

Dissecting the molecular structure of the Orion B cloud: insight from principal component analysis[★]

Pierre Gratier¹, Emeric Bron^{2,3}, Maryvonne Gerin⁴, Jérôme Pety^{5,4}, Viviana V. Guzman⁶, Jan Orkisz^{5,4}, Sébastien Bardeau⁵, Javier R. Goicoechea², Franck Le Petit³, Harvey Liszt⁷, Karin Öberg⁶, Nicolas Peretto⁸, Evelyne Roueff³, Albrecht Sievers⁵, and Pascal Tremblin⁹

¹ Laboratoire d'Astrophysique de Bordeaux, Univ. Bordeaux, CNRS, B18N, Allée Geoffroy Saint-Hilaire, 33615 Pessac, France
e-mail: pierre.gratier@u-bordeaux.fr

² ICMM, Consejo Superior de Investigaciones Científicas (CSIC), 28049 Madrid, Spain

³ LERMA, Observatoire de Paris, PSL Research University, CNRS, Sorbonne Universités, UPMC Univ. Paris 06, 92190 Meudon, France

⁴ LERMA, Observatoire de Paris, PSL Research University, CNRS, Sorbonne Universités, UPMC Univ. Paris 06, École normale supérieure, 75005, Paris, France

⁵ IRAM, 300 rue de la Piscine, 38406 Saint-Martin d'Hères, France

⁶ Harvard-Smithsonian Center for Astrophysics, 60 Garden Street, Cambridge, MA 02138, USA

⁷ National Radio Astronomy Observatory, 520 Edgemont Road, Charlottesville, VA 22903, USA

⁸ School of Physics and Astronomy, Cardiff University, Queen's buildings, Cardiff CF24 3AA, UK

⁹ Maison de la Simulation, CEA-CNRS-INRIA-UPS-UVSQ, USR 3441, Centre d'étude de Saclay, 91191 Gif-Sur-Yvette, France

Received 5 October 2016 / Accepted 13 January 2017

ABSTRACT

Context. The combination of wideband receivers and spectrometers currently available in (sub-)millimeter observatories deliver wide-field hyperspectral imaging of the interstellar medium. Tens of spectral lines can be observed over degree wide fields in about 50 h. This wealth of data calls for restating the physical questions about the interstellar medium in statistical terms.

Aims. We aim to gain information on the physical structure of the interstellar medium from a statistical analysis of many lines from different species over a large field of view, without requiring detailed radiative transfer or astrochemical modeling.

Methods. We coupled a non-linear rescaling of the data with one of the simplest multivariate analysis methods, namely the principal component analysis, to decompose the observed signal into components that we interpret first qualitatively and then quantitatively based on our deep knowledge of the observed region and of the astrochemistry at play.

Results. We identify three principal components, linear compositions of line brightness temperatures, that are correlated at various levels with the column density, the volume density and the UV radiation field.

Conclusions. When sampling a sufficiently diverse mixture of physical parameters, it is possible to decompose the molecular emission in order to gain physical insight on the observed interstellar medium. This opens a new avenue for future studies of the interstellar medium.

Key words. ISM: molecules – ISM: clouds – photon-dominated region (PDR) – ISM: individual objects: Orion B – methods: statistical

1. Introduction

Molecular clouds have a complex structure, with filaments hosting dense cores and immersed in a low density diffuse envelope. Large-scale dust continuum maps obtained with *Herschel* have provided a breakthrough, by showing the tight relationship between the filaments and the dense cores. These maps however do not provide information on the gas dynamics or its chemical composition. Furthermore the relationship between the submillimeter dust emission and the gas column density is affected by the dust temperature and possible variations of the dust emissivity. Molecular line emission maps provide alternative means to study molecular cloud structure and relate it to the flow

kinematics. Molecular line emission is linked to the underlying physical properties of the interstellar medium (ISM), such as density, gas and dust temperatures, UV radiation field, and cosmic ray ionization rate. However these relationships are complex and their detailed study is a field in itself, namely astrochemical modeling (Le Boulot et al. 2012; Agúndez & Wakelam 2013). Further complexity arises when considering radiative transfer to derive line intensities from the local chemical composition and physical structure.

The last few years have seen the installation of new wideband receivers and spectrometers at millimeter and sub-millimeter radiotelescopes. With these instruments, line surveys of several GHz bandwidth and several tens of thousands of spectral channels are the new default mode of observations. Combined with wide field imaging capabilities both for single dish and

[★] Based on observations carried out at the IRAM-30 m single-dish telescope. IRAM is supported by INSU/CNRS (France), MPG (Germany) and IGN (Spain).

Table 1. Properties of the observed spectral lines. The last six columns show the statistics of the data before and after asinh reparametrization.

Molecule	Transitions	Frequency (MHz)	Noise (K)	Original data				After asinh reparametrization			
				Min. (K)	Median (K)	Max. (K)	Std. (K)	Min. (K)	Median (K)	Max. (K)	Std. (K)
¹² CO	$J = 1 \rightarrow 0$	115 271.202	0.09	-0.39	13.40	57.11	10.18	-0.37	2.39	3.32	0.91
¹³ CO	$J = 1 \rightarrow 0$	110v201.354	0.04	-0.19	1.38	36.43	3.27	-0.19	0.97	3.03	0.74
CS	$J = 2 \rightarrow 1$	97 980.953	0.06	-0.36	0.06	15.53	0.48	-0.35	0.06	2.48	0.23
HCN	$J = 1 \rightarrow 0, F = 2 \rightarrow 1$	88 631.848	0.10	-0.58	0.15	10.32	0.39	-0.52	0.15	2.22	0.25
HCO ⁺	$J = 1 \rightarrow 0$	89 188.525	0.09	-0.45	0.26	8.07	0.47	-0.42	0.25	2.07	0.30
SO	$N = 3 \rightarrow 2, J = 2 \rightarrow 1$	99 299.870	0.06	-0.43	0.04	6.46	0.24	-0.40	0.04	1.92	0.17
CN	$N = 1 \rightarrow 0, J = 3/2 \rightarrow 1/2, F = 5/2 \rightarrow 3/2$	113 490.970	0.09	-0.58	0.10	6.33	0.27	-0.52	0.09	1.91	0.20
HNC	$J = 1 \rightarrow 0$	90 663.568	0.08	-0.49	0.07	6.01	0.27	-0.45	0.07	1.88	0.19
CCH	$N = 1 \rightarrow 0, J = 3/2 \rightarrow 1/2, F = 2 \rightarrow 1$	87 316.898	0.12	-0.62	0.08	5.72	0.22	-0.55	0.08	1.85	0.18
C ¹⁸ O	$J = 1 \rightarrow 0$	109 782.173	0.06	-0.30	0.06	5.55	0.42	-0.29	0.06	1.83	0.26
N ₂ H ⁺	$J = 1 \rightarrow 0, F1 = 2 \rightarrow 1, F = 3 \rightarrow 2$	93 173.764	0.08	-0.44	0.00	4.53	0.13	-0.41	0.00	1.70	0.10
CH ₃ OH	$J = 2 \rightarrow 1, K = 0 \rightarrow 0, (A+)$	96 741.375	0.06	-0.34	0.01	2.24	0.08	-0.32	0.01	1.26	0.08

interferometers, hyperspectral imaging is now routinely carried out with these instruments.

The analysis and interpretation of these large datasets, consisting of thousands of spatial positions and tens of thousands of spectral channels, will benefit from the use of statistical tools. Principal component analysis (PCA) is one of the most widely used multivariate analysis method. It has been used to study the ISM using molecular emission maps (Ungerechts et al. 1997; Neufeld et al. 2007; Lo et al. 2009; Melnick et al. 2011; Jones et al. 2012).

In this paper, we address the following question: can PCA provide a method to study the underlying physics of the ISM when applied to a large dataset of molecular emission, without performing either radiative transfer or astrochemical modeling? The article is divided as follows. In Sect. 2 we present the data used in this study. In Sect. 3 we describe the statistical method used in this paper and its implementation. Results are presented in Sect. 4 first by analyzing the output of the PCA, and further by comparing these outputs with independent maps of physical conditions in Orion B in Sect. 5. The last section discusses these results.

2. Data

The data used in this paper was selected from the ORION-B project (PI: J. Pety), which aims at mapping with the IRAM-30 m telescope a large fraction of the south-western edge of the Orion B molecular cloud over a field of view of 1.5 square degrees in the full 3 mm atmospheric window at 200 kHz spectral resolution. Pety et al. (2016) describe in detail the data acquisition and reduction strategies. Table 1 lists the 12 transitions selected in this paper from the already observed frequency range (from 84 to 116 GHz), based on the inspection of the full data cube.

For each line, we focused on emission coming from a limited 1.5 km s^{-1} -velocity range centered on the peak velocity (i.e., 10.5 km s^{-1}) of the main velocity component along the line of sight. Averaging the three 0.5 km s^{-1} velocity channels allowed us to get a consistent dataset from the radiative transfer and kinematics viewpoints. In particular, we avoided the need to disentangle 1) the effects of hyperfine structures of some lines, and 2) the complex velocity structure of the source (Orkisz et al. 2017).

The observed field of view covers $0.81 \text{ deg} \times 1.10 \text{ deg}$ and contains the Horsehead nebula, and the H II regions NGC 2023, NGC 2024, IC 434, and IC 435. The angular resolution ranges from 22.5 to $30.5''$. The 12 resulting maps have a common pixel size of $9''$ that corresponds to a Nyquist sampling for the highest frequency line observed (¹²CO(1–0) at 115.27 GHz). The maps thus contain 315×420 pixels. At a distance of $\sim 400 \text{ pc}$ (Menten et al. 2007), the maps give us access to physical scales between $\sim 50 \text{ mpc}$ and 10 pc .

Figure 1 shows the 12 maps of the resulting brightness temperature multiplied by an ad hoc factor in order that they can share the same color look-up table, even though the intrinsic brightness temperatures of the different lines differ by more than one order of magnitude. The relative calibration of the different lines is excellent because they were observed with the same telescope at almost the same time, since the observed bandwidth was covered in only two frequency tunings. The noise for each map, along with the minimum, median, maximum, and variance values are listed in Table 1. The noise was computed by fitting a gaussian function to the negative part of the histogram of pixel brightnesses. This enabled us to compute the noise without needing to mask out the emission.

3. Principal component analysis

3.1. Principle

We use the following standard statistical terms: the dataset is composed of “samples” each described by individual “features”. In our case, each spatial pixel is a sample and each line intensity is a feature (the full dataset thus corresponds to a data matrix of $132\,300$ samples times 12 features).

PCA is a widely-used multi-dimensional analysis technique (Jolliffe 2002), that can be defined in several mathematically equivalent ways. It aims at finding a new orthogonal basis of the feature space (whose axes are called principal components, or PCs), so that for each k , the projection onto the hyperplane defined by the first k axes is optimal in the sense that it preserves most of the variance of the dataset (or equivalently that the error caused by this projection is minimal). PCA thus defines successive approximations of the dataset by hyperplanes of increasing dimension.

This is equivalent to the diagonalization of the covariance matrix, so that the principal components are naturally uncorrelated. It can be thought of as finding the principal axes of inertia

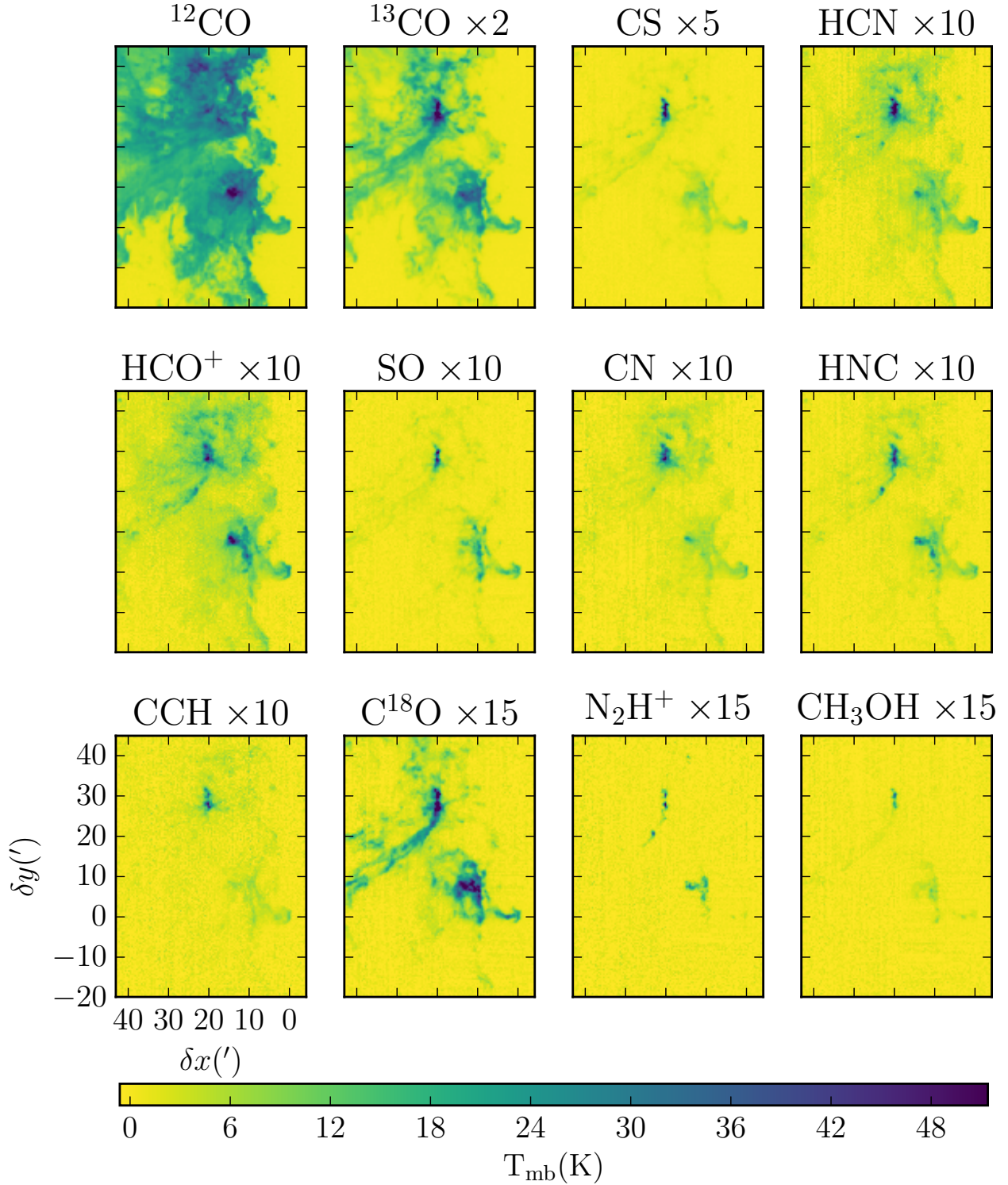


Fig. 1. Maps of molecular emission in Kelvin main beam temperatures.

of the cloud of samples about their mean in feature space, and is thus a way to analyze the covariance structure of the data. The principal components are ordered by decreasing projected variance. As a result PC1 is the axis of largest variance in the data. PC2 is then the axis of largest variance at constant PC1 (orthogonal to PC1) and so forth. Neglecting the axes of lowest variance then allows the definition of a low-dimensional hyperplane in which the dataset is approximately embedded. An important property to keep in mind is the linearity of PCA, namely that

it defines low-dimensional hyperplanes, and not general low-dimensional hypersurfaces.

A common variant¹ in the application of PCA is to normalize the variations of the dataset around the mean by the standard deviation of each features, before applying the PCA. This amounts to diagonalizing the correlation matrix instead of the

¹ This variant in the application of PCA goes back to one of the two earliest descriptions of PCA: [Hotelling \(1933\)](#).

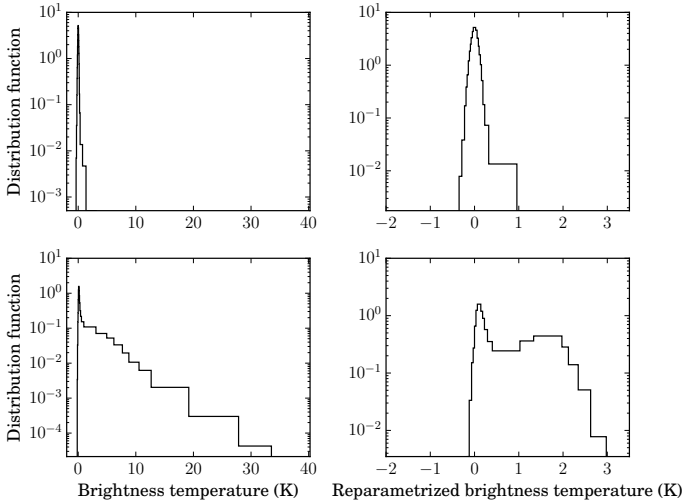


Fig. 2. Effect of the asinh renormalization on the intensity distributions. *Left column:* before renormalization, *right column:* after renormalization, *top row:* N_2H^+ , *bottom row:* ^{13}CO .

covariance matrix. The correlation-based variant allows to avoid having one feature dominate the variance, and is appropriate if the relative scales of the features are not relevant for the purpose of the analysis. As the relative intensity scales of the different molecules used here are largely affected by properties (dipole moment, elemental abundances, ...) that are not relevant for our analysis of the chemical variations across the map, we use here the correlation-based version of PCA.

In this work, we used the PCA implementation available in the Python package `scikit-learn` (Pedregosa et al. 2011), which uses a singular-value decomposition to compute the principal component axes.

3.2. Reparametrization of input data

3.2.1. The need of a reparametrization

As seen in Table 1, some of the tracers have large dynamical ranges (two orders of magnitude for $^{13}\text{CO}(1-0)$ and $^{12}\text{CO}(1-0)$).

Figure 2 shows the histogram of the brightness temperatures of two lines with contrasting behavior in our dataset, namely $^{13}\text{CO}(1-0)$, and $\text{N}_2\text{H}^+(1-0)$. As the dynamic range is large both in intensity and number of pixels per bin, these histograms use the Bayesian blocks algorithm (Scargle et al. 2013, using here the Python implementation from AstroML; see Vanderplas et al. 2012; Ivezić et al. 2014), which adapts the bin width to the underlying distribution. Although the histogram of the $\text{N}_2\text{H}^+(1-0)$ is Gaussian to first order, the histogram of $^{13}\text{CO}(1-0)$ exhibits heavy tails similar to power laws. As a result, extreme intensity values might dominate the covariance structure of the data, and hide the variations at the more common lower intensity values.

From the physical viewpoint, taking the logarithm of the brightness temperature is also desirable. PCA is a linear technique which decomposes the data as a sum of uncorrelated components. Applying it to the logarithm of the data allows a decomposition as a product of factors, and thus describes the data structure in terms of ratios, products, and power laws, which is more adapted to the underlying radiative transfer and chemical effects. Taking the logarithms of the data is the equivalent in astrochemistry to performing color–color magnitude diagram

analysis in optical or UV studies. Pety et al. (2016) show that the line-integrated brightness temperatures of our dataset are, to first order, correlated to the column density of matter along the line of sight. We expect this aspect to appear in our PCA analysis, whereas second order chemical variations around this trend, which would be revealed by line ratios, are thus better described as multiplicative (rather than additive) factors.

3.2.2. Impact of noise

The presence of noise causes the possibility of negative values in those pixels where some lines are undetected. The logarithm transform cannot be applied to these negative noise values. In addition, because the logarithm stretches the lowest values compared to the largest ones, it will also tend to stretch the positive noise values of the undetected pixels, and this gives them more weight in the covariance of the data.

There may be two different reasons for a non-detection: either the measurement is not sensitive enough to detect the line or the region just does not contain the species that emits the line. The latter case happens particularly in H II regions (e.g., IC 434), in which $^{12}\text{CO}(1-0)$ is photo-dissociated by the far UV photons. In this particular case, we could remove from our dataset all the samples (pixels) where no $^{12}\text{CO}(1-0)$ is detected. This just assumes that no $^{12}\text{CO}(1-0)$ detections at high sensitivity imply the absence of molecular gas. However, this method is not generic. For instance, $\text{N}_2\text{H}^+(1-0)$ is only detected in dense cores (Pety et al. 2016), and restricting ourselves to these regions would drastically limit the scope of our study. If we wish to use this important tracer of the molecular gas while still covering the range of different chemical regimes present in our full map, we thus need to find an alternative. Moreover, in the 3 mm band, radio recombination lines, which emit in H II regions, could in principle be added to our dataset to study the formation of molecular gas.

Adding a thresholding step before taking the logarithm will only worsen the scope of undefined values. Although there are PCA methods (e.g. Ilin & Raiko 2010) that can take missing datapoints into account, they rely on the fact that these missing points have the same statistics than the measured points (i.e., when a value is missing, it is independent of the actual value). This is clearly not the case here because the missing values (undetected lines) are missing due to them being below the sensitivity threshold. These are called “censored values” in statistics. We thus searched for a function that is linear around zero (in the noise-dominated domain), and that is asymptotically equal to a logarithm for large values (compared to the noise level). The inverse hyperbolic sinus function, $\text{asinh}(x)$, fulfills these conditions. We thus used the following function to reparametrize the data before applying PCA

$$T(x) = a \text{asinh}(x/a), \quad (1)$$

where the parameter a is the typical value for which the function’s behavior changes from a linear to a logarithmic regime (Fig. 3).

The only free parameter in the method is the threshold a . Appendix A discusses our choice, namely $a = 8 \times 0.08 = 0.64 \text{ K}$ = eight times the median noise of the dataset. In short, we select the value of a that maximizes the correlations of the first three principal components with independent known measurements of the column density, volume density, and UV illumination (Sect. 5). This appendix also demonstrates that our results are quite insensitive to the exact value of a .

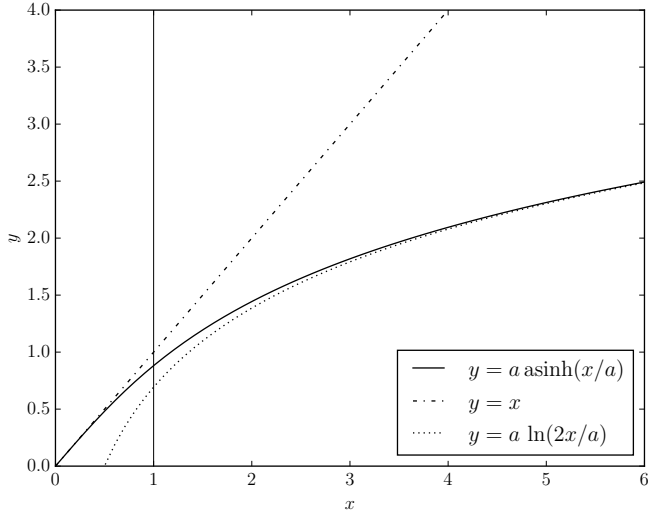


Fig. 3. Plot of the asinh function (solid line) showing the asymptotes when $x \rightarrow 0$ (dash dotted line) and $x \rightarrow +\infty$ (dotted line). The parameter a (here $a = 1$) is traced with a thin vertical line.

The right column of Fig. 2 shows the result of the asinh transformation on the intensity distributions of two transitions representative of bright ($^{13}\text{CO}(1-0)$) and weak ($\text{N}_2\text{H}^+(1-0)$) averaged lines. In the case of the bright lines, the dynamic range is drastically reduced, with the heavy tail being transformed into a second peak in the distribution but with no values above three. In the case of N_2H^+ the distributions before and after reparameterization are very similar. Figure 4 shows the 12 maps of the molecular emission after reparameterization by the asinh function, but before the normalization step of the PCA analysis. The brightness temperatures of all the maps have been compressed between about -0.5 and 3.5 , with low signal-to-noise brightness temperatures between -0.5 and 0.5 being mostly untransformed.

4. Results

The PCA method exposes the correlations between the line brightness temperatures. The derived PCs give the main axes of correlated variations in the data set. As such, PCA does not directly yield physical information underlying the dataset. In this section, we describe the results of the application of this statistical method to our dataset and we start to discuss their possible physical interpretation based on our a priori astrochemical knowledge. The possible relations between the PCs and physical variables are investigated in a later section.

4.1. Correlation fraction explained by the different principal components

Figure 5 shows the percentage of the correlation explained by each PC (as a function of the principal component number) along with the cumulative explained correlation as a function of the number of principal components kept in the decomposition.

The first principal component explains the majority (60%) of the total correlation present in the original dataset. Thus a large part of the variations in the dataset occur along a single axis (i.e. all lines are strongly correlated to each other). The second principal component accounts for about 10% of the correlation. It is significantly less than the first component, but more than any other components. PCs 3, 4 and 5 correspond to similar amounts

of correlations (around 5% each) and PC6 slightly less (3.3%). PC1 to 6 collectively explain more than 90% of the correlation in our dataset. The remaining PCs have similar low amounts of explained correlation (from 2% for PC7 to 0.9% for PC11).

4.2. Discussion of the principal components

The PCs defined by our analysis represent new axes in the feature space (the full 12 PCs are simply a rotation of the initial basis of the feature space), deduced from the data itself. They can thus be expressed in terms of the original axes, as a linear combination of the (transformed) line intensities. Figure 6 displays the quantitative contribution of each initial feature (line) to each PC. An alternative view of the relationship between the PCs and the line intensities, namely the correlation wheels, is presented in Fig. 7.

Each sample (pixel brightness) can then be projected on the new axes, providing new coordinates commonly called “component scores”. The PCA method considers the pixels as independent samples, and thus ignores the spatial structure of the molecular emission. It is nevertheless possible to reconstruct the maps of the component scores. Figure 8 shows these projected maps. The chosen color look-up table emphasizes that positive and negative values of the projected maps, which correspond to variations above and below the average along the considered axis respectively, clearly extract a different spatial pattern per principal component.

The first principal component is a linear combination of all tracers, with similar positive weights for all lines (Fig. 6). It thus describes correlated variations of all molecules, and these account for most of the variations in the dataset, which is a natural consequence of all lines being well correlated (positively) to each other. Pety et al. (2016) show that the emission of all lines is correlated to first order with the column density of matter along the line of sight. The first component is thus probably related to the total column density, whose increase causes, to first order, an increase in all lines. This is because in the linear approximation of PCA, non linear effects such as saturation of the ^{12}CO line are not captured. The corresponding component map (Fig. 8) indeed resembles a map of column density. This relation between PC1 and the total column density of matter will be investigated more quantitatively in Sect. 5. We note that as PC1 has only positive coefficients for all lines, orthogonality ensures that all other PCs will represent contrasts between different lines.

The second principal component represents the axis of largest variation at constant PC1 (orthogonal to PC1). This axis of variations is dominated by positive contributions of N_2H^+ and CH_3OH , and negative contribution of ^{12}CO and ^{13}CO . The first two tracers are chemically associated with dense and cold regions of the ISM. For instance, because N_2H^+ is easily destroyed by CO, it can thus only be abundant in the gas phase when CO has been depleted on the grain surfaces (Pety et al. 2016). The component map shows strong positive values highlighting known dense cores, including the clumps in the head and in the neck of the Horsehead (Ward-Thompson et al. 2006).

The third principal component shows positive contributions of CCH and CN that are known to be sensitive to UV illumination, and negative contributions of N_2H^+ , CH_3OH and of the CO isotopologues, that trace gas shielded from the UV field. This component thus probably traces the chemical specificities of UV illuminated gas. The component map clearly shows positive values at the eastern edge of the cloud, illuminated by σ Ori, and in the star-forming region NGC 2024.

The fourth principal component is particular in the fact that its map almost only shows large (positive or negative) values in

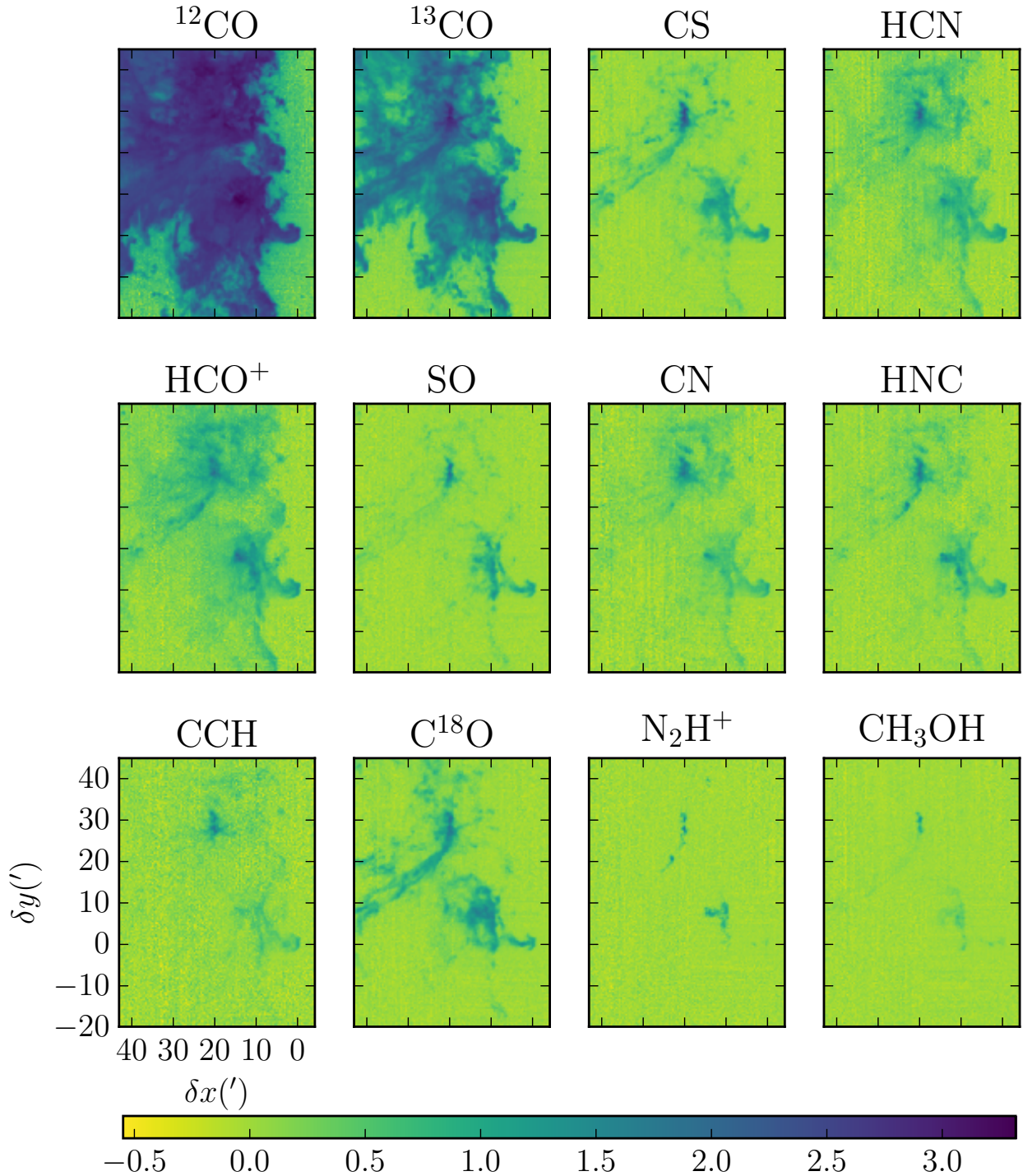


Fig. 4. Maps of molecular emission after reparametrization by the asinh function.

the regions of large positive values of PC2. It thus highlight further chemical variations inside dense cores. This component is completely dominated by opposite contributions of CH_3OH and N_2H^+ and thus traces variations in the ratio of these two lines. The component map seems to highlight smaller-sized cores embedded in some of the clumps revealed by PC2, and thus probably highlights the chemistry of the densest cores (larger N_2H^+ to CH_3OH ratios).

The fifth principal component shows positive contributions of sulfur species (CS and SO), and C^{18}O , and negative contributions of ^{12}CO , CCH, and CH_3OH . Its large positive values

highlight larger scale regions embedding the dense clouds shown by PC2 (but with negative values where PC2 is very large), and this PC could thus trace the chemistry of moderately dense gas.

The sixth principal component shows negative contributions of HCN, HCO^+ , and CN, which can all originate in photochemistry, and positive contributions from CCH, C^{18}O , ^{13}CO , and SO. Although the latter are usually associated with a larger amount of shielded gas, CCH can also be bright in UV illuminated regions. The component map shows a wide blue region around NGC 2024, similar to the large warm dust region seen in the dust temperature map of the region (Schneider et al. 2013).

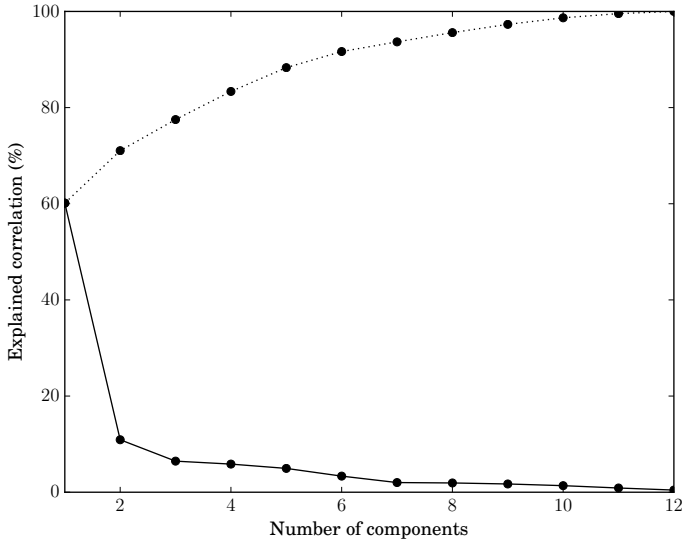


Fig. 5. Percentage of the explained correlation as a function of the number of components in the principal component analysis, dotted line: cumulative percentage.

This region could also be related to the radiation field, but trace a different aspect from PC3, that is characterized by lower CCH intensities relative to the other lines.

The remaining components are more difficult to interpret, but tend to describe opposite variations in pairs of lines that varied together in previous PCs. PCs 7, 9 and 10 display opposite variations in pairs of lines of the group HCN, HNC, CN and HCO^+ , whose variations were correlated in the previous PCs in which they had large weights (PC1 and 6). PC8 shows anticorrelated variations of SO and C^{18}O , and its component map shows a striking spatial pattern: negative values (high SO/ C^{18}O ratios) in the Horsehead, the molecular gas at the base of the Horsehead, and the small scale clumps in NGC 2024; and positive values (low SO/ C^{18}O ratios) in a dense filament stretching away from NGC 2024. PC11 is strongly dominated by CS, and thus shows specific variations of CS, mostly uncorrelated with the other lines (somewhat anticorrelated with C^{18}O), and that were not described by the previous PCs. Its component map shows small scale spots of positive values, mostly surrounding NGC 2024 and NGC 2023. The fact that it appears so late in the decomposition can be explained by the small size of the highlighted region, having little weight in the correlation matrix. PC12 is completely constrained by orthogonality to the previous PCs and is thus only an artifact of previous PCs.

4.3. Studying the effect of noise

The noisy nature of our data can have two kinds of effects. It can first induce variability in our results (the results would vary for a different realization of the random noise). We verified the stability of our results by using a bootstrapping method. Bootstrapping is a method of choice to compute uncertainties on an estimator (here the PCA components) when the distribution of estimator values cannot be assumed to follow a simple distribution (Feigelson & Babu 2012). The idea of bootstrapping is to use a Monte Carlo method to create new resampled datasets of the same size as the original dataset by sampling with replacement from the original dataset. We constructed 5000 such bootstrapped datasets and ran the PCA algorithm on each. To avoid overestimating the uncertainty, and because the PCA is

invariant through the change of sign of the PCs, we ensured that the signs were consistent before computing the distribution of the PC coefficients. The results are presented both for the eigen-spectra in Fig. 6, and the correlation wheels in Fig. 7.

The PC coefficients appear overall very stable, the variances are completely negligible for PC 1 and 2, and very small for PC 3 to 6. Only PCs 7 and 8 show higher variability. This can be understood as PCA results are particularly sensitive to noise when two PCs correspond to very close eigenvalues, and PCs 7 and 8 have the closest eigenvalues with respectively 2% and 1.9% of the total correlation. Indeed, PCs with equal eigenvalues are degenerate in the sense that any basis of the subspace they define satisfies the definition of PCA. As a result, when eigenvalues are not exactly equal but very close, the noise can result in a random rotation of this group of PCs inside their subspace. Our results, which focus on the first few PCs, are thus unaffected by noise variability.

The second possible effect of noise is to bias the results. Principal component analysis is unbiased if the noise is spherical (i.e., has equal variance in all directions) in the final dataset on which PCA is applied (i.e., after standardization in our correlation-based variant). In this case, the noise can only hide the lowest PCs (that describe variations smaller than the noise level) and make them degenerate. In our case however, the noise levels on the different molecular lines are initially close but not equal (variations by a factor of three at most, as demonstrated in Table 1). The non-linear reparametrization keeps these relative variances. Finally, the last normalization step (by the standard deviation) results in final noise variances that are proportional to the ratios of noise standard deviation to total standard deviation of the reparametrized intensities. These differ by up to a factor of 14.8 between the lines, and possible biases may be present in our results, giving higher weight to the lines with the largest ratios of noise variance to total variance. However, it was not possible with the PCA method to avoid giving higher weight to the brightest lines (which led us to use the correlation-based PCA), and at the same time ensure equal noise on all variables. We note that previous PCA studies of molecular clouds were less concerned by noise-induced bias as they only used lines that were clearly detected in all pixels. We chose to perform the PCA on the full region, which led to the identification of two PCs associated with dense core chemistry.

5. Correlation of the principal component maps with independently measured physical parameters maps

In the previous section, we combined two sources of information to interpret the main principal components: 1) astrochemistry, which teaches us that some molecules trace certain physical conditions, and 2) Orion B is an extremely well-studied source, implying that its spatial structure is well known. For instance, the molecular cloud is known to be illuminated by well-defined young massive stars (discussion in Pety et al. 2016). This allowed us to infer a link between the first three principal components and physical parameters such as the column density, the volume density, and UV illumination. In this section, we will quantitatively assert these potential relations by studying the correlation of each component map with a set of independently measured maps of physical parameters.

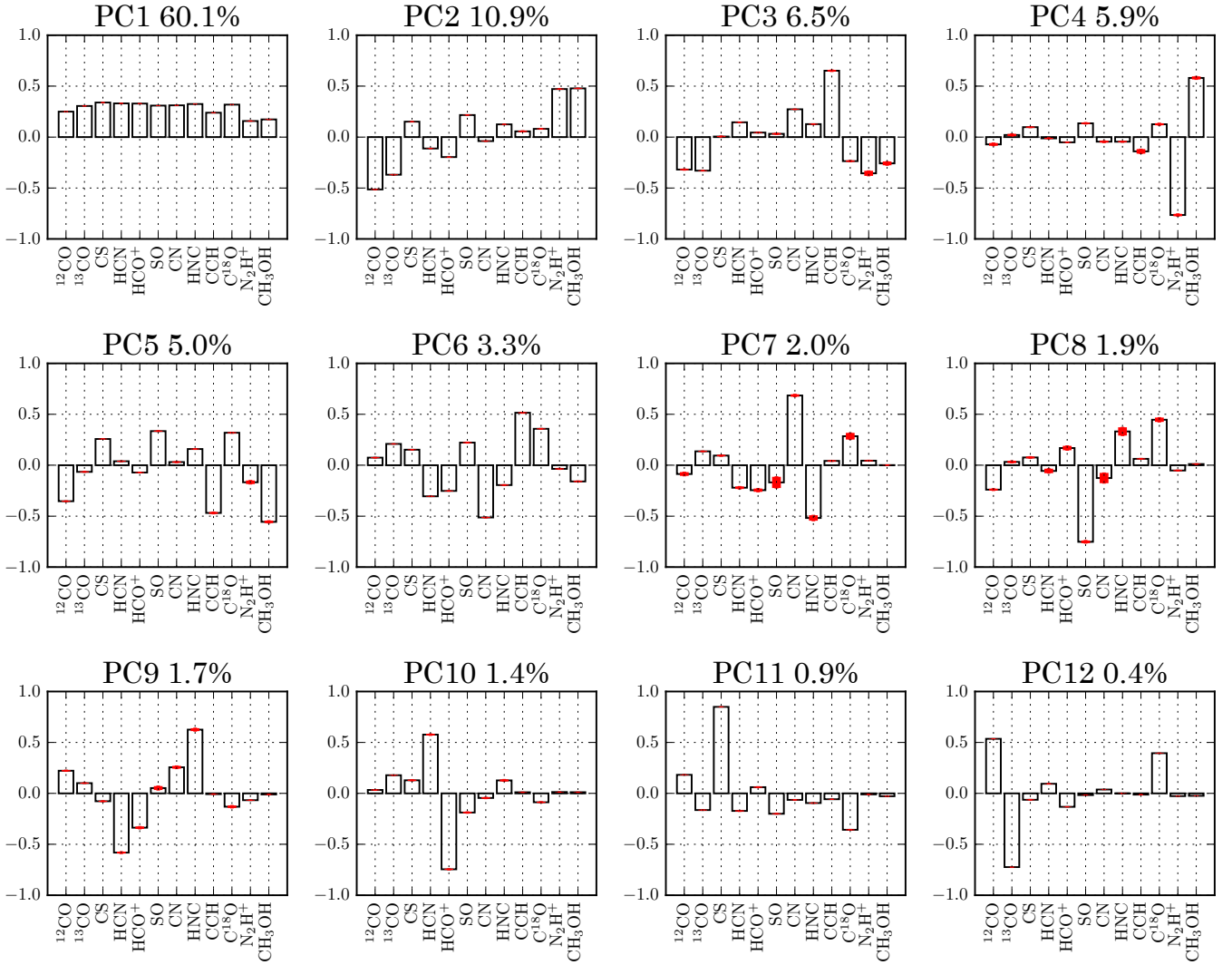


Fig. 6. Bar plots showing the contribution of each line intensity to each principal component (with the fraction of the total correlation accounted for by each PC given as a percentage). The uncertainties (standard deviations) shown in red are obtained by bootstrapping as described in Sect. 4.3.

5.1. Independent measure of the physical parameters

The goal of this section is to find the principal component that is best associated to each of the physical parameters, not to assign an absolute physical meaning to some of the components. It is therefore not necessary to have absolute values of the independently measured physical parameter maps. Only the relative variation of each physical parameter is required to compute the correlation coefficient. Figure 9 shows the different maps of the physical parameters that we will correlate with the first three principal components. This section describes how these maps were obtained.

5.1.1. Column density

The dust column density map is from the Hershel Gould Belt Survey (PI: P. Andre) Orion B map (André et al. 2010; Schneider et al. 2013)². This map was obtained by fitting the far infrared spectral energy distribution by greybodies. We applied a logarithmic scaling to the data to reduce the dynamical range. The resulting map is plotted in the left panel of Fig. 9.

² <http://www.herschel.fr/cea/gouldbelt/en/>

5.1.2. Volume density

Volume density is a difficult quantity to measure because one needs both a mass estimate and an associated volume. Density is thus dependent on the scale that it is computed at. We used the catalog of cores identified and characterized in Kirk et al. (2016) and computed masses from each cloud's 850 μm flux using their equation three. To do this, we assumed a common temperature of 17 K for all clouds. From this mass and their observed size estimates we computed a volume density for each of the dense cores in our observed field of view. In this case correlation could not be carried out over the full map but we correlated the density measured for each core with the value of the principal components measured in the nearest pixel. The data is shown as a scatter plot in the middle panel of Fig. 9.

5.1.3. UV radiation field

We computed the UV radiation field by using the fact that PAH emissivity is roughly constant per unit H and unit radiation field (Draine & Li 2007). In practice, we used the WISE (Meisner & Finkbeiner 2014) 12 μm maps divided by the column density, and clipped to a maximum value of 10^{22} cm^{-2} .

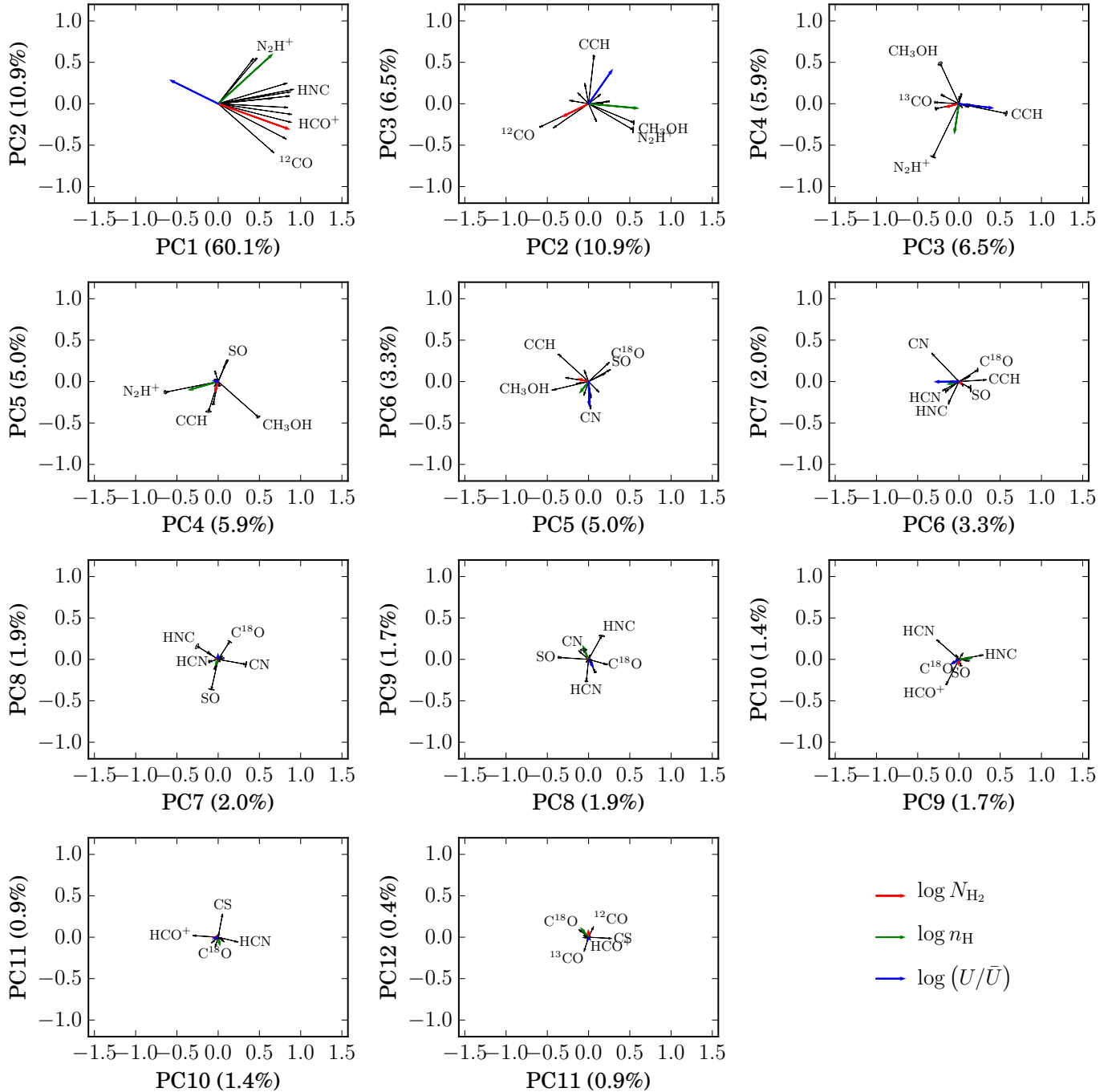


Fig. 7. Correlation wheels, showing the initial line intensities as vectors having as coordinates their correlation coefficients to each PC, represented in the planes of successive pairs of PCs. Uncertainties from our bootstrapping analysis (see Sect. 4.3) are presented as thin black contours around the arrows' heads (isodensity contours containing 68% of the distribution). Also represented in colored arrows are the correlations of our independent physical parameters with the PCs (red: $\log N_{\text{H}_2}$, green: $\log n_{\text{H}}$, blue: $\log(U/\bar{U})$).

We do not claim to have an absolute value of the UV radiation field but a quantity that should be proportional to it. The quantity $\log(U/\bar{U})$ where \bar{U} is the mean value of U is shown in the rightmost panel of Fig. 9.

The proper way to compute the UV radiation field from PAH emission would be to divide by the volume density but as we discussed in the previous paragraph, it is not possible to get a full map of volume density. We chose to use column density as a proxy for volume density even though it entails strong constraints on the spatial distribution of the gas along the line of sight. Since we are interested in relative variation of density and not absolute

values it is sufficient to assume that the matter is clustered into clouds that are of similar spatial extents.

5.2. Correlation of principal component maps with physical parameters

We computed the Spearman's rank correlation coefficient between each pair of principal component maps and physical parameters maps. We used Spearman's rank correlation instead of the Pearson linear correlation coefficient because the potential relations between the principal components and the physical

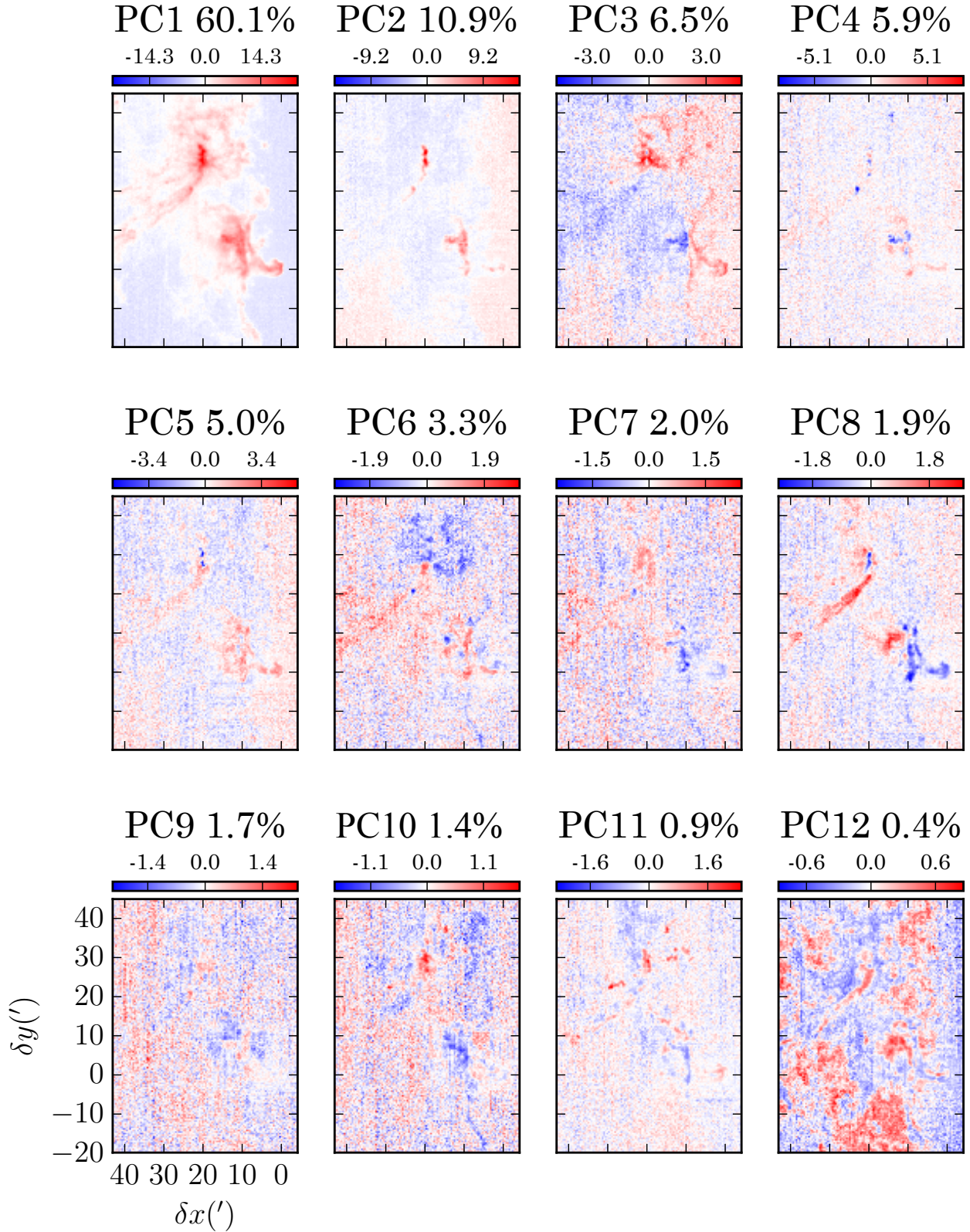


Fig. 8. Principal component maps. These maps represent the value of each observed pixel when they are projected in the space of the principal components.

parameters are most certainly non linear in nature. The rank coefficient used is only sensitive to the ordering of the values and is thus not affected by the possible non-linearities of the correlation. Table 2 summarizes all these values and Fig. 10 shows the scatter plots for the most significant correlations discussed in

the next paragraphs. An alternative way of exploring the correlations between the independent physical parameters and the PCs is to represent the correlation between each physical parameters and the PCs in the correlation wheels of Fig. 7.

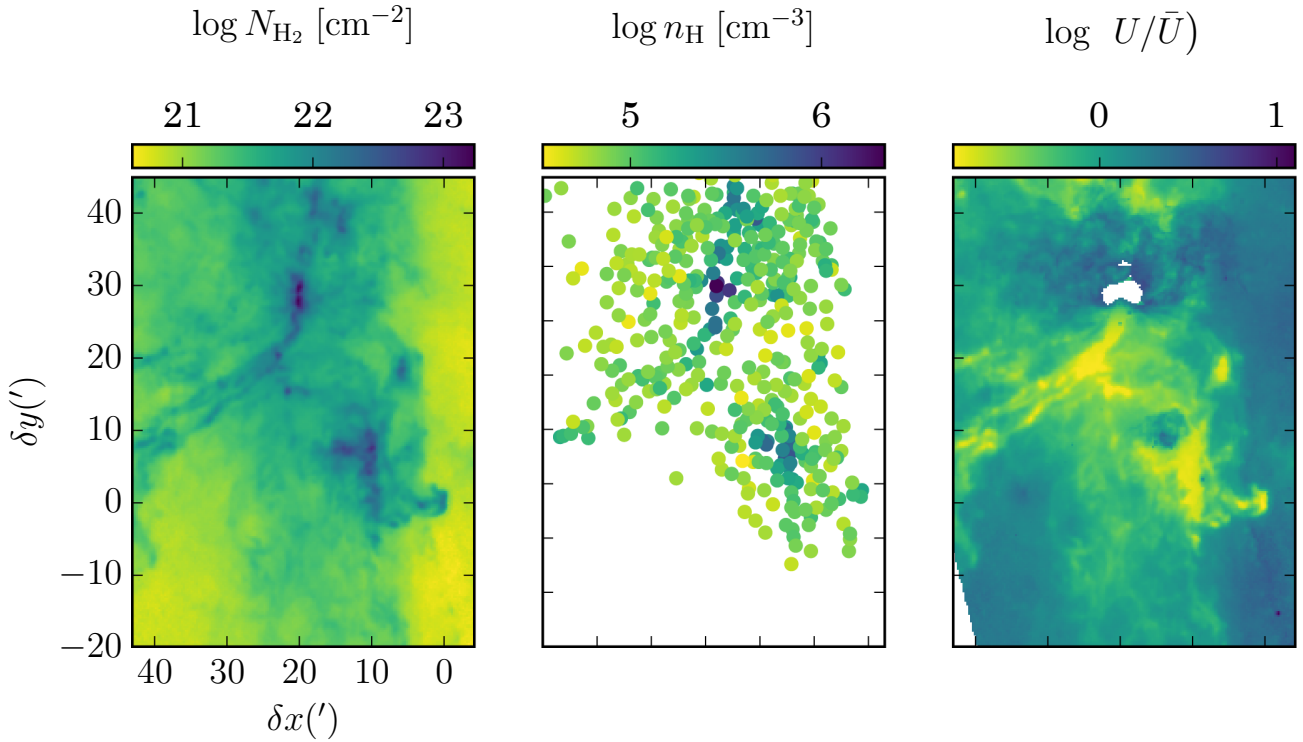


Fig. 9. Maps of the independently measured physical parameters, H_2 column density (left), volumic density (middle), UV illumination (right).

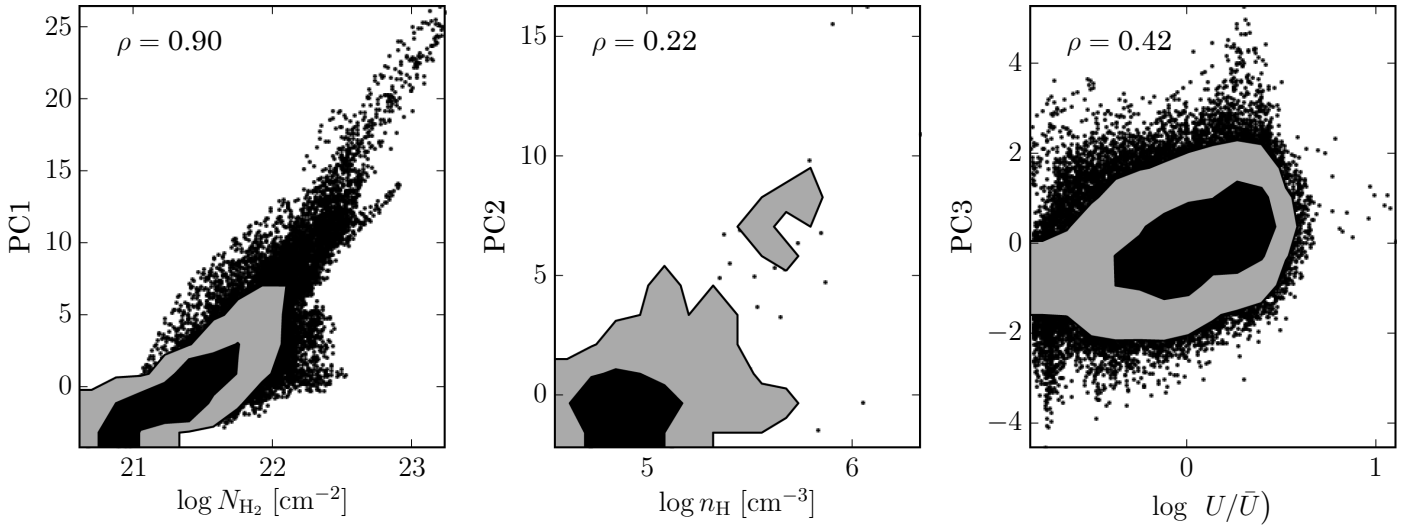


Fig. 10. Scatter plots of the first three principal components with the independent physical parameters. Contours in black and gray correspond to 68% and 95% of the samples respectively.

For this analysis, it must be kept in mind that while the principal components are necessarily uncorrelated, the physical parameters considered here are correlated: $N(H_2)$ is an integral of n_H along the line of sight and the two are thus strongly correlated, whereas U is inversely proportional to $N(H_2)$ by construction and they are thus anticorrelated. As a result, the principal components will tend to represent the uncorrelated part of the variations of the underlying physical parameters.

Column density: The component map showing the highest correlation coefficient with $N(H_2)$ is PC1. Spearman’s rank correlation coefficient is extremely high at 0.90, and the scatter plot

(Fig. 10, left panel) shows a strongly linear relation between PC1 and $\log(N(H_2))$.

Because PC1 is the first PC (axis of largest variation), it is unaffected by the decorrelation constraint that affects the other PCs. This first principal component can thus be interpreted as a global measure of total column density, as suspected in our previous discussion. Since n_H and U are positively and negatively correlated with $N(H_2)$, respectively, these physical parameters exhibit relatively strong positive and negative correlations with PC1.

Volume density: The PC most correlated to n_H is also PC1, due to the large correlation between $N(H_2)$ and n_H . The next

Table 2. Spearman’s rank correlation coefficient between the principal components and the physical parameters.

	$\log N_{\text{H}_2}$	$\log n_{\text{H}}$	$\log(U/\bar{U})$
PC1	0.90	0.43	-0.66
PC2	-0.57	0.22	0.43
PC3	-0.20	0.06	0.42
PC4	-0.01	-0.23	-0.04
PC5	-0.16	0.02	0.09
PC6	0.04	-0.07	-0.26
PC7	-0.02	-0.03	0.00
PC8	-0.02	0.05	0.12
PC9	-0.02	0.06	-0.11
PC10	-0.06	-0.10	-0.07
PC11	-0.03	-0.13	-0.04
PC12	0.04	0.11	-0.05

principal components most correlated with our limited sample of volume density measurements are PC2, which shows a Spearman’s rank correlation coefficient of 0.22, and PC4, with a Spearman’s rank correlation coefficient of -0.23 .

As was discussed in Sect. 4, PC2 and PC4 both trace chemical differences typical of dense cores. PC2 and PC4 can thus be interpreted as indicator of the presence of dense cores. We note that this comparison was only done with a limited sample of rather dense clouds. We can thus only say that PC2 traces increased density among dense clouds. Because of the opposite sign of the correlation of the density with PC4, negative values of this PC probably trace an even higher density regime. As noted before, the behavior of PC2 in less dense region is probably anticorrelated with density, and these PCs are thus only indicative of density in the high density regime.

Radiation field: For the radiation field, the most correlated PC is again PC1 (negative correlation) because high column density tends to result in highly shielded gas. Not considering PC1 and PC2, the third principal component shows the highest correlation with our estimation of the radiation field, with a Spearman’s rank coefficient of 0.42. It thus describes the part of the radiation field variations that are not correlated with the cloud column density. As a result, PC3 highlights the part of the cloud where specific sources cause increased illumination. A strong positive correlation (0.43) with PC2 is also found, which is most likely an artifact due to the positive values of PC2 in the diffuse regions surrounding the molecular cloud (where most of the lines involved in PC2 are undetected, making PC2 irrelevant). PC6 also has a significant correlation with the radiation field (-0.26). These results thus confirm our previous discussion of PC3 and PC6.

These results can be inferred graphically from the correlation wheels of Fig. 7, in which the colored arrows tracing the location of $N(\text{H}_2)$, n_{H} , and U in the PC space have a significant size only for the first four PCs. Furthermore, each arrow is roughly aligned with one of the PCs: PC1 for $N(\text{H}_2)$, PC2 for n_{H} and PC3 for U .

6. Discussion

6.1. Comparison with other works

PCA has been extensively used in astronomy as a multivariate analysis tool starting from the work of Deeming (1964)

on the classification of stellar spectra. Its use for the study of molecular maps of the ISM is more recent, starting with the work of Ungerechts & Thaddeus (1987). Notable studies include Neufeld et al. (2007), Lo et al. (2009), Melnick et al. (2011), and Jones et al. (2012).

We first discuss the common points between these studies. On a technical aspect, all these studies apply only subtraction by mean and normalization by variance, and do not attempt to introduce a non linear reparametrization of the observed intensities. The effect of noise is considered by limiting the number of observed lines to the set of brightest tracers (Lo et al. 2009) and masking regions of low emission (Jones et al. 2012). All of these studies identify the utility of PCA as a means of studying the correlations between molecular lines by studying the commonality between tracers and their variation, and as a tool to identify regions interesting for further study. With the notable exception of Ungerechts & Thaddeus (1987), rarely a discussion is made relating the principal components with the underlying physical parameters of the ISM. However, specific correlations or anticorrelations are often discussed in a chemical view or by invoking opacity effects (Lo et al. 2009).

Ungerechts et al. (1997) present a dataset of 360 spatial points in 32 lines of 20 chemical species including isotopologues toward the Orion A molecular cloud with the 14 m FCRAO telescope. Using PCA they show that the chemical abundances of most species stay similar for the Orion ridge, and that the main differences stand up for the BN-KL region. They note that their first three PCs contain 80% of the observed correlation and they use their component maps mainly to identify regions for further astrochemical study. They nevertheless discuss that data mostly lie in a 3D space spanned by the first three PCs because the molecular emission probably depends on three physical parameters of the ISM, namely the column density, volumic density, and gas temperature. Melnick et al. (2011) compared the distribution of the ground-state transition of water vapor with that of the ground state transition of N_2H^+ , CCH, HCN, CN, and $^{13}\text{CO}(5-4)$. Water vapor is found to best correlate with species like $^{13}\text{CO}(5-4)$ and CN, tracing the cloud surface up to a few magnitudes of extinction, and is poorly correlated with N_2H^+ tracing the shielded regions. Using MOPRA, Jones et al. (2012) have mapped the central molecular zone (CMZ) near the center of the Galaxy in 20 spectral lines in the 85.3 to 93.3 GHz range. They performed a PCA analysis using the strongest eight lines (HCN, HCO^+ , HNC, HNCO, N_2H^+ , SiO, CH_3CN , and HC_3N) in the restricted area around SgrB2 and SgrA, where the N_2H^+ line is stronger than 10 K km s^{-1} . The analysis recovers the overall similarity of the line maps. The main differences are found in the SgrA and SgrB2 cores between the bright lines HCN, HNC, HCO^+ , and the other species, and is attributed by Jones et al. (2012) to a difference in opacity. The other PCA components reveal specific regions where CH_3CN , HNCO, and SiO abundances are enhanced, possibly due to shocks or hot cores. Lo et al. (2009) studied the G333 molecular cloud with MOPRA. The PCA is performed on eight molecular lines with high S/N ratio (^{13}CO , C^{18}O , CS, HCO^+ , HCN, HNC, N_2H^+ , and CCH). The PCA analysis reveals differences between the regions traced by CCH and N_2H^+ . In star forming regions it also reveals an anticorrelation between ^{13}CO , C^{18}O , and N_2H^+ , and between N_2H^+ and HCO^+ . PCA is also used by Neufeld et al. (2007) to separate the different regions impacted by supernovae shock waves.

While the analysis discussed in the previous paragraph used integrated line intensities, PCA has also been used on spectral line profiles as a mean to extract information on the

spatial properties of the turbulence (Heyer & Peter Schloerb 1997; Roman-Duval et al. 2011; Brunt & Heyer 2013), to study line absorption depth (Neufeld et al. 2015), or to measure cloud properties (Rosolowsky & Leroy 2006). To our knowledge no PCA analysis takes into account the full velocity profile of the molecular emission at every spatial position. Further inquiry on this subject is required as it can add a further dimension, namely the shape of the line profiles, to study the emission correlations.

6.2. Non-linearities and multiple physical regimes

Two important properties of PCA must be kept in mind. The first is that it is a linear method. It distinguishes the axes of variations in the dataset as linear combinations of the initial variables. Thus, non-linear (approximate) relations between the variables cannot be properly captured. In this case, a single relation could be described by several PCs, one describing the best linear approximation, and additional PCs describing directions in which the non-linear relation deviates from linearity. Using our non-linear transform, equivalent to a logarithm for high S/N values, alleviates part of the problem because it allows us to describe power-law relations. However, other non-linearities (such as saturation for the ^{12}CO line) are still not captured. Non-linear extensions of PCA exist (kernel-PCA, neural network-based dimensionality reduction such as the self-organizing maps), but their results tend to be harder to interpret.

The second property is that PCA is based on the global correlation matrix of the data. If different physical regimes are present in the dataset, each with different relations between the variables, PCA will provide a global (linear) approximation including the different regimes. Depending on the fraction of samples (pixels) representing the different regimes, it may give more weight to some regimes than others, neglect some regimes, or mainly represent one of the regimes.

6.3. Reduction of dimensionality

PCA is often used as a dimensionality-reduction tool, by keeping only a subset of the PCs that account for a sufficiently large fraction of the variance in the dataset. We saw that PC1 to 6 explain more than 90% of the correlation structure of the data. Moreover, PC6 defines a transition between several PCs with similar levels, and others with different levels (5% for PC3–4–5, 1–2% for the PCs after 6). The projection on the first six PC thus define a six-dimensional hyperplane in which the data is approximately embedded. However, we saw some striking spatial features appearing in later PCs, indicating meaningful axes of variations, such as PCs 8 and 11. The pattern of small scale spots shown by PC11 (corresponding to overbright CS) is particularly interesting, and its late apparition in the decomposition could simply be a consequence of the small fraction of pixels concerned. Thus, even PCs with lower fractions of explained correlation can contain important information, such as specific variations occurring in small regions only.

6.4. A synthetic view of Orion B

Using the physical interpretation of the principal components derived from the previous section it is possible to derive a synthetic view of the Orion B cloud rendered through a color image (see Fig. 11). The principal components 1 (column density), 2 (volume density) and 3 (UV illumination) are used in the following way: the column density is used to encode the luminosity, and

the volume density and UV illumination are combined orthogonally to define a color,

$$\text{hue} = \text{atan}(uv, \text{density}). \quad (2)$$

In this way, it is possible to identify, by color only, the physical properties associated with every line of sight.

Most of the region is composed of low density gas either obscured (green) or UV illuminated (yellow). Notable features are the moderately dense (orange) photodissociation regions that are present as the surface of pillars (e.g. around the Horsehead nebula) and as globules surrounding the NGC 2024 massive H II region in the upper part of the map. A sharp illumination gradient is visible at the base of the neck of the Horse with transition from illuminated (yellow) to shielded (green) gas.

Concerning the dust lane in front of NGC 2024, there is a clear sharp frontier between the northern and southern part, the north being strongly UV illuminated (yellow, orange, and red), the southern part much more obscured (cyan and green). The variation in the density of dense cores are visible with transitions from moderately dense (cyan) to higher density (dark blue) gas.

7. Conclusion

To study the correlations between maps of the emission of 12 bright lines belonging to the 3 mm band over the southwestern edge of the Orion B molecular cloud, we applied the PCA to these data. Before this analysis, we applied a non-linear transformation that is close to linear around zero and is equivalent to a logarithmic transform at large values. The goal of this non-linear transform is two-fold.

Firstly, although ratios of brightness temperatures are easier to interpret, PCA assumes that the relations in the input data set are linear. Applying the logarithm to the input data allows us to transform ratios of brightness temperatures into subtractions well adapted to a linear analysis.

Secondly, signal is only detected on a line-dependent subset of the field of view. Applying the logarithm to noisy brightness temperatures centered around zero is mathematically ill-defined. Having a linear transform around zero solves this problem. We tuned the transition value between the linear and logarithmic value that is typically eight times the typical noise value of the dataset. We showed that the results are not very sensitive to this value.

The PCA delivers a set of maps that are a linear combination of the input brightness temperatures, taking into account their (anti-)correlations. Although PCA does not use the spatial information of the input dataset, the output maps expose well-defined structures. We thus limited our analysis to the first few principal components that expose the largest correlations present in the initial dataset. The analysis of these correlations allowed us to propose links between the first three components and physical parameters, in this case the column density, volume density, and UV radiation field. We quantified these links by computing the correlation coefficients of these principal components with independent measurements of the column density, volume density, and UV illumination. The first principal component is highly correlated to the column density measured from the dust extinction and has positive contributions from all molecules, as has been noted in Pety et al. (2016). The third principal component is well correlated to our estimation of the UV illumination, with positive contributions from CCH, CN and anticorrelations with N_2H^+ and CH_3OH . The second principal component is correlated with the volume density in the dense cores having

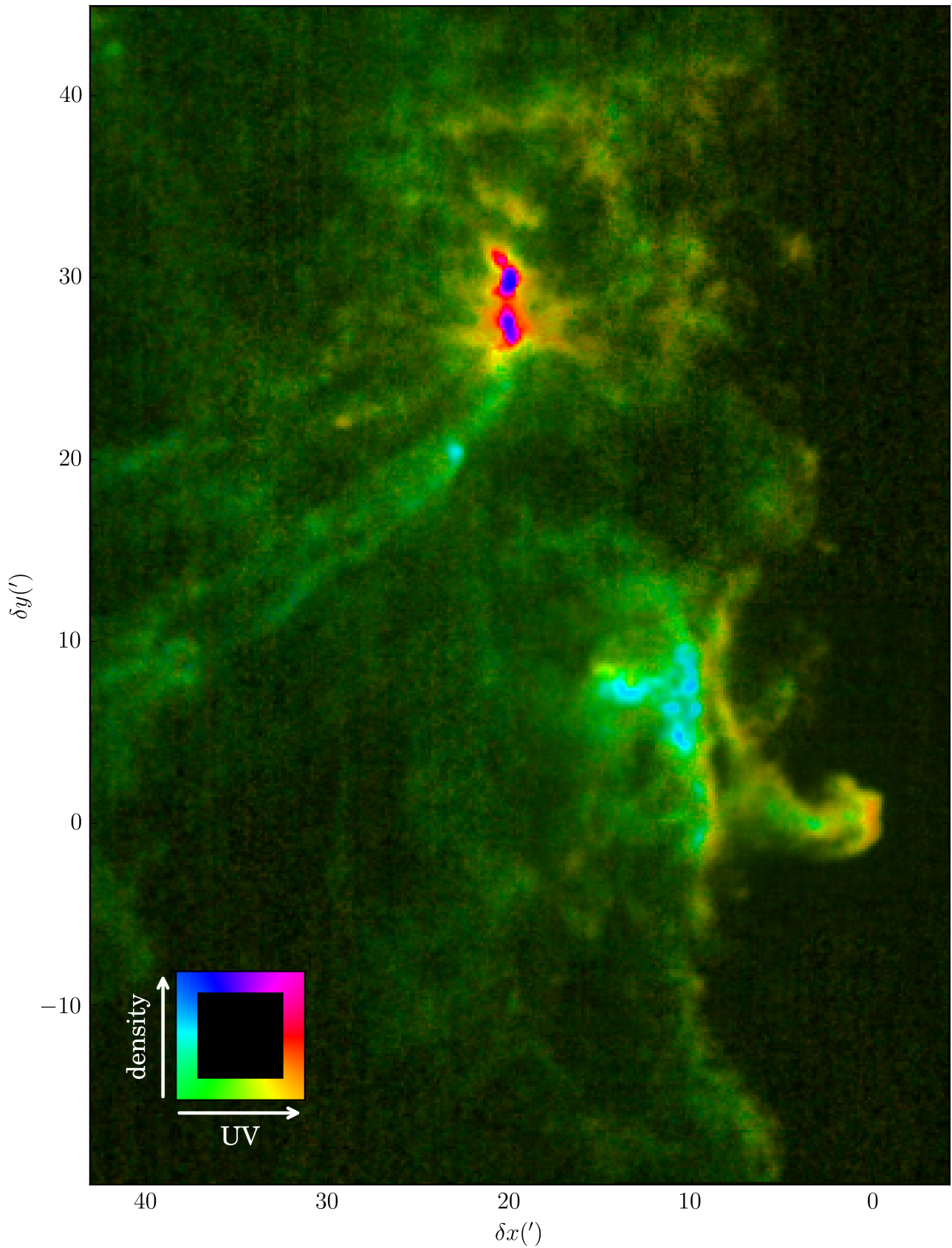


Fig. 11. Synthetic view of the Orion B molecular cloud. In this colormap, the intensity of each pixel is encoded by PC1 (column density) and the hue is encoded by the angle of the vector constructed using two orthogonal components PC2 (volume density) and PC3 (UV radiation field). It is possible to identify limiting cases. Magenta: dense PDR, yellow: diffuse PDR, green: diffuse non illuminated, blue: dense non illuminated.

a combined positive contribution from N_2H^+ and CH_3OH and a negative contribution from ^{12}CO and ^{13}CO .

The possibility of linking linear combinations of the brightness temperatures of a set of 3 mm lines to physical parameters such as the column density, volume density, or UV illumination opens an interesting avenue to analyze the large spectro-imaging data sets that (sub)-mm radioastronomy starts to produce. As PCA analysis only works on the brightness temperatures independent of their spatial relations, it also offers an easy possibility to compare with large grids of detailed 1D models of photo-dissociation regions. In future papers, we will continue to explore this with more advanced decomposition techniques that may take into account missing values, noise effects, or non-linear relations in the input dataset.

Acknowledgements. This work was supported by the French program Physique et Chimie du Milieu Interstellaire (PCMI) funded by the Conseil National de la Recherche Scientifique (CNRS) and Centre National d'Études Spatiales (CNES). We thank the CIAS for its hospitality during the two workshops devoted to this project. P.G. thanks ERC starting grant (3DICE, grant agreement 336474) for funding during this work. P.G.'s current postdoctoral position is funded by the INSU/CNRS. NRAO is operated by Associated Universities Inc. under contract with the National Science Foundation. We thank P. Andre and N. Schneider to kindly give us access to the *Herschel* Gould Belt Survey data. This research has made use of data from the *Herschel* Gould Belt survey (HGBS) project (<http://gouldbelt-herschel.cea.fr>). The HGBS is a *Herschel* Key Programme jointly carried out by SPIRE Specialist Astronomy Group 3 (SAG 3), scientists of several institutes in the PACS Consortium (CEA Saclay, INAF-IFSI Rome and INAF-Arcetri, KU Leuven, MPIA Heidelberg), and scientists of the *Herschel* Science Center (HSC). We than the anonymous referee for his/her constructive comments.

References

- Agúndez, M., & Wakelam, V. 2013, *Chem. Rev.*, **113**, 8710
 André, P., Men'shchikov, A., Bontemps, S., et al. 2010, *A&A*, **518**, L102
 Brunt, C. M., & Heyer, M. H. 2013, *MNRAS*, **433**, 117
 Deeming, T. J. 1964, *MNRAS*, **127**, 493
 Draine, B. T., & Li, A. 2007, *ApJ*, **657**, 810
 Feigelson, E. D., & Babu, G. J. 2012, *Modern Statistical Methods for Astronomy* (Cambridge: Cambridge University Press)
 Heyer, M. H., & Peter Schloerb, F. 1997, *ApJ*, **475**, 173
 Hotelling, H. 1933, *J. Educ. Psych.*, **24**
 Ilin, A., & Raiko, T. 2010, *J. Mach. Learn. Res.*, **11**, 1957
 Ivezić, Ž., Connolly, A., Vanderplas, J., & Gray, A. 2014, *Statistics, Data Mining and Machine Learning in Astronomy* (Princeton University Press)
 Jolliffe, I. 2002, *Principal Component Analysis* (New York: Springer Verlag)
 Jones, P. A., Burton, M. G., Cunningham, M. R., et al. 2012, *MNRAS*, **419**, 2961
 Kirk, H., Di Francesco, J., Johnstone, D., et al. 2016, *ApJ*, **817**, 167
 Le Bourlot, J., Le Petit, F., Pinto, C., Roueff, E., & Roy, F. 2012, *A&A*, **541**, A76
 Lo, N., Cunningham, M. R., Jones, P. A., et al. 2009, *MNRAS*, **395**, 1021
 Meisner, A. M., & Finkbeiner, D. P. 2014, *ApJ*, **781**, 5
 Melnick, G. J., Tolls, V., Snell, R. L., et al. 2011, *ApJ*, **727**, 13
 Menten, K. M., Reid, M. J., Forbrich, J., & Brunthaler, A. 2007, *A&A*, **474**, 515
 Neufeld, D. A., Hollenbach, D. J., Kaufman, M. J., et al. 2007, *ApJ*, **664**, 890
 Neufeld, D. A., Godard, B., Gerin, M., et al. 2015, *A&A*, **577**, A49
 Orkisz, J. H., Pety, J., Gerin, M., et al. 2017, *A&A*, **599**, A99
 Pedregosa, F., Varoquaux, G., Gramfort, A., et al. 2011, *J. Mach. Learn. Res.*, **12**, 2825
 Pety, J., Guzmán, V. V., Orkisz, J. H., et al. 2017, *A&A*, **599**, A98
 Roman-Duval, J., Federrath, C., Brunt, C., et al. 2011, *ApJ*, **740**, 120
 Rosolowsky, E., & Leroy, A. 2006, *PASP*, **118**, 590
 Scargle, J. D., Norris, J. P., Jackson, B., & Chiang, J. 2013, *ApJ*, **764**, 167
 Schneider, N., André, P., Könyves, V., et al. 2013, *ApJ*, **766**, L17
 Ungerechts, H., & Thaddeus, P. 1987, *ApJS*, **63**, 645
 Ungerechts, H., Bergin, E. A., Goldsmith, P. F., et al. 1997, *ApJ*, **482**, 245
 Vanderplas, J., Connolly, A., Ivezić, Ž., & Gray, A. 2012, in *Conference on Intelligent Data Understanding (CIDU)*, 47
 Ward-Thompson, D., Nutter, D., Bontemps, S., Whitworth, A., & Attwood, R. 2006, *MNRAS*, **369**, 1201

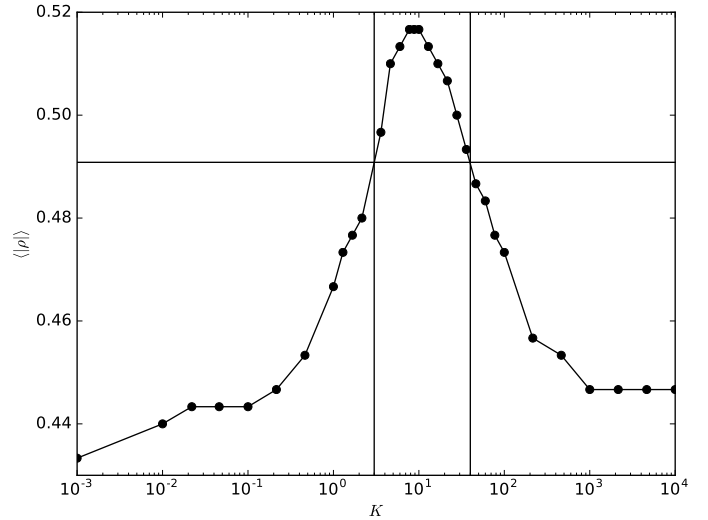


Fig. A.1. Variation of $\langle |\rho| \rangle$, the mean of the absolute values of the Spearman's correlation coefficients as a function of K , with $a = K \text{median}(\sigma)$. The optimal value is found for $K = 8$.

Appendix A: Optimal value of a in the asinh reparametrisation

The only free parameter in the asinh reparametrization is a , the parameter which marks the boundary between the linear and logarithmic regimes of the asinh function (see Fig. 3). As shown in Table 1 the noise across different lines is similar and we express a as the product of a constant factor K by the median noise 0.08 K. The quantity we choose to maximize is the mean of the absolute value of the correlation coefficient of the principal components with the physical maps $N(\text{H}_2)$, U , and n_{H} , we note this quantity $\langle |\rho| \rangle$. Figure A.1 shows the evolution of this quantity with increasing values of K , a maximum value of $\langle |\rho| \rangle$ around $K = 8$ although an acceptable range of K values (reduction of $\langle |\rho| \rangle$ by less than 5%) spans values from 3 to 40.