# Principal component analysis of spectral line data: analytic formulation

C. M. Brunt[1]★ and M. H. Heyer[2]

[1]*School of Physics, University of Exeter, Stocker Road, Exeter, UK*

[2]*Department of Astronomy, University of Massachusetts, Amherst, MA 01003, USA*

**ABSTRACT**

Principal component analysis is a powerful statistical system to investigate the structure and dynamics of the molecular interstellar medium, with particular emphasis on the study of turbulence, as revealed by spectroscopic imaging of molecular line emission. To date, the method to retrieve the power-law index of the velocity structure function or power spectrum has relied on an empirical calibration and testing with model turbulent velocity fields, while lacking a firm theoretical basis. In this paper, we present an analytic formulation that reveals the detailed mechanics of the method and confirms previous empirical calibrations of its recovery of the scale dependence of turbulent velocity fluctuations.

**Key words:** turbulence – methods: statistical – ISM: clouds – ISM: kinematics and dynamics.

## 1 INTRODUCTION

Wide field, spectroscopic imaging of molecular line emission provides a vast amount of information of the gas dynamics of interstellar clouds. To exploit this information, Heyer & Schloerb (1997, hereafter HS97) introduced the application of principal component analysis (PCA) to the position–position–velocity data cubes as a tool to investigate the structure and dynamics of molecular clouds. Brunt & Heyer (2002a, hereafter BH02) more rigorously defined HS97's method for quantifying the scale-dependence of turbulent velocity fluctuations in molecular clouds, and HS97's PCA formulation has since undergone a number of extensions and refinements (Brunt 2003a; Brunt et al. 2003; Heyer et al. 2008; Roman-Duval et al. 2011). The HS97/BH02 PCA formulation is ideally suited to the analysis of low signal-to-noise ratio data and for this reason has been most commonly applied to wide-field survey data (Brunt & Heyer 2002b; Heyer & Brunt 2004, 2012; Roman-Duval et al. 2011).

A significant limitation of the HS97/BH02 method to derive the power-law index of the velocity structure function is its reliance on an empirical calibration that establishes the relationship between the index determined from PCA and the true index of the models generated by numerical representations and computational simulations of turbulent clouds (BH02; Brunt et al. 2003; Roman-Duval et al. 2011). Therefore, the data analysis has lacked a firm theoretical underpinning upon which other statistical methods are based (Scalo 1984; Kleiner & Dickman 1985; Miesch & Bally 1994; Stutzki et al. 1998; Lazarian & Pogosyan 2000).

In this paper, we present an analytic formulation of the PCA method that validates these previous empirical calibrations. This

is a challenging task as it requires analytical representations of a complex physical process (turbulence) as measured by a complex analysis method (PCA). To simplify the problem, the formulation relies on a central assumption that the spectral line profiles in a spectral line imaging observation of a molecular cloud can be represented as an ensemble of Gaussians of fixed dispersion, with turbulent spatial correlations. The formulation predicts covariance matrices, eigenvectors, eigenvalues and eigenimage structure and enables insight into the mechanics of the PCA method that explains several empirically observed features noted in the literature (Brunt et al. 2003; Roman-Duval et al. 2011).

The layout of the paper is as follows. In Section 2 we provide a brief summary of the HS97 formulation. In Section 3 we derive covariance matrices expected from an ensemble of Gaussian line profiles with variable centroids. Sections 4–6, respectively, describe the derivation of the resultant eigenvectors, eigenvalues and eigenimages. In Section 7, we present an analytic derivation of BH02's calibration of the PCA method for the turbulent velocity fluctuation spectrum. A summary is given in Section 8.

## 2 PRINCIPAL COMPONENT ANALYSIS

In this section, we review the HS97 formulation of PCA applied to spectral line imaging observations, and summarize the key empirical findings that an analytic formulation should aim to explain.

### 2.1 The HS97 PCA formulation

A spectroscopic imaging observation comprises an ensemble of $n$ spectra each with $p$ spectroscopic channels. We write the data cube as $T(\boldsymbol{r}_i, v_j) = T_{ij}$, where $\boldsymbol{r}_i$ denotes the spatial coordinate of the $i$th spectrum.

In the formulation of HS97, the spectrum, or line profile, at each spatial grid point is taken to be the raw measurable quantity that will be subjected to PCA. From the ensemble of line profiles, the covariance matrix elements $S_{jk}$ are calculated as

$$S_{jk} = S(v_j, v_k) = \frac{1}{n} \sum_{i=1}^{n} T_{ij} T_{ik}. \qquad (1)$$

A set of eigenvectors, $u_{mj} = u_m(v_j)$, and eigenvalues, $\lambda_m$, is determined from the solution of the eigenvalue equation for the covariance matrix,

$$S_{jk} u_{mj} = \lambda_m u_{mj}. \qquad (2)$$

The eigenvalue, $\lambda_m$, equals the amount of variance projected on to its corresponding eigenvector, $u_{mj}$.

The eigenimages, $I_m(\mathbf{r}_i)$, are constructed from the projected values of the data, $T_{ij}$, on to the eigenvectors, $u_{mj}$,

$$I_m(\mathbf{r}_i) = \sum_{j=1}^{p} T_{ij} u_{mj}. \qquad (3)$$

We refer to the coupled eigenvector and eigenimage at order $m$ as the $m$th principal component (PC). In the most basic interpretation, the set of eigenvectors describe the velocity magnitude of line profile differences with the position-position-velocity (ppv) volume, as these generate varying levels of variance. Such differences arise from gas motions such as infall, outflow, rotation, turbulent velocity fluctuations and, of course, random noise of the observation. The eigenimages show where these profile differences occur within the projected position–position plane.

## 2.2 Empirical results

In their foundational work, HS97 suggested that, at each order $m$, the coupled eigenvector (as a velocity function) and eigenimage (as a spatial function) could be used to study the scale-dependence of velocity fluctuations in molecular clouds. Specifically, defining $\delta v_m$ and $\delta l_m$ as the characteristic widths of the eigenvector and eigenimage autocorrelation functions (ACFs), respectively, HS97 found power-law relations ($\delta v_m \propto \delta l_m^\alpha$) for a sample of molecular clouds subjected to PCA.

HS97's proposed method $\delta v_m$ was scrutinized by BH02, who included accounting for noise and finite resolution, and fixed $\delta v_m$ and $\delta l_m$ as the $1/e$ points of the eigenvector and eigenimage ACFs, respectively. BH02 also investigated the method's ability to recover intrinsic three-dimensional (3D) statistical information about the velocity field and established the first calibration of the method: $\alpha \approx 0.33\beta$ where $\beta$ is the spectral slope of the angular integral of the velocity power spectrum in 3D (in this representation, a Kolmogorov spectrum has $\beta = 5/3$ and a shock-dominated spectrum has $\beta = 2$). Roman-Duval et al. (2011) confirmed the BH02 calibration and examined in detail the sensitivity of the calibration to density fluctuations, using lognormal density PDFs, concluding that the calibration was stable below a critical level of (very high) density variability [$\sigma_{\ln(\rho/\rho_0)} > 2$]. Brunt et al. (2003a) and Roman-Duval et al. (2011) showed that the method is sensitive to first-order velocity fluctuations, rather than root-mean-square velocity fluctuations.

## 3 COVARIANCE MATRICES

Our analysis begins with a basic investigation of the covariance matrices that result from an ensemble of Gaussian line profiles of fixed dispersion. We initially examine the case of a single component per line of sight, and then consider the more complex case of multiple Gaussians. This analysis forms the basis of later derivations in the subsequent sections.

### 3.1 Single Gaussian component case

We first consider the covariance matrix that would be derived from an ensemble of Gaussian line profiles. Let all line profiles have the same dispersion, $\sigma_b^2$, and let the distribution of centroid velocities be drawn from a Gaussian distribution of dispersion $\sigma_c^2$ around a global mean velocity of zero. The total velocity dispersion of this ensemble is $\sigma_{\mathrm{tot}}^2 = \sigma_b^2 + \sigma_c^2$. Note that here, the subscript $b$ refers generically to 'broadening' of the line profile due to macroscopic turbulent fluctuations along the line of sight, and not just to the (typically much narrower) thermal broadening. The use of a single dispersion $\sigma_b^2$ to represent this is a simplification, as not all lines of sight will produce exactly the same broadening, though observationally linewidths do not vary significantly across a cloud.

The terms representing the $i$th spectrum, $T_{ij}$ and $T_{ik}$, in the covariance matrix equation are written as:

$$T_{ij} = T_i(v_j) = T_{0i} \exp\left(-\frac{(v_j - v_{ci})^2}{2\sigma_b^2}\right), \qquad (4)$$

$$T_{ik} = T_i(v_k) = T_{0i} \exp\left(-\frac{(v_k - v_{ci})^2}{2\sigma_b^2}\right), \qquad (5)$$

where $T_{0i}$ is the peak temperature, $v_{ci}$ is the centroid velocity and $\sigma_b^2$ is the velocity dispersion of the $i$th line profile.

For the above model, the covariance matrix equation is

$$S_{jk} = \frac{1}{n} \sum_{i=1}^{n} T_{0i}^2 \exp\left(-\frac{(v_j - v_{ci})^2}{2\sigma_b^2}\right) \exp\left(-\frac{(v_k - v_{ci})^2}{2\sigma_b^2}\right), \qquad (6)$$

where the summation is over the total number of line profiles, $n$. For large enough $n$ we can convert the normalized summation over $i$ to integrals over the probability distributions of peak temperature, $T_0$, and centroid velocity, $v_c$, to write:

$$S_{jk} = \int_0^\infty \mathrm{d}T_0 \int_{-\infty}^\infty \mathrm{d}v_c P_T(T_0) P_v(v_c) T_0^2$$
$$\times \exp\left(-\frac{(v_j - v_c)^2}{2\sigma_b^2}\right) \exp\left(-\frac{(v_k - v_c)^2}{2\sigma_b^2}\right), \qquad (7)$$

where we have assumed that $T_0$ and $v_c$ are uncorrelated, with independent probability distributions, $P_T(T_0)$ and $P_v(v_c)$, respectively. Assuming a Gaussian probability distribution for $v_c$, with dispersion $\sigma_c^2$, the integrals are easily solved to yield:

$$S_{jk} = S_0 \exp\left(-\frac{(v_j^2 + v_k^2)}{2\sigma_b^2} + \frac{(v_j + v_k)^2}{4\sigma_b^2(1 + \sigma_b^2/2\sigma_c^2)}\right), \qquad (8)$$

where

$$S_0 = \frac{\langle T_0^2 \rangle}{\sqrt{1 + 2\sigma_c^2/\sigma_b^2}}. \qquad (9)$$

Equation (8) is valid for ensembles where the peak temperature of the lines can vary with position, provided the peak temperatures are uncorrelated with the centroid velocities. Note that the contribution of a line profile to $S_{jk}$ is proportional to $T_0^2$. For consistency, this
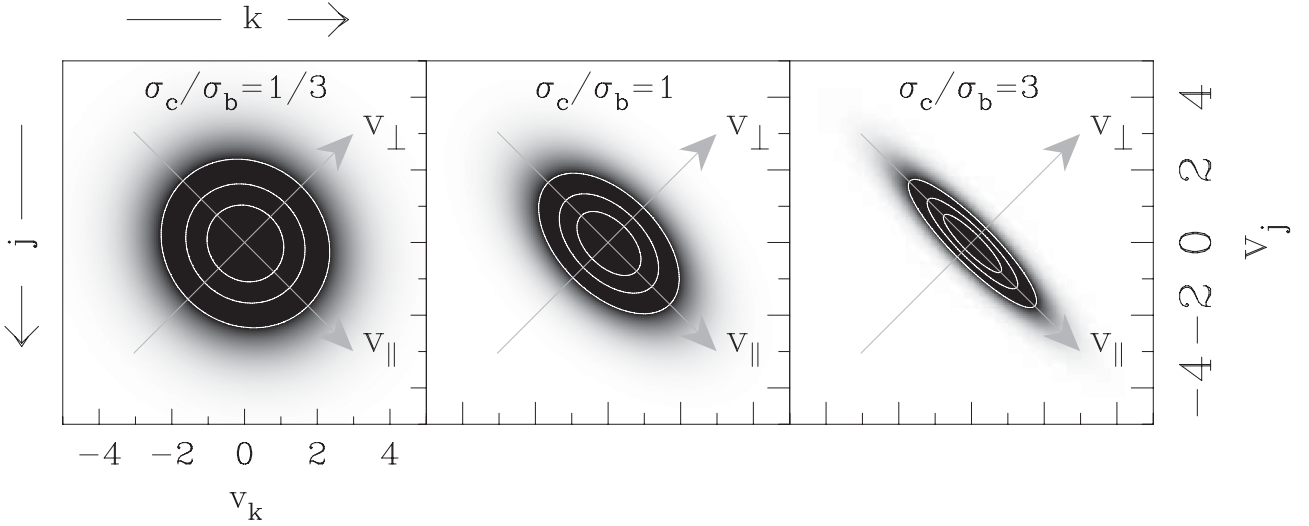
**Figure 1.** Gray-scale representations of the covariance matrix, $S$, obtained with varying $\sigma_c/\sigma_b$. The variance along the diagonal, $\sigma_{||}^2 = 2\sigma_c^2 + \sigma_b^2 = 19/9$, is the same for all plots. Contours are shown at 25 per cent, 50 per cent and 75 per cent of the peak of $S_{jk}$.

requires that $\sigma_c^2$ be defined by:

$$\sigma_c^2 = \frac{\sum\limits_{i=1}^{n} T_{0i}^2 v_{ci}^2}{\sum\limits_{i=1}^{n} T_{0i}^2} = \frac{\sum\limits_{i=1}^{n} W_{0i}^2 v_{ci}^2}{\sum\limits_{i=1}^{n} W_{0i}^2}, \tag{10}$$

where $W_{0i} = \sqrt{2\pi}T_{0i}\sigma_b$ is the integrated intensity of the $i$th line profile in the above model. Ideally, equation (7) would include a probability distribution of $\sigma_b^2$, but the simplification of a constant $\sigma_b^2$ was necessary to make the integration tractable.

To visualize equation (8) we constructed covariance matrices accordingly for varying $\sigma_b$ and $\sigma_c$. (These matrices agree with numerical realizations.) Fig. 1 shows three example covariance matrices, represented as grey-scale images. In general, the matrices will vary from a fully linearly dependent case ($\sigma_c/\sigma_b \longrightarrow 0$) to a fully diagonal case ($\sigma_b/\sigma_c \longrightarrow 0$). In a fully diagonal matrix, each row (column) is linearly independent.

We compute the dispersion of $S$ along the diagonal, $\sigma_{||}^2$, using:

$$S_{||} = S_0 \exp\left(-\frac{2v_j^2}{2\sigma_b^2} + \frac{(2v_j)^2}{4\sigma_b^2(1 + \sigma_b^2/2\sigma_c^2)}\right) = S_0 \exp\left(-\frac{2v_j^2}{2\sigma_{||}^2}\right), \tag{11}$$

obtained by setting $v_k = v_j$ in equation (4), and noting that the distance along the diagonal is $\sqrt{2}v_j$, to find

$$\sigma_{||}^2 = 2\sigma_c^2 + \sigma_b^2. \tag{12}$$

Similarly, we compute the dispersion of $S$ perpendicular to the diagonal, $\sigma_\perp^2$, using:

$$S_\perp = S_0 \exp\left(-\frac{2v_j^2}{2\sigma_b^2}\right) = S_0 \exp\left(-\frac{2v_j^2}{2\sigma_\perp^2}\right), \tag{13}$$

obtained by setting $v_k = -v_j$ in equation (4), and noting that the distance along the perpendicular is $\sqrt{2}v_j$, to find:

$$\sigma_\perp^2 = \sigma_b^2. \tag{14}$$

More generally, defining:

$$v_{||} = \frac{1}{\sqrt{2}}(v_k + v_j), \tag{15}$$

$$v_\perp = \frac{1}{\sqrt{2}}(v_k - v_j), \tag{16}$$

it is straightforward to show that:

$$S(v_{||}, v_\perp) = S_0 \exp\left(-\frac{v_\perp^2}{2\sigma_\perp^2}\right) \exp\left(-\frac{v_{||}^2}{2\sigma_{||}^2}\right), \tag{17}$$

i.e. that the covariance matrix is an elliptical Gaussian, with dispersions $\sigma_{||}^2$ and $\sigma_\perp^2$ parallel and perpendicular to the diagonal, respectively.

By fitting an elliptical Gaussian to the covariance matrix, $\sigma_{||}^2$ and $\sigma_\perp^2$ can be measured, and we can deduce the line centroid dispersion, $\sigma_c^2$, and profile dispersion, $\sigma_b^2$, via:

$$\sigma_c^2 = \frac{1}{2}(\sigma_{||}^2 - \sigma_\perp^2) \tag{18}$$

$$\sigma_b^2 = \sigma_\perp^2. \tag{19}$$

It is worth noting also that the total velocity dispersion, $\sigma_{tot}^2$, is given by

$$\sigma_{tot}^2 = \sigma_c^2 + \sigma_b^2 = \frac{1}{2}(\sigma_{||}^2 + \sigma_\perp^2). \tag{20}$$

### 3.2 Multiple Gaussian component case

We now consider a more elaborate model in which the $i$th spectrum is represented by the summation of $n_t$ spectral lines, each of dispersion $\sigma_t^2$, where we take $n_t$ to be moderately large. Let the centroid velocities of each of these components be drawn from a Gaussian probability distribution of dispersion $\sigma_b^2 - \sigma_t^2$ centred on $v_{ci}$. Here we envisage the individual narrow lines to have approximately thermal linewidths (dispersion $\sigma_t^2$) that collectively generate a broadened line profile (with dispersion $\sigma_b^2$) due to macroscopic velocity differences along the line of sight. In the limit of large $n_t$,

the single component model of the preceding section (i.e. a single Gaussian line of dispersion $\sigma_b^2$ and centroid $v_{ci}$) will be recovered. For moderate $n_t$, the line profiles could appear asymmetric and/or multiply peaked, but many profiles averaged together would appear Gaussian. The contribution of the $i$th spectrum to the covariance matrix is in this case:

$$\left[\sum_{e=1}^{n_t} T_{0ie} \exp\left(-\frac{(v_j - v_{cie})^2}{2\sigma_t^2}\right)\right]$$

$$\times \left[\sum_{f=1}^{n_t} T_{0if} \exp\left(-\frac{(v_k - v_{cif})^2}{2\sigma_t^2}\right)\right]. \tag{21}$$

The contributions for $e = f$:

$$\sum_{e=1}^{n_t} T_{0ie}^2 \exp\left(-\frac{(v_j - v_{cie})^2}{2\sigma_t^2}\right) \exp\left(-\frac{(v_k - v_{cie})^2}{2\sigma_t^2}\right) \tag{22}$$

averaged over all positions $i$, produce an overall contribution to $S_{jk}$ proportional to:

$$\exp\left(-\frac{(v_j^2 + v_k^2)}{2\sigma_t^2} + \frac{(v_j + v_k)^2}{4\sigma_t^2(1 + \sigma_t^2/2(\sigma_c^2 + \sigma_b^2 - \sigma_t^2))}\right) \tag{23}$$

(cf. equation 8). The contribution of the cross-terms ($e \neq f$) is more difficult to deal with, but we note that their contribution should recover the form of equation (8) in the limit of large $n_t$. Therefore, we write the approximate form of the covariance matrix in the multiple component case as:

$$S_{jk} \approx S_0 \eta \exp\left(-\frac{(v_j^2 + v_k^2)}{2\sigma_t^2} + \frac{(v_j + v_k)^2}{4\sigma_t^2(1 + \sigma_t^2/2(\sigma_c^2 + \sigma_b^2 - \sigma_t^2))}\right)$$

$$+ S_0(1 - \eta) \exp\left(-\frac{(v_j^2 + v_k^2)}{2\sigma_b^2} + \frac{(v_j + v_k)^2}{4\sigma_b^2(1 + \sigma_b^2/2\sigma_c^2)}\right), \tag{24}$$

where we expect $\eta \to 0$ as $n_t \to \infty$. This covariance matrix form contains an additional (small) contribution from resolvable fine structure in the line profiles, with dispersion along the diagonal of $2(\sigma_c^2 + \sigma_b^2) - \sigma_t^2$ and dispersion perpendicular to the diagonal of $\sigma_t^2$. Qualitatively, this is a weak, strongly diagonal feature in the covariance matrix, though this result is obtained only in the large $n_t$ limit.

## 4 EIGENVECTORS

In this section, we first derive the eigenvectors that result from a covariance matrix of the form given by equation (8). Next, we derive the ACFs of the eigenvectors and determine the autocorrelation scale, $\delta v_m$ (i.e. the velocity-lag of the $1/e$-point of the normalized ACF) as a function of order $m$. This is a key observable in the application of PCA to determine the turbulent energy spectrum (HS97; BH02).

### 4.1 Eigenvector structure

A valid solution of the eigenvalue equation (2) requires that:

$$\int_{-\infty}^{+\infty} dv_k \, S(v_j, v_k) u(v_k) = \lambda u(v_j), \tag{25}$$

where $u(v_k)$ is an eigenvector, $\lambda$ is its eigenvalue and we have approximated the finite sums as integrals. We now search for a valid solution of equation (25), using the form of equation (8), by setting:

$$u(v_k) = I_0 \exp(-cv_k^2), \tag{26}$$

where $I_0$ and $c$ are constants. We use the single component covariance matrix given by equation (8); an analytic solution for the multiple component case (equation 24) has not yet been found.

The terms in the exponent of equation (8) may be written as:

$$-(av_j^2 + av_k^2 - 2bv_jv_k), \tag{27}$$

where

$$a = \frac{1}{2\sigma_b^2} - \frac{1}{4\sigma_b^2(1 + \sigma_b^2/2\sigma_c^2)}, \tag{28}$$

and

$$b = \frac{1}{4\sigma_b^2(1 + \sigma_b^2/2\sigma_c^2)}. \tag{29}$$

The exponent of the integrand in equation (25) is then:

$$-(av_j^2 + av_k^2 - 2bv_jv_k + cv_k^2), \tag{30}$$

which may be regrouped as

$$-\left[\left((a+c)^{1/2}v_k - \frac{b}{(a+c)^{1/2}}v_j\right)^2 + \left(a - \frac{b^2}{(a+c)}\right)v_j^2\right]. \tag{31}$$

With a change of variable:

$$w = (a+c)^{1/2}v_k - \frac{b}{(a+c)^{1/2}}v_j, \tag{32}$$

we find that equation (25) is satisfied if:

$$a - \frac{b^2}{(a+c)} = c, \tag{33}$$

or

$$c^2 = a^2 - b^2 = \frac{1}{4}\left(\sigma_b^2(2\sigma_c^2 + \sigma_b^2)\right)^{-1/2}. \tag{34}$$

We identify the solution (equation 26) as the first eigenvector ($u_{1j} = u_1(v_j)$), and demonstrate the validity of this choice below. For simplicity, we write the solution as

$$u_1(v_j) = I_{01} \exp\left(-\frac{v_j^2}{2\sigma_1^2}\right), \tag{35}$$

where $I_{01}$ is a constant, and:

$$\sigma_1 = \sqrt{1/2c} = (\sigma_b^2(2\sigma_c^2 + \sigma_b^2))^{1/4} = \sigma_{||}^{1/2}\sigma_{\perp}^{1/2}. \tag{36}$$

To deduce the forms of the higher order eigenvectors, we make use of the orthogonality condition:

$$\int_{-\infty}^{+\infty} dv_j u_m(v_j) u_n(v_j) = I_{0m} I_{0n} \delta_{mn}, \tag{37}$$

where $I_{0m}$ and $I_{0n}$ are constants which depend on the choice of normalization of the eigenvectors, and $\delta_{mn}$ is the Kronecker delta ($\delta_{mn} = 1$ if $m = n$, and $\delta_{mn} = 0$ if $m \neq n$).

The set of functions that are orthogonal with respect to a Gaussian weight are the Hermite polynomials. The orthogonality condition for Hermite polynomials is

$$\int_{-\infty}^{\infty} dx \, H_n(x) H_m(x) \exp(-x^2) = \delta_{mn} 2^n n! \sqrt{\pi}. \tag{38}$$

Comparing equations (37) and (38), we identify the $m$th order eigenvector as the product of the first eigenvector and the $(m-1)$th order Hermite polynomial, $H_{m-1}(v_j/\sigma_1)$. Thus the $m$th order eigenvector
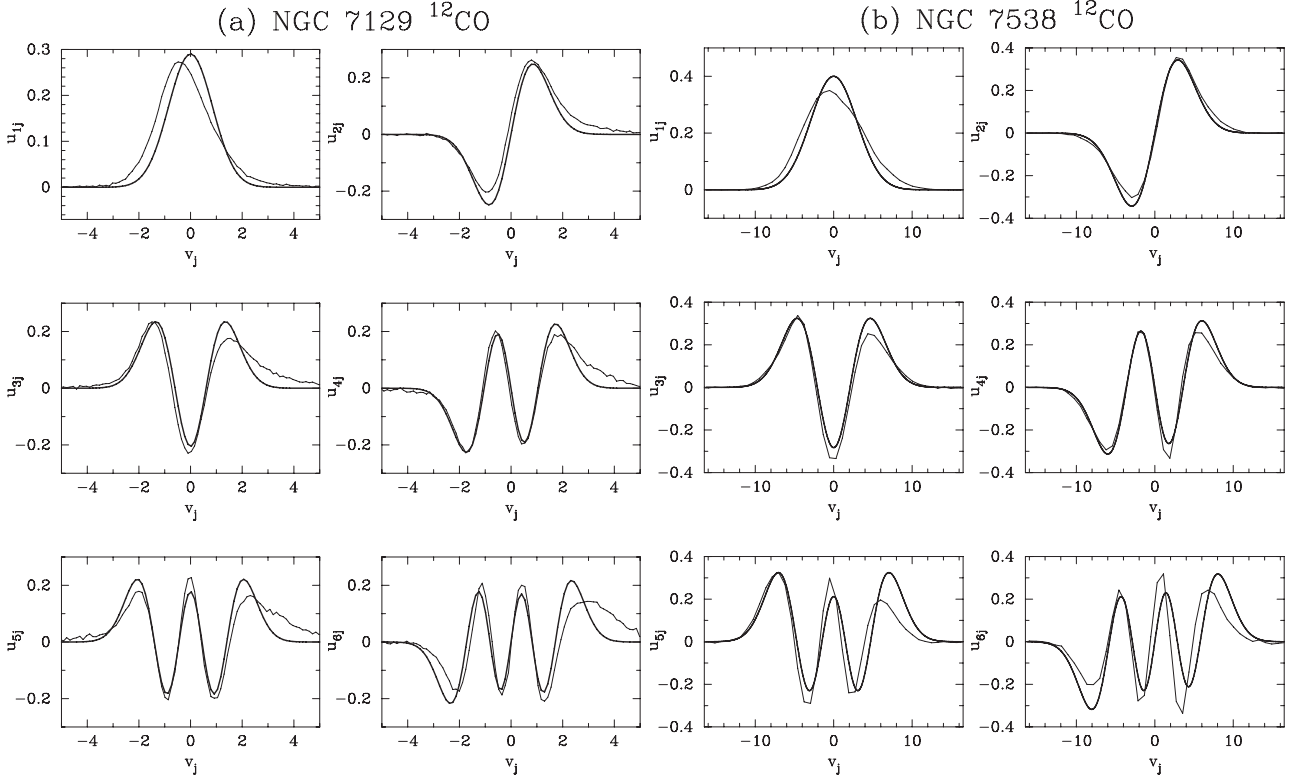
## (a) NGC 7129 $^{12}$CO    (b) NGC 7538 $^{12}$CO



**Figure 2.** Eigenvectors, $u_{nj} = u_n(v_j)$, obtained from (a) NGC 7129 $^{12}$CO and (b) NGC 7538 $^{12}$CO (lighter lines). Eyeball fits to the fourth eigenvector, $u_4$, have been made using the form given in equation (39). The heavy lines are those predicted by equation (39) with $\sigma_1$ and $I_{01}$ specified.

has the form:

$$u_{mj} = \frac{I_{01}}{\sqrt{2^{m-1}(m-1)!}} \exp\left(-\frac{v_j^2}{2\sigma_1^2}\right) H_{m-1}\left(\frac{v_j}{\sigma_1}\right), \qquad (39)$$

where $I_{01}$ is the peak amplitude of the first eigenvector.

The eigenvectors defined by equation (39) provide a reasonably good representation of eigenvectors obtained from spectral line imaging observations of CO isotopes in molecular clouds. Figs 2(a) and (b) show the first six eigenvectors obtained from PCA of $^{12}$CO emission in the NGC 7129 molecular cloud (Brunt & Mac Low 2004) and the NGC 7538 giant molecular cloud (Heyer et al. 1998), respectively. We have fitted (by eye) the fourth eigenvectors with $u_4$ from equation (39) and constructed the other eigenvectors according to $\sigma_1$ and $I_{01}$ obtained from the fit of $u_4$. The point here is not to evaluate the detailed applicability of equation (39) to real observations, which contain more sources of line profile variance than accounted for by our simple model. Line profile asymmetries, multiplicities and other non-Gaussian features will be represented in the covariance matrix and in turn will affect the detailed structure of the eigenvectors. Fig. 2 is presented to demonstrate that observed eigenvectors at order $m$ can be interpreted as the product of a $\sim$Gaussian and a polynomial of order $m-1$.

### 4.2 Eigenvector autocorrelation functions and characteristic velocity scales

The unnormalized ACF, $C_m(v)$, of the $m$th eigenvector is

$$C_m(v) = \int_{-\infty}^{\infty} dv' \, u_m(v') u_m(v' - v), \qquad (40)$$

where $u_m(v)$ at order $m$ is given by equation (39). Writing $x = v/\sigma_1$ and $y = v'/\sigma_1$, this is then

$$C_m(x) = C_m(v/\sigma_1)$$

$$= \frac{I_{01}^2}{2^{m-1}(m-1)!} \int_{-\infty}^{\infty} dy \, \exp\left(-y^2/2\right)$$

$$\times \exp\left(-(x-y)^2/2\right) H_{m-1}(y) H_{m-1}(y-x). \qquad (41)$$

We make the substitution $w = y - x/2$ to find

$$C_m(x) = C_m(v/\sigma_1)$$

$$= \frac{I_{01}^2}{2^{m-1}(m-1)!} \int_{-\infty}^{\infty} dw \, \exp(-w^2)$$

$$\times \exp(-(x/2)^2) H_{m-1}(w+x/2) H_{m-1}(w-x/2). \qquad (42)$$

The Hermite polynomial terms may be expanded as:

$$H_{m-1}(w+x/2) = \sum_{k=0}^{m-1} \frac{(m-1)!}{k!(m-1-k)!} H_k(w) x^{m-1-k},$$

$$H_{m-1}(w-x/2) = \sum_{k=0}^{m-1} \frac{(m-1)!}{k!(m-1-k)!} H_k(w)(-x)^{m-1-k}. \qquad (43)$$

Using the orthogonality of Hermite polynomials (equation 38), this then gives

$$\frac{C_m(x)}{C_m(0)} = \frac{C_m(v/\sigma_1)}{C_m(0)} = \exp\left(-(x/2)^2\right) B_{m-1}(x), \qquad (44)$$

where

$B_{m-1}(x)$

$$= \sum_{k=0}^{m-1} \frac{2^{-(m-1-k)}}{(m-1-k)!} \frac{(m-1)!}{k!(m-1-k)!}(-1)^{m-1-k}x^{2(m-1-k)}. \quad (45)$$

Note that we have also written these in normalized form.

The first five normalized ACFs are:

$$\frac{C_1(x)}{C_1(0)} = \exp\left(-(x/2)^2\right)$$

$$\frac{C_2(x)}{C_2(0)} = \exp\left(-(x/2)^2\right)(1 - x^2/2)$$

$$\frac{C_3(x)}{C_3(0)} = \exp\left(-(x/2)^2\right)(1 - x^2 + x^4/8)$$

$$\frac{C_4(x)}{C_4(0)} = \exp\left(-(x/2)^2\right)(1 - 3x^2/2 + 3x^4/8 - x^6/48)$$

$$\frac{C_5(x)}{C_5(0)} = \exp\left(-(x/2)^2\right)(1 - 2x^2 + 3x^4/4 - x^6/12 + x^8/384).$$

$$(46)$$

The first six ACFs are shown in Fig. 3 – cf. fig. 9 of HS97.

The velocity scale, $\delta v_m$, at order $m$ is given by the $1/e$-point of the normalized ACF, i.e. $C_m(\delta v_m/\sigma_1)/C_m(0) = 1/e$. While it is difficult to determine the $1/e$ points analytically, they may be determined numerically. Fig. 4 shows the measured $\delta v_m/\sigma_1$ values versus $m - 1$, which approximately obey a power-law relation:

$$\delta v_m/\sigma_1 \propto (m-1)^{-\xi}. \quad (47)$$

However, closer inspection reveals that in practice the exponent $\xi$ is dependent on the maximum number of recovered components. In Fig. 5 we plot the fitted exponent, $\xi$, as a function of the number of recovered components. For only two recovered components, $\xi \approx 0.38$, while in the (practically unachievable) limit of a very large number of recovered components, $\xi$ asymptotically approaches 0.5.
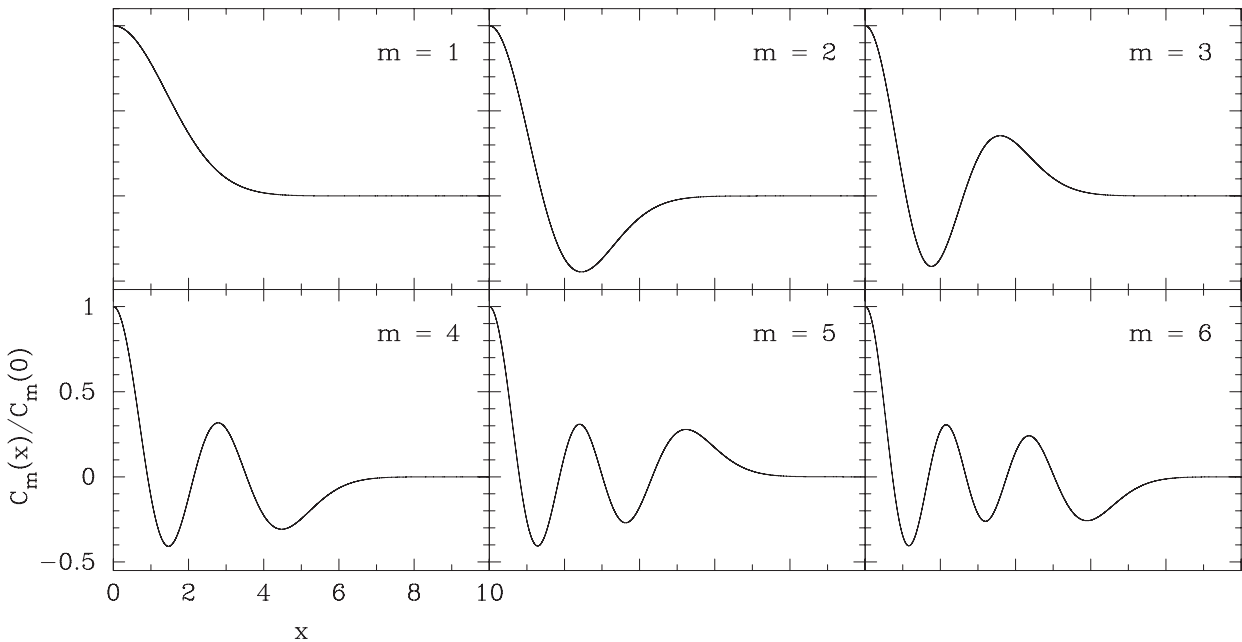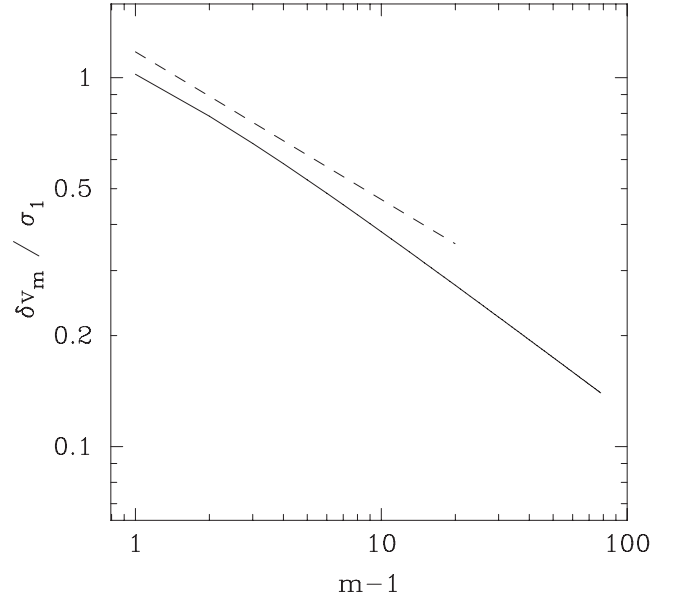


**Figure 4.** Log–log plot of the measured values of $\delta v_m/\sigma_1$ determined at the $1/e$ points of the eigenvector ACFs versus $m - 1$. For reference, the dashed line (offset) has a slope of $-\xi = -0.4$, appropriate for low orders $m$.

For a representative number of recovered components (between 3 and 20) in the calibration data of BH02, we adopt a working value of $\xi = 0.4 \pm 0.02$.

## 5 EIGENVALUES

For eigenvectors given by equation (39), it is possible to deduce the corresponding eigenvalues using equation (25). For the first two eigenvectors, equation (25) reads:

$$\int_{-\infty}^{+\infty} dv_k \, S(v_j, v_k)I_{01} \exp\left(-\frac{v_k^2}{2\sigma_1^2}\right) = \lambda_1 I_{01} \exp\left(-\frac{v_j^2}{2\sigma_1^2}\right), \quad (48)$$
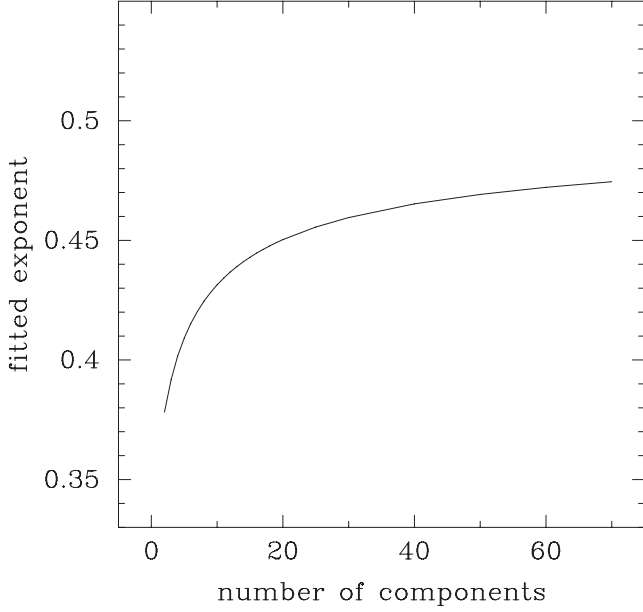


**Figure 3.** The first six eigenvector ACFs, given by equation (44).

**Figure 5.** The fitted exponent, $\xi$, from equation (47) as a function of the number of recovered components from which the fit is made.

$$\int_{-\infty}^{+\infty} \mathrm{d}v_k \, S(v_j, v_k) I_{01} \sqrt{2} \frac{v_k}{\sigma_1} \exp\left(-\frac{v_k^2}{2\sigma_1^2}\right)$$

$$= \lambda_2 I_{01} \sqrt{2} \frac{v_j}{\sigma_1} \exp\left(-\frac{v_j^2}{2\sigma_1^2}\right), \tag{49}$$

where $S(v_j, v_k)$ is given by equation (8). Making use of equations (27–34), these can be solved to find:

$$\lambda_1 = \sqrt{\frac{\pi}{a+c}} S_0, \tag{50}$$

$$\lambda_2 = \sqrt{\frac{\pi}{a+c}} \left(\frac{b}{a+c}\right) S_0, \tag{51}$$

which leads to

$$\frac{\lambda_2}{\lambda_1} = \frac{\sigma_{tot}^2}{\sigma_c^2} - \sqrt{\left(\frac{\sigma_{tot}^2}{\sigma_c^2}\right)^2 - 1}, \tag{52}$$

or

$$\frac{\sigma_c^2}{\sigma_{tot}^2} = \frac{2(\lambda_2/\lambda_1)}{1 + (\lambda_2/\lambda_1)^2}. \tag{53}$$

Equations (52) and (53), graphically represented in Fig. 6, show that, in the case of no centroid variation, all the variance of the data is contained in the first (and only) PC. The maximum value of $\lambda_2/\lambda_1 = 1$ is found in the limit where all variance in the data is caused by centroid variations. In general, the ratio $\lambda_2/\lambda_1$ can be used to provide a straightforward measurement of the ratio $\sigma_c^2/\sigma_{tot}^2$.

# 6 EIGENIMAGES

The covariance matrix and eigenvectors are independent of the spatial structure of the spectral line data. However, each eigenvector has an associated spatial map, the 'eigenimage', formed by projection of the data on to the eigenvector via equation (3). This can be alternatively viewed as the integration of the data over the velocity axis
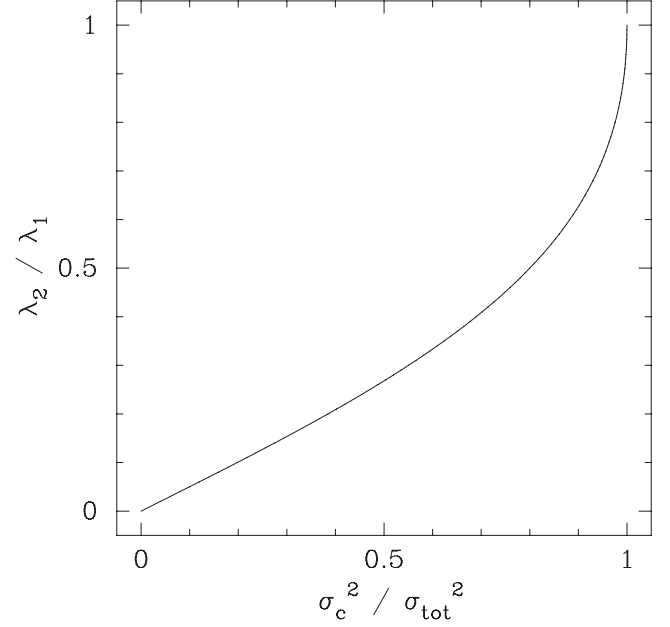


**Figure 6.** Relation between $\sigma_c^2/\sigma_{tot}^2$ and the ratio of the first two eigenvalues $\lambda_2/\lambda_1$.

with the eigenvector acting as a weighting or windowing function. For low-order eigenvectors, there is a straightforward interpretation of this procedure and it is possible to relate the resulting eigen-images to physical measures of the medium from which the line profiles originate. Below, we first derive the eigenimage structures for the two lowest order eigenvectors. Following this, we derive the asymptotic form of the eigenimages at high order.

## 6.1 Eigenimage structure

The form of the eigenimages, given by equation (3), is

$$I_m(\mathbf{r}) = \frac{I_{01}}{\sqrt{2^{m-1}(m-1)!}} \int_{-\infty}^{\infty} \mathrm{d}v \, T(\mathbf{r}, v)$$

$$\times \exp\left(-\frac{v^2}{2\sigma_1^2}\right) H_{m-1}(v/\sigma_1), \tag{54}$$

which can be interpreted as a generator of moments over the brightness temperature, subject to an overall windowing function $I_{01} \exp(-v^2/2\sigma_1^2) = u_1(v)$.

The first few Hermite polynomials are:

$$H_0(x) = 1$$

$$H_1(x) = 2x$$

$$H_2(x) = 4x^2 - 2$$

$$H_3(x) = 8x^3 - 12x$$

$$H_4(x) = 16x^4 - 48x^2 + 12, \tag{55}$$

so the first two eigenimages are closely related to the 0th and 1st moments – i.e. the integrals of $T(v)$ and $T(v)v$, respectively – but with the additional velocity-windowing provided by $u_1(v)$. Brunt (2003a) and Brunt, Heyer & Mac Low (2009) have made use of this to probe the outer scale of turbulence in molecular clouds, since the 0th moment and 1st moment are proportional to the column density and the projected momentum, respectively (see e.g.

Brunt & Federrath 2013, in preparation), with both subject to the same windowing function.

Writing $T(\mathbf{r}, v) = T_0(\mathbf{r}) \exp(-(v - v_c(\mathbf{r}))^2/2\sigma_b^2)$ and choosing the convenient normalization $I_{01} = 1$, the first eigenimage, $I_1(\mathbf{r})$ is

$$I_1(\mathbf{r}) = \sqrt{2\pi}\sigma_b T_0(\mathbf{r}) F(v_c), \tag{56}$$

where

$$F(v_c) = F_0^{1/2} \exp\left(-\frac{v_c^2}{2(\sigma_1^2 + \sigma_b^2)}\right) \tag{57}$$

is the integrated effect of the windowing function (with $F_0 = \sigma_1^2/(\sigma_1^2 + \sigma_b^2)$). Note that $\sqrt{2\pi}\sigma_b T_0(\mathbf{r}) = W_0(\mathbf{r})$ is the integrated intensity (0th moment) of the emission. The effect of $F(v_c)$ is to attenuate the eigenimage intensity for line profiles with high $|v_c|$.

The second eigenimage, $I_2(\mathbf{r})$, is

$$I_2(\mathbf{r}) = \sqrt{2\pi}\sigma_b T_0(\mathbf{r}) v_c(\mathbf{r}) \frac{F_0^{3/2}}{\sqrt{2}\sigma_1} F(v_c)$$
$$= \frac{F_0}{\sqrt{2}\sigma_1} I_1(\mathbf{r}) v_c(\mathbf{r}), \tag{58}$$

which is seen to be the 1st moment of the intensity again subject to the integrated effect of the windowing function.

Higher order eigenimages combine higher order moments, again with windowing by $u_1(v)$, but become increasingly difficult to interpret except in a statistical way. An approximate form for higher order eigenimages may be arrived at by making use of the following expansion at high $n$:

$$\exp\left(-\frac{x^2}{2}\right) H_n(x) \approx \frac{2^n}{\sqrt{\pi}} \Gamma\left(\frac{n+1}{2}\right) \cos\left(x\sqrt{2n} - n\frac{\pi}{2}\right). \tag{59}$$

Inserting this expression into equation (3) yields, after some manipulation:

$$I_m(\mathbf{r}) \approx G(m)\sqrt{2\pi}\sigma_b T_0(\mathbf{r}) \cos\left(\frac{v_c(\mathbf{r})}{\sigma_1}\sqrt{2(m-1)}\right) \text{ for odd } m,$$

$$I_m(\mathbf{r}) \approx G(m)\sqrt{2\pi}\sigma_b T_0(\mathbf{r}) \sin\left(\frac{v_c(\mathbf{r})}{\sigma_1}\sqrt{2(m-1)}\right) \text{ for even } m, \tag{60}$$

where $G(m)$ is an unimportant (constant) $m$-dependent multiplicative factor. While strictly only accurate at high $m$, these expressions provide a reasonably good representation of the eigenimage structure even at the lowest $m$-values (though quantitatively, the differences are important as we discuss in the next section). Note that for small $v_c/\sigma_1$, both equations (58) and (60) give $I_2 \propto T_0 v_c$. In addition, the windowing term, $F_{v_c}$, in equation (58) crudely approximates the roll-off in $I_2$ caused by the sinusoidal behaviour in equation (60).

The structure of the eigenimages predicted by equation (60) is as follows. The overall amplitude (at any order $m$) is controlled by the column density [$\sqrt{2\pi}\sigma_b T_0(\mathbf{r})$], and this is modulated by a common multiplicative factor (dependent on $m$) and, more importantly, a sine or cosine factor, dependent on the centroid velocity, $v_c(\mathbf{r})$. Therefore, as the order $m$ increases, the eigenimage values cycle through a sine or cosine variation. This provides the key to understanding their characteristic spatial scale lengths needed for the measurement of the turbulent velocity spectrum, as described in the next section.

## 7 ANALYTIC CALIBRATION OF THE PCA METHOD FOR THE TURBULENT VELOCITY SPECTRUM

Our procedure here is to generate a coupled sequence of characteristic spatial and velocity scales $(\delta l_m, \delta v_m)$ at order $m$, for a specified spectral index $\beta$ of the 3D velocity field. The dependence of the predicted exponent $\alpha$ (where $\delta v_m \propto \delta l_m^\alpha$) on the intrinsic $\beta$ will then establish the calibration (see Section 2.2).

We have already established the $m$-dependence of $\delta v_m$ in Section 4.2, where it was found that $\delta v_m \propto (m-1)^{-\xi}$ with $\xi \approx 0.4$. It still remains to determine the corresponding sequence $\delta l_m$. Here, however, while we have a functional form for the asymptotic eigenimage structure (equation 60) we do not have a definite expression for the field $v_c(\mathbf{r})$, but instead only have a statistical knowledge of its properties, which may be quantified via structure functions.

The $p$th-order structure function of a velocity field is written as:

$$S_p(l) = \langle|\Delta v(l)|^p\rangle \propto l^{\zeta_p}, \tag{61}$$

where $\Delta v(l)$ represents the ensemble of velocity fluctuations measured on spatial scale $l$ in the field, and angle brackets denote spatial averaging. The function $\zeta_p$ describes the dependence of the scaling exponent on the order $p$. Alternatively, one may write:

$$(S_p(l))^{1/p} = \langle|\Delta v(l)|^p\rangle^{1/p} \propto l^{\gamma_p}, \tag{62}$$

where $\gamma_p = \zeta_p/p$. For velocity fields produced by fBm, $\gamma_p$ is independent of $p$ (e.g. Brunt et al. 2003). For now, we will assume that the centroid velocity field, $v_c(x, y)$, can be described by a scaling exponent $\gamma_c$ (valid at all $p$), allowing us to write:

$$\langle|\Delta v_c(l)/\sigma_1|^p\rangle^{1/p} = (l/l_1)^{\gamma_c}, \tag{63}$$

where $l$ is the 2D spatial scale and $l_1$ is the spatial scale corresponding to a mean velocity fluctuation of $\sigma_1$.

The original calibration established by BH02 used uniform density fields (and therefore uniform column density fields) so that only the effect of the (co)sine term in equation (60) need be inspected. The (co)sine term leads to an oscillatory eigenimage structure with a characteristic spatial wavelength $L_m$ set by the condition that *the typical velocity fluctuation between points separated by a distance $L_m$ generates a phase difference of $2\pi$ in the argument of the (co)sine term*. That is:

$$\sqrt{2(m-1)}\langle\Delta v_c(L_m)\rangle/\sigma_1 \approx 2\pi. \tag{64}$$

Referring to equation (60), note that because the $\sqrt{2(m-1)}$ factor effectively *amplifies* the $v_c$ field, progressively smaller velocity fluctuations are capable of inducing a $2\pi$ phase difference as the order $m$ increases [i.e. the typical $v_c$ fluctuation required falls proportionally to $(m-1)^{-1/2}$]. Consequently, there is a corresponding reduction in the characteristic spatial wavelength, governed by equation (63), such that

$$L_m/l_1 \approx \langle|\Delta v_c(L_m)/\sigma_1|\rangle^{1/\gamma_c} \approx \left(\frac{\sqrt{2}\pi}{(m-1)^{1/2}}\right)^{1/\gamma_c}, \tag{65}$$

meaning that the characteristic wavelength of eigenimage structure decreases with order $m$ as $L_m \propto (m-1)^{-1/2\gamma_c}$.

The characteristic spatial scale, $\delta l_m$, of the $m$th-order eigenimage is determined by the $1/e$ point of the eigenimage ACF, and it is straightforward to show that for a (co)sinusoid:

$$\delta l_m = \left(\frac{\text{acos}(1/e)}{2\pi}\right) L_m \approx 0.19 L_m. \tag{66}$$

Therefore, the $m$-dependence of characteristic eigenimage scales, in the asymptotic approximation, is

$$\delta l_m \propto (m-1)^{-1/2\gamma_c}, \tag{67}$$

where $\gamma_c$ is the scaling exponent of the centroid velocity field. However, this is slightly inaccurate as the asymptotic expansions are not strictly applicable at low-order $m$. We note first that, crudely approximating the $v_c$ field as a $\sim$ linear gradient, the exact equation (58) predicts a scale $\delta l_2$ that is 20 per cent larger than that predicted by equation (60). Since as the order $m$ increases, the asymptotic formula becomes increasingly more accurate; this in effect means that $\delta l_m$ falls faster with $m-1$ than equation (67) predicts. Assuming a smooth transition from an $\sim$20 per cent over-estimation at low $m$ to accurate representation at, say, $m \gtrsim 10$, we estimate that the effective $m$-dependence of $\delta l_m$ is better represented by:

$$\delta l_m \propto (m-1)^{-1.1/2\gamma_c}, \tag{68}$$

i.e. an increase of the exponent, by a factor of 1.1 ($\pm 0.03$), describing the reduction of characteristic spatial scale as the order increases.

Combining equation (68) with the $m$-dependence of the characteristic velocity scales (equation 47), we arrive at a calibration of the PCA $\alpha$ exponent to the centroid velocity scaling exponent, $\gamma_c$, via:

$$\delta v_m \propto \delta l_m^\alpha \propto \delta l_m^{2\xi\gamma_c/1.1}, \tag{69}$$

so that:

$$\alpha \approx 2\xi\gamma_c/1.1, \tag{70}$$

and taking the representative value $\xi = 0.4 \pm 0.02$, as discussed in Section 4.2, this leads to

$$\alpha \approx 0.72\gamma_c. \tag{71}$$

It remains to relate $\gamma_c$ to the spectral index, $\beta$, of the 3D velocity field. This is a general question (not restricted to the PCA method) but one that has a simple answer in the uniform density conditions assumed by BH02 in the original calibration. As explained in Brunt & Mac Low (2004, and references therein), the following relation holds for uniform density and optically thin conditions:

$$\gamma_c = \frac{\beta}{2}. \tag{72}$$

Some discussion of this equation is warranted, as the scaling exponent of the velocity field in 3D ($\gamma_{3D}$, here assumed independent of $p$, appropriate for the non-intermittent fBm fields used by BH02) is given by

$$\gamma_{3D} = \frac{\beta-1}{2}, \tag{73}$$

and therefore:

$$\gamma_c = \gamma_{3D} + \frac{1}{2} = \frac{\beta}{2}. \tag{74}$$

The increase in the exponent upon projection (by $1/2$) is known as 'projection smoothing', and can be qualitatively understood by considering that large-scale velocity fluctuations suffer proportionally less line-of-sight averaging than small-scale fluctuations.

Using equations (72) and (71) we arrive at the analytic calibration of the PCA $\alpha$ exponent:

$$\alpha \approx 0.36\beta. \tag{75}$$

This relation is close to, though slightly steeper than, the empirically determined $\alpha \approx (0.33 \pm 0.04)\beta$ (BH02; Roman-Duval et al. 2011).

This is encouraging analytic support for the empirical calibration, and the small difference in exponent ($0.36 \pm 0.04$ versus $0.33 \pm 0.04$) is not too concerning, given the approximations used in the derivations above.

In the above, we have not explicitly included the effects of opacity, and it is worth considering how this may affect the result. Previously, it has been found empirically that opacity/saturation does not have a drastic effect on $\alpha$ (Brunt et al. 2003; Roman-Duval et al. 2011). It is also observed that application of the method to $^{12}$CO and $^{13}$CO data on the same cloud yields very similar $\delta v(\ell)$ spectra and similar values of $\alpha$ (e.g. Brunt 2003b; Brunt & Mac Low 2004; Brunt et al. 2009). A likely reason for this insensitivity is that the centroid velocity field is not strongly affected by saturation *if* the saturation is symmetric about line centre. Brunt & Mac Low (2004) demonstrate directly that the centroid fields derived in their observations from $^{12}$CO and $^{13}$CO are almost indistinguishable statistically. A secondary effect of saturation may be to move the line profiles to a flat-topped appearance, invalidating the Gaussian form assumed above. However, the requirement of orthogonality in the eigenvectors essentially ensures a polynomial sequence similar to the derived Hermite polynomials, so any deviations from our scaling result will likely be small. However, we cannot analytically assess this at present, and must rely on the empirical/observational results.

Finally, we comment on two other aspects of the PCA method for which a better understanding is now available in light of the above analysis. First, Brunt et al. (2003) found that PCA appears to operate at first order – i.e. in the case of an intermittent field when $\gamma_1 \neq \gamma_2$, the PCA exponent $\alpha$ is better-correlated with the first-order index $\gamma_1$. This can be now understood to be related to the 'phase-rolling' effect [i.e. the $m$-dependent amplification of velocity fluctuations to roll the (co)sinusoid phase of the eigenimage structure] discussed above, which is a first-order effect rather than a root-mean-square effect. Secondly, it has been shown empirically that the recovered PCA exponent $\alpha$ is not strongly affected by (column) density fluctuations (BH02; Roman-Duval et al. 2011). While a full analysis of this effect is beyond the scope of the current paper, a preliminary understanding of why this is can be arrived at by considering the eigenimage structure given by equation (60). An eigenimage of order $m$ is the product of the column density ($m$-independent) and the (co)sinusoid ($m$-dependent). The ACF of such an eigenimage is the Fourier transform of its power spectrum, which in turn is the square of its Fourier transform. A product in direct space transforms to a convolution in Fourier space, so the quantity of interest [the Fourier transform of the (co)sinusoid] is convolved with the Fourier transform of the column density – a function that is independent of order $m$. In the case of uniform column density, this function is a delta function and the transform of the (co)sinusoid is unchanged. As column density fluctuations become more important, a broadening of the column density transform is induced, but as long as this remains narrow (in Fourier space) relative to the (co)sinusoid transform's Fourier-space width, no significant effect on the combined power spectrum (and therefore ACF) will be induced. However, for an extremely variable column density field with a broad Fourier space extent (as examined by Roman-Duval et al. 2011) this must eventually break down. Roman-Duval et al. (2011) determine that a density field with a lognormal PDF with $\sigma_{\ln(\rho/\rho_0)} > 2$ is required for this to occur (see their fig. 5).

## 8 SUMMARY

In this paper, we have derived and discussed analytic expressions for covariance matrices, eigenvectors, eigenvalues and eigenimages

expected from PCA of molecular cloud emission lines, in the limit where these can be represented by a collection of Gaussian line profiles with turbulent spatial correlations. Previous to this study, the PCA method was based almost entirely on empirical analysis and lacked a firm theoretical basis.

We have derived an analytic calibration of the PCA method for measuring the spectrum of turbulent velocity fluctuations, which agrees reasonably well with previous empirical calibrations. However, given the level of approximation in the analysis, we see the analytic calibration more as a validation of the empirical calibration, rather than a replacement. We have also gained significant insight into the mechanisms by which PCA operates, allowing us to explain more esoteric aspects of the method, such as its preferential operation at first order and its general robustness against (column) density fluctuations.

## ACKNOWLEDGEMENTS

## REFERENCES

Brunt C. M., 2003a, ApJ, 583, 280
Brunt C. M., 2003b, ApJ, 584, 293
Brunt C. M., Heyer M. H., 2002a, ApJ, 566, 276 ( BH02)
Brunt C. M., Heyer M. H., 2002b, ApJ, 566, 289
Brunt C. M., Mac Low, 2004, ApJ, 604, 196
Brunt C. M., Heyer M. H., Vázquez-Semadeni E., Pichardo B., 2003, ApJ, 595, 824
Brunt C. M., Heyer M. H., Mac Low, 2009, A&A, 504, 883
Heyer M. H., Brunt C. M., 2004, ApJ, 615, L45
Heyer M. H., Brunt C. M., 2012, MNRAS, 420, 1562
Heyer M. H., Schloerb F. P., 1997, ApJ, 475, 173 ( HS97)
Heyer M. H., Brunt C. M., Snell R. L, Howe J. E., Schloerb F. P., Carpenter J. M., 1998, ApJS, 115, 241
Heyer M. H., Gong H., Ostriker E., Brunt C. M., 2008, ApJ, 680, 420
Kleiner S. C., Dickman R. L., 1985, ApJ, 295, 466
Lazarian A., Pogosyan D., 2000, ApJ, 537, 720
Miesch M. S., Bally J., 1994, ApJ, 429, 645
Roman-Duval J., Federrath C., Brunt C. M., Heyer M. H., Jackson J. M., Klessen R. S., 2011, ApJ, 740, 120
Scalo J. M., 1984, ApJ, 277, 556
Stutzki J., Bensch F., Heithausen A., Ossenkopf V., Zeilinsky M., 1998, A&A, 336, 697

This paper has been typeset from a TEX/LATEX file prepared by the author.