

De la possibilité d'ignorer un bruit de calibration lors de la construction d'un estimateur du maximum de vraisemblance en radio-astronomie

Léontine SÉGAL^{1,2} Antoine ROUEFF² Claude JAUFFRET² Jérôme PETY^{1,3} Maryvonne GERIN³ et ORION-B⁴

¹Institut de Radioastronomie Millimétrique, 300 rue de la Piscine, 38400 Saint-Martin-d'Hères, France

²Université de Toulon, Aix Marseille Univ, CNRS, IM2NP, Toulon, France

³LUX, France

⁴<https://www.iram.fr/~pety/ORION-B/team.html>

Résumé – Nous présentons une méthodologie permettant de quantifier la perte de précision sur l'estimation de paramètres d'intérêt lorsqu'un bruit de calibration est négligé. L'approche repose sur l'utilisation des bornes de Cramér-Rao (BCR) et BCR modifiée (BCRM), cette dernière étant une borne inférieure sur la matrice de covariance d'un estimateur construit avec un modèle mal spécifié. La méthode est illustrée sur un cas théorique simple avant d'être appliquée au cas de données astrophysiques, où on estime la température cinétique et la densité volumique du gaz interstellaire. La perte de précision mesurée, inférieure à 13%, est négligeable compte-tenu de la complexité du problème astrophysique considéré.

Abstract – We propose a method for quantifying the loss of precision in the estimation of parameters of interest when a calibration noise is neglected. The method is based on the use of the Cramér-Rao Bound (CRB) and the modified CRB, this latter being a lower bound on the covariance matrix of an estimator built with a misspecified model. The approach is illustrated on a simple theoretical case before being applied to the case of astrophysical data, where we estimate the kinetic temperature and the volume density of the interstellar medium. The accuracy loss, less than 13%, is negligible given the complexity of the problem under consideration.

1 Introduction

On s'intéresse aux émissions du gaz composant les nuages interstellaires, régions où se forment les étoiles. Ces nuages moléculaires froids (≤ 100 K) sont observés dans le domaine des ondes millimétriques. Chaque espèce chimique émet un spectre d'émission, ou « raie moléculaire », observé sur N canaux fréquentiels. Comme proposé dans [2], les observations se modélisent par

$$\mathbf{x} = c \mathbf{s}(\boldsymbol{\theta}) + \mathbf{b} \quad (1)$$

où $\mathbf{x} = (x_1 \cdots x_N)^\top$ est entâché de deux bruits statistiquement indépendants : un bruit additif $\mathbf{b} = (b_1 \cdots b_N)^\top$ tel que $b_n \in [1, N] \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma_b^2)$ et un bruit de calibration multiplicatif $c \in \mathbb{R}$. Le signal d'intérêt $\mathbf{s} = (s_1 \cdots s_N)^\top$ est fonction du paramètre $\boldsymbol{\theta} \in \mathbb{R}^d$ caractérisant le milieu. Le bruit additif est dominé par l'interaction entre l'atmosphère et les photons avant qu'ils ne parviennent au télescope terrestre. Le bruit de calibration résulte de la stratégie d'observation et des incertitudes liées aux éléments successifs de la chaîne de mesure. Dans d'autres travaux, ce terme multiplicatif est aussi fonction de $\boldsymbol{\theta}$ et contribue donc à caractériser la source [8].

Pour le jeu de données acquis dans le cadre de la campagne d'observation du Large Program ORION-B (P.I. : J. Pety & M. Gerin), on mesure une dispersion du bruit de calibration $\sigma_c \sim 0.05$ à 0.1 selon la fenêtre atmosphérique considérée (de 3 mm à 1 mm) [2]. Pour analyser ces données ORION-B, nous avons construit dans [12] un estimateur du maximum de vraisemblance (MV) de $\boldsymbol{\theta}$ en négligeant le bruit de calibration, car de faible dispersion. Dans cet article, nous présentons une méthodologie afin de quantifier la perte de précision sur l'estimation de $\boldsymbol{\theta}$ induite par ce choix.

Ayant négligé le bruit de calibration, nous considérons la borne de Cramér-Rao Modifiée (BCRM) [3]. En effet, comme illustré dans les travaux en géolocalisation [7] et applications RADAR [9], la BCRM permet de prendre en compte les erreurs de modèle dans la quantification des performances d'estimation. La BCR hybride [6], adaptée à l'estimation jointe de paramètres déterministes et aléatoires, aurait pu être considérée. Toutefois, le paramètre de nuisance aléatoire c peut être marginalisé pour de faibles σ_c , afin de n'estimer que le paramètre d'intérêt déterministe $\boldsymbol{\theta}$. La borne de Cramér-Rao (BCR) sur ce modèle marginalisé nous fournit une précision de référence. Notre analyse de précision repose sur la comparaison de la BCRM à cette BCR.

Cet article est construit comme suit. La Sec. 2 présente une étude de précision théorique menée dans le cadre d'un modèle simplifié de Eq. (1) permettant d'effectuer des calculs analytiques. La Sec. 3 présente l'application de la méthodologie aux observations astrophysiques du projet ORION-B, dans le cadre de l'estimation de la température cinétique et de la densité volumique du gaz interstellaire. Les résultats obtenus et perspectives sont discutés en Sec. 4.

2 Étude théorique

Dans cette étude, on montre que la distribution du bruit de calibration peut s'approximer par une loi normale pourvu que l'écart-type soit inférieur à ~ 0.3 . Ensuite, nous calculons analytiquement la perte de précision relative induite lorsque ce bruit est négligé pour un modèle simplifié de Eq. (1).

2.1 Approximation Gaussienne de c

Comme le montre [4], la loi log-normale, notée $\text{log-}\mathcal{N}$, est appropriée pour décrire un bruit multiplicatif tel que c . La Fig. 1a présente trois densités $\text{log-}\mathcal{N}$ de paramètres $\{\mu, \sigma^2\}$ variables, indiqués en légende. Une loi normale $\mathcal{N}(\mu_c, \sigma_c^2)$ avec $\sigma_c = 0.1$ et $\mu_c = 1$ se confond à la $\text{log-}\mathcal{N}$ de plus faible dispersion. L'approximation Gaussienne devient critiquable pour $\sigma_c \gtrsim 0.3$, compte-tenu de l'asymétrie apparente de la $\text{log-}\mathcal{N}$ (courbe orange continue). La caractérisation du bruit de calibration faite par [2] indique que $0.05 \leq \sigma_c \leq 0.1$ pour les données ORION-B. Dans cette étude, nous considérerons donc que $c \sim \mathcal{N}(\mu_c = 1, \sigma_c^2)$ et gardons à l'esprit que le modèle de mesures construit n'est valable que pour $\sigma_c \leq 0.3$.

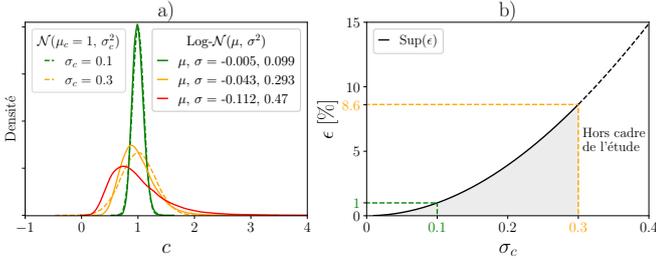


FIGURE 1 : a) $\text{Log-}\mathcal{N}$ (courbes continues) et \mathcal{N} (discontinues). b) Écart relatif ϵ sur la précision d'estimation (%) (défini Sec. 2.2.3) en fonction de σ_c pour différentes valeurs de $\{N, \sigma_b, \theta\}$ et $\mu_c = 1$ (région grisée). La borne supérieure $\text{Sup}(\epsilon)$ est en noir (continue).

2.2 Illustration sur un modèle simplifié

2.2.1 Description du modèle

A titre didactique, on considère le modèle suivant

$$\mathbf{x} = c\theta\mathbf{1}_N + \mathbf{b} \quad (2)$$

avec $\mathbf{1}_N = (1 \dots 1)^\top$ de taille N , $\theta \in \mathbb{R}$, $c \sim \mathcal{N}(\mu_c, \sigma_c^2)$ et $\mathbf{b} \sim \mathcal{N}(\mathbf{0}_N, \sigma_b^2 \mathbf{I}_N)$ où $\mathbf{0}_N = (0 \dots 0)^\top$ de taille N et \mathbf{I}_N la matrice identité de taille N . Les bruits \mathbf{b} et c étant Gaussiens et indépendants, le couple d'éléments aléatoires (\mathbf{X}, C) suit aussi une distribution Gaussienne. La marginalisation $p(\mathbf{x}|\theta) = \int p_{\mathbf{X}, C}(\mathbf{x}, c|\theta) dc$, souvent difficile à calculer analytiquement, s'exprime simplement ici car \mathbf{X} est la somme de deux vecteurs aléatoires Gaussiens. Dans notre cas, la distribution marginale est

$$p(\mathbf{x}|\theta) = \frac{1}{\sqrt{\det(2\pi\Sigma)}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\} \quad (3)$$

de moyenne $\boldsymbol{\mu} = \mu_c\theta\mathbf{1}_N$ et de matrice de covariance $\Sigma = \sigma_c^2\theta^2\mathbf{1}_N\mathbf{1}_N^\top + \sigma_b^2\mathbf{I}_N$. De plus, on a

$$\det(\Sigma) = \sigma_b^{2(N-1)}\gamma^2, \quad \text{avec} \quad \gamma^2 = \sigma_c^2 N\theta^2 + \sigma_b^2 \quad (4)$$

et la formule de *Sherman-Morrison* [1] permet de calculer

$$\Sigma^{-1} = \frac{1}{\sigma_b^2} \left(\mathbf{I}_N - \frac{\sigma_c^2\theta^2}{\gamma^2} \mathbf{1}_N\mathbf{1}_N^\top \right) \quad (5)$$

et ainsi d'avoir une formulation explicite de la vraisemblance.

2.2.2 Calcul des précisions sur l'estimation de θ

La BCR est une borne inférieure de la matrice de covariance d'un estimateur sans biais [5]. Notée \mathcal{B} , elle est l'inverse de la matrice d'information de Fisher, soit $\mathcal{B}(\theta) = \mathcal{I}(\theta)^{-1}$, où

$$\mathcal{I}(\theta) = \mathbb{E} \left[\frac{\partial \ln p(\mathbf{x}|\theta)}{\partial \theta} \frac{\partial \ln p(\mathbf{x}|\theta)}{\partial \theta}^\top \right]. \quad (6)$$

Pour le modèle Eq. (2), $\mathcal{I}(\theta) = \frac{N}{\gamma^4} (\mu_c^2 \gamma^2 + 2\sigma_c^4 N\theta^2)$ et

$$\mathcal{B}(\theta) = \frac{\gamma^2}{N} \frac{1}{\mu_c^2 + \frac{2\sigma_c^4 N\theta^2}{\gamma^2}}. \quad (7)$$

Ignorer la présence du bruit de calibration, comme cela a été fait dans [11, 12], équivaut à attribuer aux mesures la distribution approchée $p_f(\mathbf{x}|\theta)$,

$$p_f(\mathbf{x}|\theta) = \frac{1}{\sqrt{(2\pi\sigma_b^2)^N}} \exp\left\{-\frac{1}{2\sigma_b^2} \sum_{n=1}^N (x_n - \theta)^2\right\}. \quad (8)$$

Pour $\mu_c = 1$, l'approximation de Eq. (3) par (8) n'induit pas de biais d'estimation. La BCRM construite à partir de $p_f(\mathbf{x}|\theta)$, notée $\mathcal{B}_M(\theta)$, est définie alors par $\mathcal{B}_M(\theta) = A^{-1}JA^{-1}$ [10, 9], où

$$A = -\mathbb{E}_{p(\mathbf{x}|\theta)} \left[\frac{\partial^2 \ln p_f(\mathbf{x}|\theta)}{\partial^2 \theta} \right] \quad (9)$$

et

$$J = \mathbb{E}_{p(\mathbf{x}|\theta)} \left[\frac{\partial \ln p_f(\mathbf{x}|\theta)}{\partial \theta} \frac{\partial \ln p_f(\mathbf{x}|\theta)}{\partial \theta}^\top \right]. \quad (10)$$

Ici, $A = -\frac{N}{\sigma_b^2}$ et $J = \frac{N}{\sigma_b^2} (\sigma_c^2\theta^2 N + \sigma_b^2)$, d'où

$$\mathcal{B}_M(\theta) = \frac{\gamma^2}{N}. \quad (11)$$

Comme montré dans [6], $\mathcal{B}_M(\theta) \geq \mathcal{B}(\theta)$:

$$\mathcal{B}(\theta) = \mathcal{B}_M(\theta) \frac{1}{\mu_c^2 + \frac{2\sigma_c^4 N\theta^2}{\gamma^2}}, \quad \mu_c = 1 \quad \text{et} \quad \frac{2\sigma_c^4 N\theta^2}{\gamma^2} \geq 0.$$

2.2.3 Perte de précision sur l'estimation de θ

On quantifie la perte de précision sur l'estimation de θ avec l'écart relatif

$$\epsilon(\tau) \triangleq \frac{\sqrt{\mathcal{B}_M} - \sqrt{\mathcal{B}}}{\sqrt{\mathcal{B}}} \quad (12)$$

Ici, $\epsilon(\tau) = \sqrt{\mu_c^2 + \frac{2\sigma_c^2}{1+\tau}} - 1$, où $\tau = \frac{\sigma_b^2}{\sigma_c^2 N\theta^2}$. La fonction $\epsilon(\tau)$, strictement décroissante sur \mathbb{R}^* , est majorée par

$$\text{Sup}(\epsilon) = \lim_{\tau \rightarrow 0} \epsilon(\tau) = \sqrt{\mu_c^2 + 2\sigma_c^2} - 1. \quad (13)$$

Pour ce modèle particulier Eq. (2), la perte de précision induite par la négligence du bruit de calibration est essentiellement régie par la dispersion σ_c . On retrouve en particulier, que lorsque $2\sigma_c^2 \ll \mu_c^2$, l'écart relatif tend à s'annuler, traduisant l'équivalence des deux modèles. La Fig. 1b montre en fonction de σ_c l'écart relatif ϵ , pour $\mu_c^2 = 1$. On constate que l'écart maximal, $\text{Sup}(\epsilon)$, demeure inférieur à 10% pour $\sigma_c \leq 0.3$. L'écart relatif augmente à mesure que σ_c augmente. On observe ainsi un lien entre la perte de précision et la qualité de l'approximation Gaussienne d'une loi $\text{log-}\mathcal{N}$.

3 Étude sur les données astrophysiques

Après avoir illustré notre méthodologie sur le modèle théorique, nous examinons une partie des données ORION-B, correspondant aux observations de la nébuleuse de la Tête de Cheval [12].

3.1 Analyse du milieu interstellaire

Notre compréhension de la formation stellaire s'appuie essentiellement sur l'ajustement d'un modèle physique aux observations, ce qui permet d'estimer notamment la température cinétique (T en K) et la densité volumique (n en cm^{-3}) du gaz. L'observation de la Tête de Cheval sur une large bande fréquentielle a permis de mesurer une vingtaine de raies moléculaires par pixel. Comme le montre la Fig. 2a, observer simultanément différentes raies moléculaires permet de sonder l'évolution des propriétés du gaz le long de la ligne de visée.

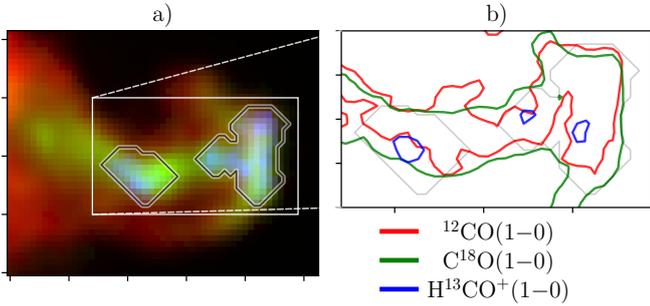


FIGURE 2 : a) Émissions intégrées normalisées ($\sum_n x_n / \max x$) des raies (1-0) de ^{12}CO (canal rouge), C^{18}O (vert) et H^{13}CO^+ (bleu) pour les lignes de visées où $\text{RSB} \geq 3$. Les régions les plus denses sont délimitées en noir. b) Zones où le $\text{RSB} \geq 250, 50, 10$ pour ^{12}CO (rouge), C^{18}O (vert), et H^{13}CO^+ (bleu), respectivement.

L'émission de la raie $^{12}\text{CO}(1-0)$ (en rouge) permet de sonder l'enveloppe chaude et diffuse ($\geq 30 \text{ K}$, $\leq 10^3 \text{ cm}^{-3}$). L'émission de $\text{C}^{18}\text{O}(1-0)$ (en vert) provient d'une région plus compacte d'aspect filamentaire, de température et densité intermédiaires. La raie $\text{H}^{13}\text{CO}^+(1-0)$ (en bleu) n'est détectée que dans deux régions de la nébuleuse, appelées « cœurs denses » car la densité y est la plus élevée ($\geq 10^5 \text{ cm}^{-3}$). Ces régions sont aussi les plus froides ($\sim 5 - 15 \text{ K}$) car protégées des radiations extérieures. C'est au sein des cœurs denses que se formeront des étoiles.

Le code de transfert radiatif RADEX [13] permet de simuler une raie moléculaire à partir des composantes du paramètre θ (dont T et n) d'un milieu homogène. Afin de tenir compte de l'hétérogénéité dans la ligne de visée, nous avons proposé dans [12] un modèle de nuage à 3 couches, disposées en « sandwich » (voir Fig. 3). Cette approche permet de décrire les régions denses de la nébuleuse par un cœur dense et froid, enveloppé dans du gaz plus diffus et chaud ([12], Fig. 16). Comme l'illustre la Fig. 2b, le rapport signal sur bruit (RSB), défini par $\frac{\max x}{\sigma_b}$, diffère de plus d'un ordre de grandeur pour différentes raies moléculaires au sein d'un même pixel. Cette disparité est problématique car la simplicité du modèle sandwich pour décrire le nuage induit inévitablement des erreurs de modèle. Nous modérons donc la confiance attribuée aux

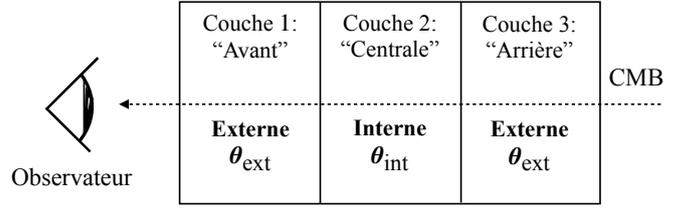


FIGURE 3 : Modèle « sandwich » de la ligne de visée. L'émission mesurée résulte des contributions de l'avant et arrière de la couche externe, de la couche interne et du fond diffus cosmologique (Cosmic Microwave Background (CMB)). Les couches externe et interne sont respectivement caractérisées par $\theta_{\text{ext}} = (T_{\text{ext}}, n_{\text{ext}}, \dots)$ et $\theta_{\text{int}} = (T_{\text{int}}, n_{\text{int}}, \dots)$.

raies de fort RSB au bénéfice des raies plus bruitées, ces dernières pouvant sonder les cœurs denses du nuage. Pour cela, nous saturons le RSB à 10 en augmentant la variance du bruit additif dans le calcul de la vraisemblance. Les données restent intouchées.

Le bruit de calibration ayant été négligé dans nos travaux [11, 12], nous examinons ici les conséquences de ce choix à l'aide de la méthodologie présentée Sec. 2.

3.2 Quantification de la perte de précision

Le modèle sandwich construit dans [12], nécessite l'estimation de 10 inconnues physico-chimiques supplémentaires, en plus de celle des températures et densités volumiques. Dans cet article, nous nous limitons à l'examen de T et n des couches externe et interne du nuage, indicées respectivement ext et int (voir Fig.3). La Fig. 4 montre l'écart relatif ϵ (cf. Eq.(12)) pour les estimations de $\log(T)$ et $\log(n)$ obtenues dans [12].

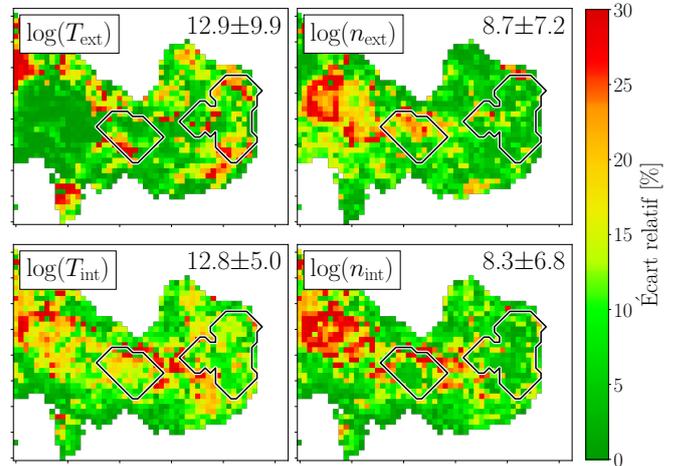


FIGURE 4 : Écart relatif ϵ (%) sur la précision d'estimation de $\log(T)$ et $\log(n)$ dans les couches externe (1^{ère} ligne) et interne (2^{ème} ligne). La moyenne et l'écart-type indiqués sont calculés dans les zones denses, délimitées en noir.

Nous observons un écart relatif des précisions d'estimation $\sim 13\%$ pour les zones les plus denses. Cet écart atteint au plus 30% dans les régions plus diffuses, où le modèle à trois couches est le moins adapté. Ces précisions asymptotiques étant fournies par les BCRs, il convient d'étudier l'efficacité de l'estimateur utilisé, qui néglige le bruit de calibration. Cette analyse est menée dans la section suivante.

3.3 Performances de l'estimateur sur simulations Monte-Carlo

Afin de construire un jeu de spectres correspondant à l'émission provenant d'un cœur dense et froid ($n_{\text{int}}, T_{\text{int}} = 10^5 \text{ cm}^{-3}, 13 \text{ K}$) enveloppé par un gaz chaud et diffus ($n_{\text{ext}}, T_{\text{ext}} = 10^{3.8} \text{ cm}^{-3}, 28 \text{ K}$), nous utilisons le modèle de mesures Eq. (1), où s est décrit par le sandwich (voir [12], Eq. 8). La Fig. 5 présente sous forme d'histogrammes les résultats d'estimation obtenus par Monte-Carlo avec l'estimateur du MV.

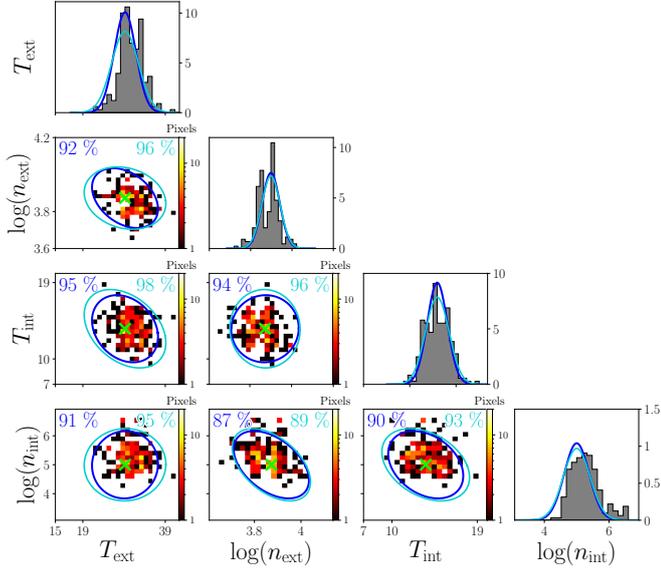


FIGURE 5 : Histogrammes des estimations sur jeu de données simulées (200 réalisations). La valeur vraie est indiquée par la croix verte. Les ellipses de confiance à 99% décrites par la BCR et BCRM sont respectivement tracées en bleu et cyan. Les proportions (%) des estimations contenues dans ces ellipses sont indiquées en haut de chaque panneau.

Comme obtenu dans les études précédentes, les BCR et BCRM sont proches (voir ellipses, Fig.5). La proportion de points contenus dans les ellipses est de 89 à 98% alors qu'elle serait de 99% pour un estimateur efficace. Les histogrammes traduisent un comportement Gaussien de l'estimateur, excepté celui de $\log(n_{\text{int}})$, où un deuxième mode apparaît à $n_{\text{int}} \sim 10^6 \text{ cm}^{-3}$.

4 Conclusion et perspectives

La méthode présentée dans ce travail a permis de quantifier que pour une faible dispersion ($\sigma_c \leq 0.1 \times \mu_c$), négliger un bruit de calibration lors de la construction de l'estimateur du MV induit une perte de précision $\lesssim 13\%$ sur l'estimation des paramètres d'intérêt. Compte-tenu des ordres de grandeur de ces variables astrophysiques et de la complexité du problème considéré, cette perte semble négligeable. Ce constat repose toutefois sur l'analyse des cas particuliers du modèle théorique Eq. (2) et d'un modèle physique adapté aux données radio-astronomiques. Il serait donc intéressant d'examiner le cas de systèmes de mesures différents. En particulier, mener à bien le calcul de marginalisation lorsque le bruit de calibration présente une forte dispersion telle que l'approximation Gaussienne n'est plus vérifiée (≥ 0.3) constituera un défi important.

5 Remerciements

Ce travail a été soutenu par l'Agence Nationale de la Recherche sous la subvention DAOISM (ANR-21-CE31-0010).

Références

- [1] M. S. BARTLETT : An inverse matrix adjustment arising in discriminant analysis. *The Annals of Mathematical Statistics*, 22(1):107–111, 1951.
- [2] L. EINIG et le consortium ORION-B : Deep learning denoising by dimension reduction : Application to the orion-b line cubes. *A&A*, 677:A158, 2023.
- [3] F. GINI, R. REGGIANNINI et U. MENGALI : The modified cramer-rao bound in vector parameter estimation. *IEEE Trans. on Communications*, 46(1):52–60, 1998.
- [4] N L. JOHNSON, S. KOTZ et N. BALAKRISHNAN : *Continuous multivariate distributions*, volume 1, page 209. Wiley New York, 1972.
- [5] S. M. KAY : *Fundamentals of statistical signal Proc. : estimation theory*, chapitre 3. Prentice-Hall, Inc., USA, 1993.
- [6] Y. NOAM et H. MESSER : Notes on the tightness of the hybrid cramer - rao lower bound. *IEEE Trans. on Signal Proc.*, 57(6):2074–2084, 2009.
- [7] L. ORTEGA, J. VILÀ-VALLS et E. CHAUMETTE : Theoretical evaluation of the gnss synchronization performance degradation under interferences. *In ION GNSS+ 2022*.
- [8] P. PALUD et le consortium ORION-B : Mélange de bruits et échantillonnage de posterior non log-concave. *In Actes du GRETSI 2022*.
- [9] C. REN, M. N. EL KORSO, J. GALY, E. CHAUMETTE, P. LARZABAL et A. RENAUX : Performance bounds under misspecification model for mimo radar application. *In 2015 23rd EUSIPCO*, pages 514–518, 2015.
- [10] C. D. RICHMOND et L. L. HOROWITZ : Parameter bounds under misspecified models. *In 2013 Asilomar Conference on Signals, Systems and Computers*, pages 176–180, 2013.
- [11] A. ROUEFF et le consortium ORION-B : Bias versus variance when fitting multi-species molecular lines with a non-LTE radiative transfer model - Application to the estimation of the gas temperature and volume density. *A&A*, 686:A255, 2024.
- [12] L. SÉGAL et le consortium ORION-B : Toward a robust physical and chemical characterization of heterogeneous lines of sight : The case of the horsehead nebula. *A&A*, 692:A160, 2024.
- [13] F. F. S. van der TAK, J. H. BLACK, F. L. SCHÖIER, D. J. JANSEN et E. F. van DISHOECK : A computer program for fast non-LTE analysis of interstellar line spectra. With diagnostic plots to interpret observed line intensity ratios. *A&A*, 468(2):627–635, juin 2007.