GRAPH SPECTRAL CLUSTERING OF CONVOLUTION ARTEFACTS IN RADIO INTERFEROMETRIC IMAGES

Matthieu Simeoni[†]*, *Paul Hurley*[‡]

[†]IBM Zurich Research Laboratory, *École Polytechnique Fédérale de Lausanne (EPFL), [‡]Western Sydney University

ABSTRACT

The starting point for deconvolution methods in radio astronomy is an estimate of the sky intensity called a *dirty image.* These methods rely on the telescope *point*spread function so as to remove artefacts which pollute it. In this work, we show that the intensity field is only a partial summary statistic of the matched filtered interferometric data, which we prove is spatially correlated on the celestial sphere. This allows us to define a sky covariance function. This previously unexplored quantity brings us additional information that can be leveraged in the process of removing dirty image artefacts. We demonstrate this using a novel unsupervised learning method. The problem is formulated on a graph: each pixel interpreted as a node, linked by edges weighted according to their spatial correlation. We then use spectral clustering to separate the artefacts in groups, and identify physical sources within them.

Index Terms— Graph Spectral Clustering, Unsupervised Learning, Radio Interferometry, Dirty Image

1. INTRODUCTION

Radio interferometry is concerned with the sensing and analysis of electromagnetic fields produced by stars and celestial objects [1, 2, 3]. Since star radiation typically fluctuate randomly [4], astronomers make extensive use of summary statistics such as *moments* to analyse them. In particular, the second order moment –or *intensity* of radiation– is most often the quantity of interest. It is assessed by *matched filtering* the random radiation recorded on the ground by a network of antennas called an *interferometer* [5, 6]. The variance of the matched filtered output is then estimated for gridded directions in the sky, yielding the so-called *dirty image* [3]. This image can in general be shown to be the convolution between the intensity map of the underlying sources and the point-spread function of the matched filtering imaging procedure, called *dirty beam* [1]. This function, which depends only on the geometry of the interferometer in use, is in general poorly localised in space, typically



(a) True sky.

(b) Dirty beam. (c) Dirty image.

Fig. 1: The dirty image (c) is the result of the convolution between the true sky (a) and the dirty beam (b).

composed of a main *central lobe* surrounded by multiple *sidelobes* with smaller magnitudes (see fig. 1b). As a result, dirty images are most often polluted by strong *convolution artefacts*, which could be confused with actual stars (see fig. 1).

In this work, we take a radically new perspective on the problem, and propose a fully automatic unsupervised learning method permitting to cluster dirty image artefacts and extract their associated parent source, whose convolution with the dirty beam most likely generated said artefact. The clustering is based on a custom notion of *similarity* between features in the dirty image. To assess this similarity, we compute the covariance existing between matched filtered scans for any two directions in the sky. This covariance function is derived by reinterpreting the radio astronomy processing pipeline in a rigorous statistical framework. Since the proposed notion of similarity is highly non-local, we use *spectral* clustering [7] and define a graph structure on the dirty image: each pixel is interpreted as a node, and edges between pixels are weighted according to their similarity. Finally and once the dirty image features separated in clusters, we identify physical sources as the nodes with maximum connectivity within each cluster.

There is, to our knowledge, no precedent in the radio interferometry community to our approach. The deconvolution algorithms –such as CLEAN [8, 9] or related compressed sensing methods [10, 11]– traditionally used for removing dirty image artefacts, are unaware of the existence of a sky covariance function, and exploit only the intensity field. This necessarily handicaps the quality of results, and this work paves the road towards covariance-aware algorithms for post-processing dirty images in radio astronomy.

2. THE SKY COVARIANCE FUNCTION

In this section we reformulate radio interferometric imaging in a rigorous statistical framework. This allows us to derive a new summary statistic of the source field in the form of a *covariance kernel*, whose diagonal is the traditional dirty image [1, 3]. We shall use this covariance kernel in section 4.1 to assess similarities between features in the dirty image.

Interferometers [5, 6] are antenna networks for sensing random emissions from unknown celestial objects in the *far-field* [1]. Typically, the *source field* $S : \Omega \times \mathbb{S}^2 \to \mathbb{C}$ is modelled as a sum of Q point sources lying on the *celestial sphere* \mathbb{S}^2 and with randomly fluctuating amplitudes $\xi_q : \Omega \to \mathbb{C}$:

$$S(\boldsymbol{r}) = \sum_{q=1}^{Q} \xi_q \,\delta(\boldsymbol{r} - \boldsymbol{r}_q), \quad \forall \boldsymbol{r} \in \mathbb{S}^2, \tag{1}$$

where Ω is the sample space of some probability space, and $\{r_1, \ldots, r_Q\} \subset \mathbb{S}^2$ are the *unknown* sources locations. The amplitudes of the sources are generally assumed to be independent complex Gaussian random variables with mean zero and *unknown* variances σ_q [4]. When excited with a narrowband waveform of wavelength λ , the source field generates a far-field radiation pattern which is sampled by *L* antennas at locations $\{p_1, \ldots, p_L\} \subset \mathbb{R}^3$ on the ground. The Fraunhofer equation [12] provides us with an approximate expression of the radio telescope random measurements $\boldsymbol{y} : \Omega \to \mathbb{C}^L$:

$$\boldsymbol{y} := \sum_{q=1}^{Q} \xi_{q} \begin{bmatrix} e^{-\frac{2\pi j}{\lambda} \langle \boldsymbol{r}_{q}, \boldsymbol{p}_{1} \rangle} \\ \vdots \\ e^{-\frac{2\pi j}{\lambda} \langle \boldsymbol{r}_{q}, \boldsymbol{p}_{L} \rangle} \end{bmatrix} = \sum_{q=1}^{Q} \xi_{q} \boldsymbol{\varphi}(\boldsymbol{r}_{q}), \quad (2)$$

where $\varphi(r_q)$ is the steering vector with direction $r_q \in \mathbb{S}^2$.

Radio interferometric imaging [2, 3] can then be seen as characterising the stochastic behaviour of the source field *S* from independent observations $\{y_1, \ldots, y_N\} \subset \mathbb{C}^L$ of y. For uncorrelated Gaussian point sources as in eq. (1), this amounts to estimating the source positions r_q and intensities σ_q . Observe in eq. (2) that point sources have very characteristic signatures in the dataspace \mathbb{C}^L : they each contribute to the overall data in the form of a weighted steering vector. Evidence of the template vector $\varphi(r)$ in the data should hence reveal likely positions of the underlying sources. Such evidence can be gathered by *correlating* the measurements y with $\varphi(r)$:

$$\hat{S}(\boldsymbol{r}) = \langle \boldsymbol{y}, \boldsymbol{\varphi}(\boldsymbol{r}) \rangle = \varphi(\boldsymbol{r})^H \boldsymbol{y}, \quad \boldsymbol{r} \in \Theta$$

for some set $\Theta \subset \mathbb{S}^2$ of candidate locations¹. This signal processing technique is known as *matched filtering* [13] and is in radio astronomy key to the concept of the *dirty image* [3, 8], which can be seen as a statistical summary of the random object² $\hat{S} : \Omega \times \Theta \to \mathbb{C}$. Since \boldsymbol{y} has mean zero, it is indeed entirely characterised by its *covariance kernel*³ $\hat{\kappa} : \Theta \times \Theta \to \mathbb{C}$:

$$\hat{\kappa}(\boldsymbol{r},\boldsymbol{s}) := \mathbb{E}\left[\hat{S}(\boldsymbol{r})\hat{S}^{*}(\boldsymbol{s})\right] = \boldsymbol{\varphi}(\boldsymbol{r})^{H}\boldsymbol{\Sigma}\boldsymbol{\varphi}(\boldsymbol{s}), \qquad (3)$$

where $\Sigma := \mathbb{E}[\boldsymbol{y}\boldsymbol{y}^H] \in \mathbb{C}^{L \times L}$ is the *population* covariance matrix of the data, also called *visibility matrix* [1]. In practice this quantity is of course unavailable to us, and must be replaced it by its *empirical* counterpart: $\hat{\Sigma} = (1/N) \sum_{n=1}^{N} \boldsymbol{y}_n \boldsymbol{y}_n^H$. Radio astronomy has not exploited the off-diagonal part of $\hat{\kappa}$, to date only considering the diagonal, or *dirty image*,

$$\hat{I}(\boldsymbol{r}) := \hat{\kappa}(\boldsymbol{r}, \boldsymbol{r}) = \boldsymbol{\varphi}(\boldsymbol{r})^H \Sigma \boldsymbol{\varphi}(\boldsymbol{r}), \quad \boldsymbol{r} \in \Theta.$$
 (4)

This is because of a simplification in modelling, that overlooks the existence of the sky covariance function. In this work we argue that this function actually contain crucial information content for ridding the dirty image of convolution artefacts mentioned in section 1 and investigated in greater detail in the next section.

3. THE PARENTING PROBLEM

In this section we introduce the *parenting problem*, a classification problem aiming at associating features in the dirty image to their physical *parent source* in the underlying source field. As motivation, we investigate the structure of the dirty image, and make apparent the convolutional nature of the artefacts polluting it. Using eqs. (1) and (2) as well as the definition of Σ and the steering vector, we can indeed rewrite the dirty image in eq. (4) as

$$\hat{I}(\boldsymbol{r}) = \boldsymbol{\varphi}(\boldsymbol{r})^{H} \left[\sum_{q=1}^{Q} \sigma_{q} \boldsymbol{\varphi}(\boldsymbol{r}_{q}) \boldsymbol{\varphi}(\boldsymbol{r}_{q})^{H} \right] \boldsymbol{\varphi}(\boldsymbol{r})$$

$$= \sum_{q=1}^{Q} \sigma_{q} \left| \zeta(\boldsymbol{r} - \boldsymbol{r}_{q}) \right|^{2} \qquad (5)$$

$$:= I(\boldsymbol{r})$$

$$= \int_{\mathbb{S}^2} |\zeta(\boldsymbol{r} - \boldsymbol{s})|^2 \left[\sum_{q=1}^Q \sigma_q \delta(\boldsymbol{s} - \boldsymbol{r}_q) \right] d\boldsymbol{s}, \quad (6)$$

¹Typically a uniform tessellation of the sphere.

²As per the usual nomenclature in statistics, the random quantity \hat{S} is called random vector, random process or random field for Θ respectively finite, countably infinite or uncountably infinite.

³Again $\hat{\kappa}$ will either be a matrix, an infinite dimensional matrix or a kernel depending on the cardinality of Θ .



Fig. 2: Synthetic example of dirty image (plain line) for a dirty beam given by $\zeta(x) = \operatorname{sinc}(x)$ and Q = 3 sources, with intensities $\sigma_q = 2, 6, 7$ and locations $x_q = 0, 7, 4$. The contributions (dashed lines) and regions of dominance (ROD) of the three sources (cyan, blue and magenta) are displayed.

where $\zeta(\mathbf{r}) := \sum_{i=1}^{L} e^{\frac{2\pi j}{\lambda} \langle \mathbf{r}, \mathbf{p}_i \rangle}$, $\mathbf{r} \in \mathbb{S}^2$. Equations (5) and (6) offer two dual perspectives on the dirty image. Equation (6) writes the dirty image as the *convolution* between the *intensity field I* of the underlying source field *S* and the so-called *dirty beam* $|\zeta|^2$ [3]. As previously noted, this function is in general poorly localised in space, hence incurring severe convolution artefacts within the dirty image. Such artefacts can be better understood by looking at eq. (5), which decomposes the dirty image in a sum of *Q* independent contributions from each source in the field. In practice, each source will typically not contribute equally to the overall intensity. We hence associate to a source *q* a *region of dominance* $\mathcal{R}_q \subset \Theta$, defined as the level sets of the *parenting function*:

Definition 1 (Parenting Function & Region of Dominance). We call the parenting function the function $\pi: \Theta \rightarrow \{1, \dots, Q\}$ defined as

$$\pi(\boldsymbol{r}) := \operatorname*{argmax}_{q=1,\dots,Q} \left\{ \sigma_q \, |\zeta^2(\boldsymbol{r} - \boldsymbol{r}_q)|^2 \right\}, \quad \boldsymbol{r} \in \Theta.$$
 (7)

The parenting function associates each point $\mathbf{r} \in \Theta$ to its parent source, which contributed the most to the value of the dirty image at this location. The region of dominance (ROD) $\mathcal{R}_q \subset \Theta$ for a given source $q \in \{1, \dots, Q\}$ is

$$\mathcal{R}_q := \{ \boldsymbol{r} \in \Theta : q = \pi(\boldsymbol{r}) \}.$$

In fig. 2 we computed the regions of dominance for a simulated example with three parent sources. Observe that in this example, the regions of dominance form a partition of Θ and classified each source in a different cluster, which is general behaviour ⁴.

Our goal will hence be to estimate the regions of dominance within a given dirty image and obtain estimates of their associated parent sources according to a process described in section 4.2. Solving for this classification problem is equivalent to estimating the parenting function eq. (7). For this reason, we will refer to it as the *parenting problem*. This problem is non-trivial since the parenting function involves unknown quantities, namely the source intensities and locations. Moreover, since the number of sources *Q* is in general unknown, we must also learn the number of groups, transforming a classification problem into a *clustering* one.

4. SOLVING THE PARENTING PROBLEM

4.1. Computing the Regions of Dominance by Spectral Clustering

In this section, we assume Θ to be a finite and discrete set of size N, and estimate the regions of dominance by reinterpreting the problem as a *graph clustering problem* [7]. To this end, we initially define an *undirected fully connected graph* $\mathcal{G} = (V, E)$, with *node set* $V = \Theta$ and *edge set* $E = V \times V$. To encode in the network potential similarities existing between nodes with common parent sources, we attribute weights $w_{ij} \in \mathbb{R}_+$ to each edge $e_{ij} = (\mathbf{r}_i, \mathbf{r}_j) \in E$,

$$w_{ij} = \sigma(\mathbf{r}_i, \mathbf{r}_j) \ge 0, \quad i, j = 1, \dots, N,$$

where σ is a suitable *similarity measure* $\sigma : \Theta^2 \to \mathbb{R}_+$, assessing the *degree of kinship* between any two nodes in the graph. Edges with weights close to zero are not related and hence discarded from the edge set. The goal is then to cluster this kinship network in order to reconstruct the regions of dominance and identify the parent sources within each of them. Not surprisingly, the success of this operation will heavily depend on the chosen similarity measure [7]. Given the close link between the dirty image and the covariance function in eq. (3), it seemed natural to us to define the similarity of any two points in Θ as the modulus of their covariance

$$\sigma(\boldsymbol{s}, \boldsymbol{r}) = |\hat{\kappa}(\boldsymbol{s}, \boldsymbol{r})|, \quad \forall (\boldsymbol{s}, \boldsymbol{r}) \in \Theta^2.$$

The clustering step is then performed by means of spectral clustering (see [7] for more details on the algorithm). Roughly speaking, this algorithm aims to partition the graph into K connected components of comparable size, such that the sum of the weights of the intercomponents edges is minimised. Various heuristics can be used to estimate the number of clusters K, but the most popular one is certainly the *eigengap heuristic* [7], which we used here. In fig. 3, we constructed the kinship graph for the scenario described in fig. 2, and performed spectral clustering to recover the regions of dominance. The eigengap heuristic for this example yields K = Q = 3 clusters, or exactly as much as the number of point sources in the field (as expected). Accuracy wise, more than 90% of the points in Θ are correctly classified in their actual region of dominance.

⁴The regions of dominance could actually intersect in very degenerate cases, but this is virtually impossible to happen in practice.



(a) The kinship graph and the clusters obtained with spectral clustering (cyan, magenta and blue). The sum of the weights of the edges (in purple here) than need to be removed to separate the graph in three disconnected components is minimised.



(b) Estimates of the regions of dominance for each source, using spectral clustering. The misclassified portions are highlighted in orange ($\simeq 10\%$ of the total region). The centers of the clusters (coloured pentagrams) as well as the estimated intensities of the parent sources (coloured squares) are also displayed.

Fig. 3: Estimating the regions of dominance with the kinship graph formulation and spectral clustering.

4.2. Estimating the Parent Sources

Disposing of an algorithm to recover the respective regions of dominance, we are now interested in finding their associated parent source. This task can be accomplished by analysing the connection network within each cluster. Indeed, it makes intuitive sense that the parent source should be very connected with its children nodes, or mathematically speaking have the highest kinship degree within its corresponding region of dominance. We then propose to recover the parent sources $\hat{\boldsymbol{r}}_k \in V$ of each cluster \mathcal{R}_k by solving the following Koptimisation problems

$$\hat{\boldsymbol{r}}_k = \operatorname*{argmax}_{\boldsymbol{r}_j \in \mathcal{R}_k} \sum_{\boldsymbol{r}_i \in \mathcal{R}_k} w_{ij}, \quad k = 1, \dots, K.$$
 (8)

The vertex set V being in bijection with the set $\Theta \subset \mathbb{S}^2$, we have then equivalently recovered the positions of each of the sources in the field. Intensities of the sources can then be trivially recovered by solving a linear problem, as described in [14, Section 2.3.2].

5. EXPERIMENTAL RESULTS

Figure 4 shows a slightly more realistic scenario with five point sources of various intensities. We used eq. (2) to simulate N = 800 realisations of the random measurement vector y as sensed by the first 24 core stations of the LOFAR interferometer [5]. For more realistic experimental conditions, the data was furthermore corrupted by white additive complex Gaussian noise,







(b) Dirty image.

(c) Kinship graph.

(a) Actual sky with 5 point sources.







(d) Regions dominance magenta, orange and locations source source (white squares).

green, magenta, locations (white dots).

of (e) Actual regions (f) Correctly (green) (cyan, of dominance (cyan, and wrongly (red) green, classified portions of yellow) orange and yellow) the field (accuracy and estimated parent and actual parent rate of $\simeq 86.7\%$).

Fig. 4: Solving for the parenting problem with spectral clustering in radio-astronomy. Eight hundreds simulated samples from 24 stations of the LOFAR telescope were used to estimate the dirty image and covariance function. The peak signal to noise ratio for this experiment is of -23 dB.

with PSNR around -23 dB. Despite the very high noise level, the results remain very satisfactory: spectral clustering empowered by our covariance-based similarity measure could recover the regions of dominance with an accuracy of 86.7%. The source locations and intensities were also recovered almost perfectly.

6. CONCLUSIONS

We derived, for what we believe to be the first time, the sky covariance function from the matched filtered output of a radio interferometer. This function measures the correlation between matched filtered scans for various focus directions in the sky. Existing deconvolution algorithms have not made use of this information. We, on the other hand, proposed a novel method to locate sources within dirty images which leverages this additional information. Our technique clusters strongly correlated convolution artefacts together and identify the parent source from which they originate within each cluster. To perform the clustering step, we constructed a kinship network encoding correlations between various features in the dirty image, on which we performed spectral clustering to learn the regions of dominance of the respective sources, and identify the actual sources in the field. Our initial tests are extremely promising, and further experiments are planned as future work.

7. REFERENCES

- [1] A Richard Thompson, James M Moran, and George W Swenson Jr, *Interferometry and synthesis in radio astronomy*, John Wiley & Sons, 2008.
- [2] Matthieu Simeoni, "Towards more accurate and efficient beamformed radio interferometry imaging," M.S. thesis, EPFL, Spring 2015.
- [3] Alle-Jan van der Veen and Stefan J Wijnholds, "Signal processing tools for radio astronomy," in *Handbook of Signal Processing Systems*, pp. 421–463. Springer, 2013.
- [4] Hamid Krim and Mats Viberg, "Two decades of array signal processing research: the parametric approach," *IEEE Signal processing magazine*, vol. 13, no. 4, pp. 67–94, 1996.
- [5] MP Van Haarlem, MW Wise, AW Gunst, George Heald, JP McKean, JWT Hessels, AG De Bruyn, Ronald Nijboer, John Swinbank, Richard Fallows, et al., "Lofar: The low-frequency array," Astronomy & Astrophysics, vol. 556, pp. A2, 2013.
- [6] Peter E Dewdney, Peter J Hall, Richard T Schilizzi, and T Joseph LW Lazio, "The square kilometre array," *Proceedings of the IEEE*, vol. 97, no. 8, pp. 1482– 1496, 2009.
- [7] Ulrike Von Luxburg, "A tutorial on spectral clustering," *Statistics and computing*, vol. 17, no. 4, pp. 395–416, 2007.
- [8] JA Högbom, "Aperture synthesis with a nonregular distribution of interferometer baselines," *Astronomy and Astrophysics Supplement Series*, vol. 15, pp. 417, 1974.
- [9] Cyril Tasse, S van der Tol, J van Zwieten, Ger van Diepen, and S Bhatnagar, "Applying full polarization a-projection to very wide field of view instruments: An imager for lofar," Astronomy & Astrophysics, vol. 553, pp. A105, 2013.
- [10] Jérôme Bobin, Jean-Luc Starck, and Roland Ottensamer, "Compressed sensing in astronomy," *IEEE Journal of Selected Topics in Signal Processing*, vol. 2, no. 5, pp. 718–726, 2008.
- [11] Yves Wiaux, Laurent Jacques, Gilles Puy, Anna MM Scaife, and Pierre Vandergheynst, "Compressed sensing imaging techniques for radio interferometry," *Monthly Notices of the Royal Astronomical Society*, vol. 395, no. 3, pp. 1733–1742, 2009.

- [12] T Douglas Mast, "Fresnel approximations for acoustic fields of rectangularly symmetric sources," *The Journal of the Acoustical Society of America*, vol. 121, no. 6, pp. 3311–3322, 2007.
- [13] Mehrdad Soumekh, Synthetic aperture radar signal processing, vol. 7, New York: Wiley, 1999.
- [14] Hanjie Pan, Matthieu Simeoni, Paul Hurley, Thierry Blu, and Martin Vetterli, "Leap: Looking beyond pixels with continuous-space estimation of point sources," *Astronomy & Astrophysics*, vol. 608, pp. A136, 2017.