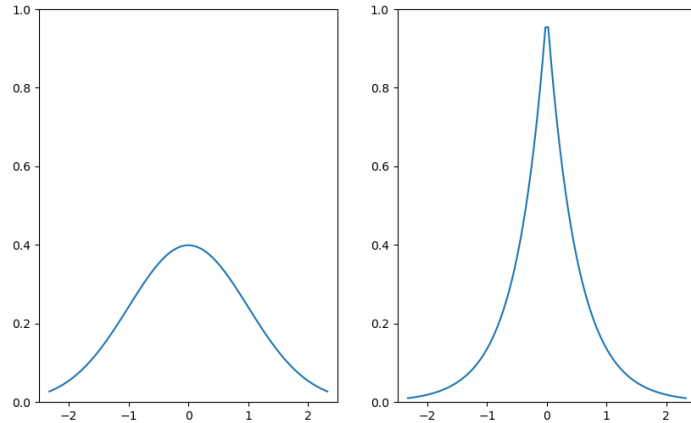


FIGURE 1. Probability density functions of the standard normal distribution (left) and the Laplace distribution, each with mean 0 and covariance 1 (right).

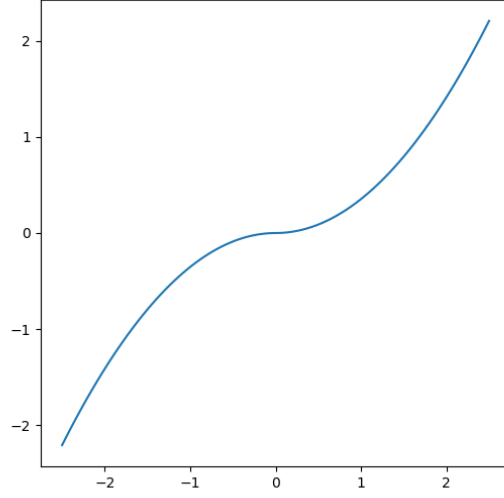


FIGURE 2. The transform τ defined in (3.6). It transports a standard normal distribution to a Laplace distribution with mean 0 and covariance 1.

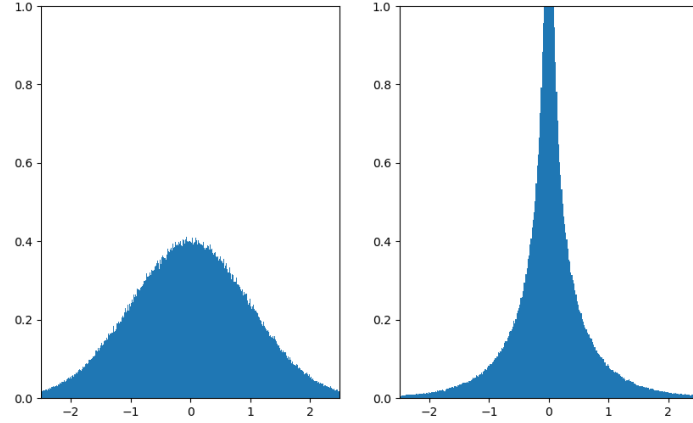


FIGURE 3. Left: Histogram of 10^6 independent samples of a standard normal distribution. Right: The image of those samples under the function τ .

The transform τ can easily be extended to \mathbb{R}^n . We define the multivariate transform $\mathbf{T} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ by

$$\mathbf{T}(\mathbf{x}) = \Sigma^{1/2} \tau(\Sigma^{-1/2} \mathbf{x}), \quad (3.7)$$

where $\tau(\mathbf{x}) = [\tau(x_1), \dots, \tau(x_n)]$ denotes componentwise application of τ . Proposition 3.5 in [16] then implies that with this choice for the transform T , problems (3.4) and (3.5) are equivalent. (Note that this is not true for arbitrary transformations. It is a well-known fact (and problem) of mode-based inference that maximum-a-posteriori estimators are not invariant under general reparametrization.)

If we apply this reparametrization trick to our particular problem (3.1)-(3.3), we obtain the constrained nonlinear least-squares problem

$$\min_{\mathbf{g}} \quad \frac{1}{2\delta^2} \|\mathbf{y} - \mathcal{G}(\mathbf{T}(\mathbf{g}), \boldsymbol{\theta}_v)\|_2^2 + \frac{1}{2} \|\tilde{\Sigma}_1^{-1/2} \mathbf{g}\|_2^2 + \frac{1}{2} \|\Sigma_2^{-1/2} (\boldsymbol{\theta}_v - \bar{\boldsymbol{\theta}}_v)\|_2^2, \quad (3.8)$$

$$\text{subject to} \quad \Phi^{-1} \mathbf{T}(\mathbf{g}) \geq \mathbf{0}, \quad (3.9)$$

$$\sum_i w_i T(\mathbf{g})_i = M_{\text{total}}. \quad (3.10)$$

where $\mathbf{g} = \mathbf{T}^{-1}(\tilde{\mathbf{f}})$. This can be written more concisely as

$$\min \quad \frac{1}{\sigma^2} \|\mathbf{y} - \mathcal{H}(\mathbf{z})\|_2^2 + \left\| \Sigma^{-1/2} \mathbf{z} \right\|, \quad (3.11)$$

$$\text{subject to} \quad \mathbf{T}(\mathbf{z}) \in \mathbb{F}, \quad (3.12)$$

where $\mathbf{z} = [\mathbf{z}_1, \mathbf{z}_2] = [\mathbf{g}, \boldsymbol{\theta}_v]$, $\mathcal{H}(\mathbf{z}) = \mathcal{G}(\mathbf{T}(\mathbf{z}_1), \mathbf{z}_2)$, \mathbb{F} is defined as above, and

$$\Sigma = \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{bmatrix}.$$

3.2. The constrained Gauss-Newton method. So, we have reduced our inverse problem to the solution of the finite-dimensional constrained nonlinear least-squares problem (3.8)-(3.10). Probably the most popular method for solving nonlinear least-squares problems is the Gauss-Newton method, and it can be generalized to problems with constraints. One such generalization that has seen a lot of applications in parameter estimation is the *constrained Gauss-Newton method* (CGN) (usually called the generalized Gauss-Newton method ("verallgemeinertes Gauß-Newton-Verfahren") by the people from Heidelberg, although there are a couple of different methods under that name). It is an extension of the classical Gauss-Newton method and was developed by Bock and colleagues for parameter estimation with ODEs (see [3] and [4]).

Consider a generic nonlinear least-squares problem with equality- and inequality constraints.

$$\begin{aligned} \min_{\mathbf{x}} \quad & \|\mathbf{F}_1(\mathbf{x})\|_2^2, \\ \text{subject to} \quad & \mathbf{F}_2(\mathbf{x}) = \mathbf{0}, \\ & \mathbf{F}_3(\mathbf{x}) \geq \mathbf{0}. \end{aligned}$$

(Note that our transformed problem (3.8)-(3.10) fits into this framework.) The CGN iteration is given by

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \Delta \mathbf{x}_k,$$

where \mathbf{x}_0 is an initial guess and $\Delta \mathbf{x}_k$ is a solution of the linearized subproblem

$$\begin{aligned} \min_{\Delta \mathbf{x}} \quad & \|\mathbf{F}_1(\mathbf{x}_k) + \mathbf{F}'_1(\mathbf{x}_k) \Delta \mathbf{x}\|_2^2, \\ \text{such that} \quad & \mathbf{F}_2(\mathbf{x}_k) + \mathbf{F}'_2(\mathbf{x}_k) \Delta \mathbf{x} = \mathbf{0}, \\ & \mathbf{F}_3(\mathbf{x}_k) + \mathbf{F}'_3(\mathbf{x}_k) \Delta \mathbf{x} \geq \mathbf{0}. \end{aligned}$$

This is a constrained linear least-squares problem, which can be solved using, for example, an active set approach. We omit a detailed description of the CGN method here, but it is available, for example, in the dissertation of Bock [4] (you can ask me if you want the pdf).

3.3. Uncertainty quantification. Suppose we have computed a MAP estimator \mathbf{z}^* , i.e. a solution of (3.11)-(3.12). Then, the transformed estimator

$$\begin{bmatrix} \tilde{\mathbf{f}}^* \\ \boldsymbol{\theta}_v^* \end{bmatrix} = \begin{bmatrix} \mathbf{T}(\mathbf{z}_1^*) \\ \mathbf{z}_2^* \end{bmatrix}$$

is a MAP estimate for our original problem.

We want to perform uncertainty quantification by computing Bayesian *highest probability-density regions* (HPD regions). The definition of the HPD region is distinct from the frequentist confidence interval, although both regions often overlap in practice. Given an unknown parameter $\mathbf{x} \in \mathbb{R}^n$ and a constant $\alpha \in (0, 1)$, an $(1 - \alpha)$ -HPD region C_α is any smallest subset of \mathbb{X} such that

$$\mathbb{P}(\mathbf{x} \in C_\alpha \mid \mathbf{y}) \geq 1 - \alpha.$$

That is, C_α is a minimal region that contains the parameter \mathbf{x} with probability $1 - \alpha$, given the measurement \mathbf{y} . Normally, such a region could be computed using a Markov chain Monte Carlo approach. However, since we do not want to use Monte Carlo methods (due to the fact that they do not scale well enough), we

have to use the next-best thing, which consists in mode-based inference. That is, we try to approximate C_α using local properties of the posterior distribution at the MAP estimator \mathbf{x}^* .

The limitation of this strategy comes from the fact that we are approximating global properties of the posterior distribution using only local information at the mode. Intuitively, this performs reasonably well if the greater part of the posterior probability mass is concentrated around a single point.

Mode-based covariance estimation for unconstrained maximum-a-posteriori inference is well-known and widely-used. If \mathbf{x}^* is a MAP estimator for the statistical inverse problem

$$\begin{aligned}\mathbf{y} &= \mathcal{G}(\mathbf{x}) + \boldsymbol{\xi}, \\ \mathbf{x} &\sim \mathcal{N}(\mathbf{x}_0, \boldsymbol{\Sigma}), \\ \boldsymbol{\xi} &\sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Gamma}),\end{aligned}$$

then the posterior covariance $\boldsymbol{\Sigma}_{\text{post}}$ can be approximated by

$$\begin{aligned}\hat{\boldsymbol{\Sigma}}_{\text{post}} &= \boldsymbol{\Sigma}^{1/2} \left(\boldsymbol{\Sigma}^{1/2} \mathcal{G}'(\mathbf{x}^*)^* \boldsymbol{\Gamma}^{-1} \mathcal{G}'(\mathbf{x}^*) \boldsymbol{\Sigma}^{1/2} + \mathbf{I} \right)^{-1} \boldsymbol{\Sigma}^{1/2} \\ &= (\boldsymbol{\Sigma}^{-1} + \mathcal{G}'(\mathbf{x}^*)^* \boldsymbol{\Gamma}^{-1} \mathcal{G}'(\mathbf{x}^*))^{-1}.\end{aligned}$$

One can then compute an approximation $\tilde{C}_\alpha^{(1)}$ of C_α using $\hat{\boldsymbol{\Sigma}}_{\text{post}}$ (see [10] and [1, section 2.3]). The approximate HPD region $\tilde{C}_\alpha^{(1)}$ is defined as

$$\tilde{C}_\alpha^{(1)} = \left\{ \mathbf{x} \in \mathbb{R}^n : \left\| \hat{\boldsymbol{\Sigma}}_{\text{post}}^{-1/2} (\mathbf{x} - \mathbf{x}^*) \right\| \leq \chi_n^2(1 - \alpha) \right\},$$

where $\chi_n^2(1 - \alpha)$ is the $(1 - \alpha)$ -quantile of the chi-squared distribution with n degrees of freedom. These estimates can be generalized to the constrained case as in [5].

Once we have computed $\hat{\boldsymbol{\Sigma}}_{\text{post}}$, we have to transform it back to a covariance for the original problem. This can be done by a Monte Carlo approach. If $\text{Cov}(\mathbf{g}) = \hat{\boldsymbol{\Sigma}}_{\text{post}}$ and $\mathbf{f} = \boldsymbol{\Phi}^{-1} \mathbf{T}(\mathbf{g})$, we can simply sample $\mathbf{g}_1, \dots, \mathbf{g}_K \sim \mathcal{N}(\mathbf{g}_*, \hat{\boldsymbol{\Sigma}}_{\text{post}})$, set $\mathbf{f}_k = \boldsymbol{\Phi}^{-1} \mathbf{T}(\mathbf{g}_k)$ for $k = 1, \dots, K$, and estimate

$$\text{Cov}(\mathbf{f}) \approx \frac{1}{K} \sum_{k=1}^K (\mathbf{f}_k - \bar{\mathbf{f}})(\mathbf{f}_k - \bar{\mathbf{f}})^\top,$$

where $\bar{\mathbf{f}} = \frac{1}{K} \sum_{k=1}^K \mathbf{f}_k$ is the sample mean.

An alternative estimator that works for linear problems is based on the *concentration of measure* phenomenon. This term describes a collection of mathematical and empirical results which state that the mass of high-dimensional probability distributions tends to be concentrated on lower-dimensional manifolds. Pereyra [13] used this phenomenon to derive an approximation $\tilde{C}_\alpha^{(2)}$ of the $(1 - \alpha)$ -HPD region. It is given by

$$\tilde{C}_\alpha^{(2)} = \left\{ \mathbf{x} \in \mathbb{R}^n : \Phi(\mathbf{x}) \leq \Phi(\mathbf{x}^*) + n(\tau_\alpha + 1) \right\},$$

where Φ is the MAP cost functional and $\tau_\alpha = \sqrt{\frac{16 \log(3/\alpha)}{n}}$. It is then shown in [13] that $\tilde{C}_\alpha^{(2)}$ is a conservative estimate of C_α , i.e. we have

$$C_\alpha \subset \tilde{C}_\alpha^{(2)}.$$

This is a very nice result, because it guarantees that we do not underestimate the uncertainty in our results. Alas, it only applies to log-concave probability distributions (e.g. Gaussians), which in our case is equivalent to the assumption that \mathcal{G} is linear. (But I mailed Pereyra and he said that if \mathcal{G} is not too nonlinear around \mathbf{x}_* , $\tilde{C}_\alpha^{(2)}$ should still be a useful approximation of C_α , especially if the regularization term is l^1 , as is the case in our situation! He suggested to use a convex approximation of the cost functional around \mathbf{x}_* , which we obtain for example by linearizing the observation operator.) We note also that in our situation, $C_\alpha^{(2)}$ can be computed for the non-transformed problem, since all we need is the MAP estimator.

To summarize, we have presented two mode-based approximations of the posterior HPD region. Both essentially assume that the model is linear, but we hope that they still give useful results in our mildly nonlinear case. The idea behind the first approximation $C_\alpha^{(1)}$ is classical (it can be found in research from 60

years ago [2]) and well-established as a tool in nonlinear regression [1]. The second approximation $C_\alpha^{(2)}$ is very modern (2017) and there are yet no papers that have applied it to nonlinear inverse problems, and only a couple studies in the linear case (for example, [6] and [14]). Since both $C_\alpha^{(1)}$ and $C_\alpha^{(2)}$ can be essentially computed for free once we have obtained the MAP estimator \mathbf{x}_* , we suggest to use both and see how they differ. It would also be wise to perform a preliminary experiment where MCMC-sampling is still doable. Then, the credible region computed by MCMC can serve as a gold standard with respect to which we can compare the performance of the approximations $C_\alpha^{(1)}$ and $C_\alpha^{(2)}$.

4. WHAT REMAINS TO BE DONE

If we agree on the computational method, we can finally start implementing.

- *Prior hyperparameters:* We have some hyperparameters that have to be tuned. Among them are the prior covariances Σ_1 and Σ_2 and the prior mean $\bar{\theta}_v$ for the hyperparameters of the Gauss-Hermite interpolation. We should start with something simple, like choosing Σ_1 and Σ_2 to be diagonal matrices, and then adapt if we have first numerical results.
- *Curvelet transform and optimization method:* There is already a good implementation of the fast curvelet transform in CurveLab (www.curvelet.org/software), which is freely available for non-commercial use. For the CGN method, I know some people that I could mail for code, but it will probably be easier in the long run if we implement it by ourselves.
- *Uncertainty quantification:* As discussed above, computation of confidence intervals are something that can easily be added once we have implemented the CGN method. Both of the presented approximations do not add much computational cost to the overall method, and so we can simply use both and compare them, either with each other, or with the results from the already existing Hamiltonian Monte Carlo implementation.

4.1. Open questions.

- (i) How does the discretization of the measurement P_{N_y} look like (see Section 1.2)? Is it accurate to write the discretized measurement $\bar{\mathbf{y}}$ as

$$\bar{y}_j = \int_{I_j} \frac{1}{1+v/c} S\left(\frac{\lambda}{1+v/c}; f\right) L(v; \boldsymbol{\theta}) dv,$$

where $I_j = [\lambda_j, \lambda_{j+1}) \subset \mathbb{R}^+$ is the j -th wavelength bin? How did you implement it?

- (ii) As an alternative to curvelets, we could also use a Gaussian mixture model, i.e. try to represent \mathbf{f} by a superposition of Gaussians, which can be generically written as

$$\sum_{k=1}^K \phi_k \mathcal{N}(\mathbf{m}_k, \boldsymbol{\Sigma}_k), \quad \phi_1, \dots, \phi_K \in \mathbb{R}.$$

This would simplify the later implementation, since for example the nonnegativity condition $f \geq 0$ is much easier to ensure for a Gaussian mixture (just choose $\phi_k \geq 0$ for all k). Furthermore, the parameters of a Gaussian mixture model are much easier to interpret than the parameters of a curvelet model. So, Gaussian mixtures could also be a viable alternative. But this hinges of course on the following question:

Is the distribution function of a galaxy well-modelled by a mixture of Gaussians?

4.2. Wishlist. This is a non-exhaustive list of things that we need, or that would be at least very helpful, before we can start implementing our approach.

- (i) *Observation operator:* An implementation of the input-output operator \mathcal{G} given by (1.4). Ideally an implemented function that takes \mathbf{f} and $\boldsymbol{\theta}_v$ as input and returns the corresponding measurement $\bar{\mathbf{y}}$. It would be helpful if this part is well-documented, so that we can adapt it quickly to the curvelet formulation.

- (ii) *Fake data simulator*: It would also be vital for testing to have a program that can simulate distribution functions of a galaxy. Ideally, the program would randomly generate realizations of \mathbf{f} that look (at least approximately) like distribution functions of real galaxies. Only then can we check how well the sparsity-promoting regularization is working.
- (iii) *Access to the Monte Carlo implementation*: It would also be very helpful as a reference to have access to the already existing Hamiltonian Monte Carlo implementation for the problem, including a description of the made choices (what are the prior distributions that you used, etc.). The program should be in a form where we can give it a noisy measurement \mathbf{y} as input (and maybe also tune the hyperparameters of the prior) and it automatically generates samples from the posterior distribution of \mathbf{f} and $\boldsymbol{\theta}_v$, given \mathbf{y} .

REFERENCES

- [1] D. M. Bates and D. G. Watts. “Nonlinear Regression Analysis and Its Applications”. John Wiley & Sons, Inc., 1988. DOI: [10.1002/9780470316757](https://doi.org/10.1002/9780470316757) (cited on pages 11, 12).
- [2] E. M. L. Beale. “Confidence Regions in Non-Linear Estimation”. In: *Journal of the Royal Statistical Society: Series B (Methodological)* 22.1 (1960), pp. 41–76. ISSN: 0035-9246. DOI: [10.1111/j.2517-6161.1960.tb00353.x](https://doi.org/10.1111/j.2517-6161.1960.tb00353.x) (cited on page 12).
- [3] H. G. Bock. “Numerical Treatment of Inverse Problems in Chemical Reaction Kinetics”. In: *Springer Series in Chemical Physics*. Springer Berlin Heidelberg, 1981, pp. 102–125. DOI: [10.1007/978-3-642-68220-9_8](https://doi.org/10.1007/978-3-642-68220-9_8) (cited on page 10).
- [4] H. G. Bock. “Randwertproblemmethoden zur Parameteridentifizierung in Systemen nichtlinearer Differentialgleichungen”. PhD thesis. Universität Bonn, 1985 (cited on page 10).
- [5] H. G. Bock, E. Kostina, and O. Kostyukova. “Covariance Matrices for Parameter Estimates of Constrained Parameter Estimation Problems”. In: *SIAM Journal on Matrix Analysis and Applications* 29.2 (2007), pp. 626–642. ISSN: 0895-4798. DOI: [10.1137/040617893](https://doi.org/10.1137/040617893) (cited on page 11).
- [6] X. Cai, M. Pereyra, and J. D. McEwen. “Uncertainty quantification for radio interferometric imaging: II. MAP estimation”. In: *Monthly Notices of the Royal Astronomical Society* 480.3 (2018), pp. 4170–4182. ISSN: 0035-8711. DOI: [10.1093/mnras/sty2015](https://doi.org/10.1093/mnras/sty2015) (cited on page 12).
- [7] E. Candès, L. Demanet, D. Donoho, and L. Ying. “Fast discrete curvelet transforms”. In: *Multiscale Modeling & Simulation* 5 (2006), pp. 861–899. ISSN: 1540-3459 (cited on page 6).
- [8] E. J. Candès and D. L. Donoho. “Continuous curvelet transform II: Discretization and frames”. In: *Applied and Computational Harmonic Analysis* 19.2 (2005), pp. 198–222. ISSN: 1063-5203. DOI: [10.1016/j.acha.2005.02.004](https://doi.org/10.1016/j.acha.2005.02.004) (cited on page 6).
- [9] P. Lambert, S. Pires, J. Ballot, R. A. García, J.-L. Starck, and S. Turck-Chièze. “Curvelet analysis of asteroseismic data”. In: *Astronomy & Astrophysics* 454.3 (2006), pp. 1021–1027. DOI: [10.1051/0004-6361:20054541](https://doi.org/10.1051/0004-6361:20054541) (cited on page 4).
- [10] D. Lu, M. Ye, and M. C. Hill. “Analysis of regression confidence intervals and Bayesian credible intervals for uncertainty quantification”. In: *Water Resources Research* 48.9 (2012). ISSN: 0043-1397. DOI: [10.1029/2011wr011289](https://doi.org/10.1029/2011wr011289) (cited on page 11).
- [11] J. Ma and G. Plonka. “The Curvelet Transform”. In: *IEEE Signal Processing Magazine* 27.2 (2010), pp. 118–133. DOI: [10.1109/msp.2009.935453](https://doi.org/10.1109/msp.2009.935453) (cited on page 4).
- [12] S. G. Mallat. “A wavelet tour of signal processing: the sparse way”. 3rd ed. Amsterdam; Boston: Elsevier/Academic Press, 2009 (cited on page 5).
- [13] M. Pereyra. “Maximum-a-Posteriori Estimation with Bayesian Confidence Regions”. In: *SIAM Journal on Imaging Sciences* 10.1 (2017), pp. 285–302. ISSN: 1936-4954. DOI: [10.1137/16m1071249](https://doi.org/10.1137/16m1071249) (cited on page 11).
- [14] A. Repetti, M. Pereyra, and Y. Wiaux. “Scalable Bayesian Uncertainty Quantification in Imaging Inverse Problems via Convex Optimization”. In: *SIAM Journal on Imaging Sciences* 12.1 (2019), pp. 87–118. ISSN: 1936-4954. DOI: [10.1137/18m1173629](https://doi.org/10.1137/18m1173629) (cited on page 12).

-
- [15] J. L. Starck, D. L. Donoho, and E. J. Candès. “Astronomical image representation by the curvelet transform”. In: *Astronomy & Astrophysics* 398.2 (2003), pp. 785–800. DOI: [10.1051/0004-6361:20021571](https://doi.org/10.1051/0004-6361:20021571) (cited on page [4](#)).
- [16] C. A. Zarzer. “On Tikhonov regularization with non-convex sparsity constraints”. In: *Inverse Problems* 25 (2009), p. 025006. ISSN: 0266-5611 (cited on pages [8](#), [9](#)).