Detection of Long-Range Concerted Motions in Protein by a Distance Covariance

Amitava Roy* and Carol Beth Post

Medicinal Chemistry and Molecular Pharmacology, Purdue University, West Lafayette, Indiana, United States

Supporting Information

ABSTRACT: We asses the ability of a distance correlation coefficient (DiCC), calculated from distance covariance, for detecting long-range concerted motion in proteins. We establish a set of criteria for ideal correlation coefficient values based on the coefficient of determination in multidimension, \mathbf{R}^2 . We compare in detail DiCC and conventional correlation coefficients against these criteria. We demonstrate that, in contrast to conventional correlation coefficients, which capture long-distance correlation adequately only with certain restrictions in multidimension, DiCC reflects appropriate correlation in both one-dimension and multidimensions. Finally, we demonstrate the usefulness of DiCC for assessing long-distance correlated fluctuation in protein dynamics.

1. INTRODUCTION

Concerted, low-frequency motions are inherent to large or multidomain proteins and can be essential for proteins to carry out their function;¹ particularly those involving allosteric processes. Large-scale, concerted motion implies correlated fluctuation of different parts of the protein separated by relatively long-distance. Atomistic molecular dynamics (MD) simulation of proteins has the potential for revealing concerted motions in great detail. Nonetheless the assessment of long-range correlated motion from simulations has so far been elusive except for a small number of cases.² A correlation coefficient (CC) can be defined to quantify correlation between two random variables, including atomic fluctuations in the case of proteins. The most widely used CC between scalar variables is Pearson's correlation coefficient (PCC). The displacement vector correlation coefficient (VCC) is an extension of PCC to quantify correlation between two positional vectors. Some recent insightful usages of VCC are reported in refs 3–5. VCC depends on the cosine of the angle between the vectors and is most sensitive when the vectors are parallel.^{2,6,7} To overcome this shortcoming of VCC, a few studies have used the generalized correlation coefficient (GCC)⁷⁻⁹ or radial correlation coefficient $(RCC)^2$ to detect correlation between atomic fluctuations of proteins. In previous work,² we exploited the radial symmetry of icosahedral viral capsids and found long-range correlated motions between residues 55 Å apart in human rhinovirus using RCC, which is a PCC on the norm of position vectors. RCC is highly useful when applied to systems with radial symmetry, but it is insensitive to azimuthal fluctuation. GCC is an excellent CC between scalar random variables; however, in multidimensions, GCC does not combine the one-dimensional CCs in a suitable way to investigate concerted motions. In this article, we asses the ability of a distance correlation coefficient (DiCC),^{10,11} calculated from distance covariance, to capture correlation without imposing any assumption on the time series of the vectors. A comparison of DiCC with VCC, RCC, and GCC elucidates the merit and weaknesses of each and the potential of DiCC for detecting long-range concerted motion in proteins.

2. RESULTS

2.1. Correlation Coefficients. DiCC between two vector series, $\{A\}$ and $\{B\}$, is defined as

$$DiCC = \frac{\nu(\mathbf{A}, \mathbf{B})}{\sqrt{\nu(\mathbf{A}, \mathbf{A})\nu(\mathbf{B}, \mathbf{B})}}$$
(1)

where $\nu(\mathbf{A},\mathbf{B})$ is the distance covariance between the vectors. Let us assume that the vector series, $\{A\}$ and $\{B\}$ have *n* entries each and the *i*th entry in $\{A\}$ is denoted by A^i . If $\{A\}$ and $\{B\}$ are position vectors of two atoms from a simulation study then \mathbf{A}^i is the *i*th saved position vector of one atom. Distance covariance is defined as

$$u(\mathbf{A}, \mathbf{B}) = \sqrt{\frac{1}{n^2}} \sum_{ij}^{ij} \alpha_{ij} \beta_{ij}$$

where

$$\alpha_{ij} = a_{ij} - a_{i.} - a_{.j} + a_{..} \tag{2}$$

The following steps are needed to calculate α_{ii} from {A}.

- 1. Build the $n \times n$ matrix, a, from {A}, where a_{ii} is the distance between the *i*th and *j*th entries of $\{A\}$: $a_{ii} = |A^i - A^j|$
- 2. Average the rows of **a**: $a_i = (1/n)\sum_j a_{ij}$ 3. Average the columns of **a**: $a_j = (1/n)\sum_j a_{ij}$
- 4. Average all elements of **a**: $a_{i} = (1/n^2) \overline{\sum}_{ij} a_{ij}$
- 5. Build the $n \times n$ matrix $\boldsymbol{\alpha}$ from **a** where $\alpha_{ii} = a_{ii} a_{i} a_{j} + a_{i}$

VCC and GCC between the vector series $\{A\}$ and $\{B\}$, and RCC between $\{A_r\}$ and $\{B_r\}$, the norms of **A** and **B**, respectively, are defined as follows:

$$VCC = \frac{\langle (\mathbf{A} - \langle \mathbf{A} \rangle) (\mathbf{B} - \langle \mathbf{B} \rangle) \rangle}{\sqrt{\langle (\mathbf{A} - \langle \mathbf{A} \rangle)^2 \rangle \langle (\mathbf{B} - \langle \mathbf{B} \rangle)^2 \rangle}}$$
$$GCC = \sqrt{1 - e^{-2I/d}}$$

Received: July 4, 2012 Published: August 21, 2012

Journal of Chemical Theory and Computation

$$RCC = \frac{\langle (A_r - \langle A_r \rangle)(B_r - \langle B_r \rangle) \rangle}{\sqrt{\langle (A_r - \langle A_r \rangle)^2 \rangle \langle (B_r - \langle B_r \rangle)^2 \rangle}}$$
(3)

where $\langle ... \rangle$ is the ensemble average or average over all entries in $\{A\}$ and $\{B\}$, *I* is the mutual information between $\{A\}$ and $\{B\}$, calculated using the method developed by Kraskov, Stogbauer, and Grassberger, ^{12,13} and *d* is the dimension of vectors **A** and **B**.

It should be noted that calculation of VCC and GCC require the dimensions of **A** and **B** to be the same, while the calculation of RCC and DiCC does not impose any such restriction.

2.2. Coefficient of Determination. If the dependency between two scalar random variables is known, then the coefficient of determination,¹⁴ R^2 , can be considered a measure of correlation between the variables. In the case of a linear dependency, R^2 is 1 when the variation in one of the variables can be determined exactly by the variation in the other, and R^2 is zero when the variation in one cannot be determined at all by the variation in the other. If the dependency is nonlinear, we can refer to Nagelkerke.¹⁵

 R^2 between two scalar variables is a scalar quantity. In the case of two vectors of dimension *m* and *n*, we can define \mathbf{R}^2 as a $m \times n$ matrix where R_{ij}^2 , the *ij*th component of the matrix, is the coefficient of determination between the *i*th component of one vector and the *j*th component of another vector. An example of such a matrix is given later in the article.

A CC indicates the strength of the relationship between random variables. For a CC to be practical, physically meaningful and robust in the context of atomic fluctuations, it should satisfy the following criteria:

- be a scalar quantity
- equal 1 when **R**² is a unity matrix, and the dependency between the random variables is linear
- equal 0 when \mathbf{R}^2 is a null matrix
- if **R**² between a pair of vectors is identical to **R**² between another pair of vectors, then the CC should be similar in both cases
- be independent of coordinate system

While we were developing this assessment based on R^2 of coefficients for detecting concerted motions in proteins, a study appeared¹⁶ in which similar criteria were proposed to establish associations between scalar data sets. Reshef et al. showed that, for nonlinearly dependent random variables, none of the established CCs becomes 1 even when R^2 is 1, and the sensitivity of CC calculated with different methods depends on the specific functional form of the dependency.¹⁶ Accordingly, for practical purposes, we demanded the second criteria stated be true only for linearly dependent random variables, although we would like it to be true in general.

2.3. Correlation Coefficients in Multidimensions. To compare the performances of different CCs, we calculated the CC between the positions of two particles *A* and *B* specified by their two-dimensional position vectors, **A** and **B**, as shown in Figure 1. We can write

$$\mathbf{A} = A_r \hat{r} + A_\theta \hat{\theta} = A_x \hat{i} + A_y \hat{j}$$
$$\mathbf{B} = B_r \hat{r} + B_\theta \hat{\theta} = B_x \hat{i} + B_y \hat{j}$$
(4)

where \hat{i} and \hat{j} are unit vectors in Cartesian coordinate system and \hat{r} and $\hat{\theta}$ are unit vectors in the spherical coordinate system. The value of the CC obtained from the different parameters are compared to the known coefficient of determination between the



Figure 1. *A* and *B* are two particles with their position vectors **A** and **B** respectively. $B_r = A_r + \Delta r + \delta r$ and $B_\theta = A_\theta + \Delta \theta$ where Δr and $\Delta \theta$ are two constants and δr is a normally distributed noise with mean zero and variance σ_r^2 .

components of the vector. If B_r can be expressed as a linear function of A_{rr} $f(A_r)$, then the coefficient of determination, $R^2(B_{rr}A_r)$ is¹⁴

$$R^2(B_r, A_r) \equiv 1 - \frac{\sigma_{\rm err}^2}{\sigma_{\rm tot}^2}$$

where

$$\sigma_{\text{tot}}^2 = \langle B_r - \langle B_r \rangle \rangle^2 \quad \sigma_{\text{err}}^2 = \langle B_r - f(A_r) \rangle^2 \tag{5}$$

In the two-dimensional model, we define

$$B_r = f(A_r) = A_r + \Delta r + \delta r$$

$$B_\theta = f(A_\theta) = A_\theta + \Delta \theta$$
(6)

where Δr and $\Delta \theta$ are constants and δr is a random variable normally distributed, with mean zero and variance σ_r^2 .

To build a series {A}, we generated 100 000 normally distributed values of A_r with mean value of 10 and variance of $\sigma_{A_r}^2 = 36$. We fixed the value of A_{θ} to $\pi/4$. We independently generated another 100 000 normally distributed values, with a mean of 0 and variance of $\sigma_r^2 = 16$, to build { δr }. For a particular value of $\Delta \theta$, we built {B} from eq 6 with $\Delta r = 3.0$. We generated 90 such series of {A} and {B} while varying the value of $\Delta \theta$ from 0 to $\pi/2$. In all 90 series, $\sigma_{tot}^2 = \sigma_{A_r}^2 + \sigma_r^2$ and $\sigma_{err}^2 = \sigma_r^2$ in B_r . Hence, $R^2(B_{rr}A_r) = 1 - \sigma_r^2/(\sigma_{A_r}^2 + \sigma_r^2) = 0.69$. Variances of B_{θ} and A_{θ} are zero as their values are fixed. Also, σ_{err}^2 is zero in B_{θ} . We can still define $R^2(B_{\theta}A_{\theta})$ in such a case from the limit $\sigma_{tot}^2 \rightarrow 0$, $R^2(B_{\theta}A_{\theta})$ becomes 1 when $\sigma_{err}^2 = 0$. Since angular and radial components are independent of each other, $R^2(B_{\theta A_r})$ and $R^2(B_{rr}A_{\theta})$ are zero. So, the expected R^2 between \hat{r} and $\hat{\theta}$ components are

$$\mathbf{R}^{2}(\mathbf{B}, \mathbf{A}) = \frac{B_{r}}{B_{\theta}} \begin{pmatrix} A_{r} & A_{\theta} \\ 0.69 & 0 \\ 0 & 1 \end{pmatrix}$$

We calculated the DiCC, VCC, and GCC of (**B**,**A**) and DiCC, RCC, and GCC of ($B_{rr}A_{r}$). The values of the correlation coefficients between (**B**,**A**) are plotted as a function of $\Delta\theta$ in 2. For reference, the PCC of two linearly dependent random scalars is equal to (R^2)^{1/2}, the square root of the coefficient of determination between them, which is 0.83 here. DiCC of ($B_{rr}A_{r}$) and (**B**,**A**), dotted and solid blue lines in Figure 2, respectively, have identical values of 0.81. Uncertainty in determining B_r from A_r in one dimension and **B** from **A** in multidimension appears only due to the random variable δr , and the DiCC values in one dimension and multidimensions correctly reflect that. In Figure 2, RCC of ($B_{rr}A_r$) (red dotted line) and VCC when $\Delta\theta = 0$ (green solid line)



Figure 2. Red dotted line represents GCC and RCC of $(B_{\mu}A_{r})$. As both the values are very close to each other, only one line is drawn for clarity. Blue solid and dotted lines represent DiCC of $(\mathbf{B}_{r}\mathbf{A})$ and $(\mathbf{B}_{rr}\mathbf{A}_{r})$. Dotted blue line is hardly visible, as it overlaps with the solid blue line. Solid red and green lines represent GCC and VCC of $(\mathbf{B}_{r}\mathbf{A})$ respectively. VCC of $(\mathbf{B}_{r}\mathbf{A})$ depends on $\Delta\theta$, the angle between **B** and **A**. The inset shows how the RCC of $(B_{rr}A_{r})$, when angle between $(\mathbf{B} \text{ and } \mathbf{A})$ is $\pi/2$, changes as the origin of the coordinate axis moves along *x*-axis. As the origin changes the radial component of the vectors decreases, as does RCC.

become exactly $(R^2)(B_r,A_r)$. VCC, however, decreases monotonically to 0 as $\Delta\theta$ increases from 0 to $\pi/2$ and changes sign for $\pi/2 < \Delta\theta < \pi$. In multidimensions, VCC between two random vectors is the VCC value when the vectors are parallel multiplied by the cosine of the angle between them.

RCC is independent of $\Delta\theta$ and reproduces *R*; however, it depends on the position of the origin of the coordinate system as the definition of the radial component of motion depends on the position of the origin. To illustrate this limitation, we used the series of {A} and {B} for $\Delta\theta = \pi/2$ and calculated RCC of (B_r,A_r) , while moving the origin along the *x*-axis. The inset of Figure 1 shows how RCC changes as a function of the position of the origin on the *x*-axis of original coordinate frame.

That a GCC-like quantity can be used to define correlation between Gaussian random scalars was first suggested by Joe.¹⁷ For this case, the GCC of $(B_{r,}A_r)$ becomes exactly $R(B_{r,}A_r)$, as evident in Figure 2. However, in multidimensions, the GCC of (\mathbf{B},\mathbf{A}) is much higher, even though the source of uncertainty in determining B_r from A_r is the same in determining **B** from **A**. If the \mathbf{R}^2 matrix is diagonal and λ_i^2 are it is diagonal elements, then GCC = $(1 - (\prod_i (1 - \lambda_i^2))^{1/d})^{1/2}$, where *d* is the dimension of the vectors. The derivation and the physical meaning of the above relation is explained in the Supporting Information, Note 1. Accordingly, the GCC of $(\mathbf{B},\mathbf{A}) = (1 - ((1 - 0.69)(1 - 1))^{1/2})^{1/2} = 1$. Irrespective of the value of $R^2(B_{r,}A_r)$, the GCC of (\mathbf{B},\mathbf{A}) is 1, as $R^2(B_{\theta i}A_{\theta i}) = 1$. The scheme with which GCC combines one-dimensional PCC values is not suitable to find association between positional fluctuations.

In the model illustrated with Figure 1, $\rho_{A_xA_y}$ and $\rho_{B_xB_y}$ are close to one. In protein dynamics, however, CCs between the components of a vector are usually much smaller. Distribution of CCs between the components of position vectors calculated from a protein dynamics simulation is given in the Supporting Information, Note 2 and Figure S1. To check the performance of DiCC and GCC as the correlation between the components changes, we generated (**B**,**A**) with

$$\mathbf{R}^{2}(\mathbf{B}, \mathbf{A}) = \frac{B_{x}}{B_{y}} \begin{pmatrix} A_{x} & A_{y} \\ 0.49 & 0.49 \\ 0 & 0 \end{pmatrix}$$

while varying the correlation between $\rho_{A_{x}A_{y}}$ from 0.0 to 0.99.

DiCC of (**B**,**A**), shown in the solid blue line in Figure 3, changes very slightly as $\rho_{A_xA_y}$ varies from 0 to 0.99. The GCC of



Figure 3. Two random variables **A** and **B** are generated with $R^2(A_{xy}B_x) = 0.49$, $R^2(A_{yy}B_x) = 0.49$, and $R^2(A_{xy}B_y) = 0$, and $R^2(A_{yy}B_y) = 0$, while varying $\rho_{A_xA_y}$ from 0 to 0.99. Solid blue and red lines represents DiCC and GCC of (**B**, **A**), respectively. DiCC of (**B**, **A**) changes very slightly as $\rho_{A_xA_y}$ varies from 0 to 0.99. GCC of (**B**, **A**) is 0.93 when $\rho_{A_xA_y}$ is equal to 0 and approaches DiCC of (**B**, **A**) as $\rho_{A_xA_y}$ tends toward 1.

(**B**,**A**), shown in solid red line in Figure 3, is 0.93 for $\rho_{A_xA_y}$ equal to 0 and approaches DiCC of (**B**,**A**) as $\rho_{A_xA_y}$ tends toward 1. For the same **R**² matrix in Figure 3, DiCC of (**B**,**A**) varies slightly, from 0.63 to 0.52, while GCC by contrast varies greatly from 0.93 to 0.52.

2.4. Long-Range Correlated Fluctuations in Protein. We further compared the capability of the various CC parameters using the example of Src SH2 domain in complex with a conformationally constrained mimetic of a phosphotyrosyl tetrapeptide ligand pYEEI,^{18,19} and show that DiCC detects longrange concerted motions that are underestimated by VCC. VCC and DiCC were calculated between 106 C_{α} atoms from the cumulative 80 000 conformations. In Figure 4a, VCC and DiCC are plotted against the average distance between the C_{α} pairs. The DiCC values are overall much greater than VCC values. No C_{α} pairs, with average distances between them greater than 7.5 Å, have a VCC value greater than 0.6. On the other hand, there are more than 40 pairs of C_{α} pairs, shown as circles in Figure 4a, separated by more than 7.5 Å and have DiCCs greater than 0.6. One C_{α} pair, shown as a diamond in Figure 4a, with an average distance equal to 24.8 Å has a DiCC value of 0.58.

The GCC values are also much greater than VCC values, and the distribution is less disperse (Figures 4b and 5). The GCC value is also greater in general than the DiCC value (Figures 4c and 5). The increased value of GCC arises because GCC is dominated by the largest element of \mathbf{R}^2 calculated in a coordinate system where \mathbf{R}^2 is diagonal, as explained earlier and in the Supporting Information, Note 1. Accordingly GCC, without an effective scheme of combining one-dimensional PCC values, does not characterize correlated behavior in a manner suitable to find association between positional fluctuations. GCC reflects some kind of correlation between random vectors, but it is not clear given the tight distribution in Figures 4b and the variation with respect to nonindependent vector components (Figure 3), how useful it is to detect long-range concerted fluctuations in protein dynamics.

2.5. Convergence of Correlation Coefficient Values. We investigated convergence behavior of different CCs of five pairs of C_{α} atoms whose DiCC values fall within 0.9–1.0, 0.8–0.9,

Wonpil: Here, GCC is Grubmuller's generalized correl coeff. GCC is not good for correlating a vector quantity such as atomic position. The problem with GCC is it is dominated by the largest correlation of individual vector components. So, for example, if the x component of A and x component of B are highly correlated, but the correlation of the y components is near zero, you still get a high correlation for the vector using GCC. That behavior is not a good property for a correlation coefficient. Journal of Chemical Theory and Computation



Figure 4. (a) VCC (blue dots) and DiCC (red dots) between C_a atoms of from 40 × 2 ns long trajectory of Src SH2 domain in complex with the ligand pYEEI (see text for details) plotted against average distance between the C_a atoms. Distance correlation reveals more than 40 pairs of C_a atoms with average distances between them >7.5 Å and DiCC > 0.6 (circles). One pair of C_a atoms, showed in diamond, has an average distance of 24.78 Å and DiCC of 0.58. DiCC reveals long distance correlations, which are underestimated by VCC. (b) VCC (blue dots) and GCC (orange dots) of C_a atoms plotted against average distances between the C_a atoms. GCC does not reflect correlation suitable to investigate concerted fluctuation of positional vectors of C_a atoms. (c) DiCC and GCC of C_a atoms. Comparisons between VCC and DiCC and VCC and GCC are given in Supporting Information, Figure S2.

0.7–0.8, 0.6–0.7, and 0.5–0.6, respectively and have the highest intrapair average distances among all C_{α} pairs with DiCC in their respective ranges. The five pairs of C_{α} atoms are from residues 164 and 165, residues 153 and 154, residues 206 and 216, residues 206 and 216, and residues 193 and 203, with average intrapair distances 3.87 Å, 3.85 Å, 11.13 Å, 11.08 Å, and 24.71 Å, respectively. We combined *n* ps from the beginning of the 40 trajectories and calculated VCC, GCC, and DiCC, and their bootstrap standard deviation of the five C_{α} pairs from the combined data, while varying *n* from 10 to 2000 with a step of 20. Mean values of CCs calculated from the combine data stabilizes with 40 × 400 ps of data (Figure 6). Standard deviations

calculated from 400 bootstrap sample change by less than 0.01 for all CCs during last 40 \times 1 ns of data (Figure 6). Convergence behavior of all CCs are similar, and they converge well with 40 \times 1 ns of data.

3. DISCUSSION

A correlation coefficient should be able to characterize correlation between displacement vectors due to concerted motion in a protein regardless of the distance of separation. VCC depends on the angle between the position vectors and hence underestimates the correlation when vectors are not parallel. While RCC is highly suitable for detecting radial motion in spherically symmetric system, it depends on the position of origin of the coordinates axis and is insensitive to azimuthal correlation. When radial symmetry does not dominate the concerted motion, RCC does not reflect the full correlation between positional vectors. GCC is dominated by largest element of diagonalized \mathbf{R}^2 matrix and not suitable to find association between vectors.

The CC best matching the criteria outlined is DiCC calculated from distance covariance. DiCC was found here to capture the true correlation between positional vectors based on agreement with R^2 , is insensitive to the angle between the displacement vectors, and has limited sensitivity to the dependence between the vector components. Further DiCC reflects both linear and nonlinear correlation.¹⁶ Using DiCC we observe long-distance concerted motions in a protein that was not revealed by VCC. Detection of such collective motion, which has mostly been elusive in analyses of molecular dynamics simulation, can be insightful for understanding allosteric function and other long-distance effects in proteins.

4. METHODS

4.1. Generating Correlated Gaussian. We determine $\{B\}$ and $\{A\}$ with a specified covariance matrix C between $A_{xx} A_{yy} B_x$ by defining

$$\mathbf{C} = \mathbf{W}\mathbf{W}^{\dagger}$$

$$\begin{pmatrix} A_x \\ A_y \\ B_x \end{pmatrix} = \mathbf{W} \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix}$$
(7)

where v_1 , v_2 , and v_3 are three independent Gaussian random variables with variance one and \mathbf{W}^{\dagger} is the transpose of \mathbf{W} . The actual $R^2(A_x, B_x)$ and $R^2(A_y, B_x)$ calculated from the generated {**B**} and {**A**} were 0.48 to 0.50.

4.2. Molecular Dynamics of Src SH2. In the simulation system, the natural phosphotyrosine (pY) residue was replaced by one with the main-chain amide nitrogen, C_{α} and C_{β} substituted by a cyclopropane moiety, which effectively constrains the side-chain conformation of the residue to that of the proteinbound state.^{18,19} The details of the MD simulations of the Src SH2 complex have been reported previously.¹⁹ Briefly, five sets of initial coordinates for the complex in explicit water were obtained from the multiple copies of the complex in the crystallographic asymmetric unit (PDB identifier 1IS0 and 1SPS). Eight simulations were initiated from each conformation by varying the initial velocities, yielding a total of 40 independent simulations. Each simulation was equilibrated for 500 ps and extended for 2 ns of production MD under constant temperature (298 K) and pressure (1 atm). Coordinates were saved at 1 ps interval from the production period.



Figure 5. ||VCC|| (left), GCC (middle), and DiCC (right) between C_a atoms of Src SH2 domains complexed with ligand pYEEI. Last four residues, starting from L1, are from ligands.



Figure 6. Correlation coefficient of the position vectors of five pairs of C_{α} atoms were calculated using 10 to 2000 ps long trajectory from each of the 40 molecular dynamics simulation of Src SH2 domain. Five pairs of C_{α} atoms are from residues 164 and 165 (red), residues 153 and 154 (orange), residues 206 and 215 (yellow), residues 206 and 216 (green), and residues 193 and 203 (blue). Standard deviation of correlation coefficients were calculated from 400 bootstrap samples. Panels a and d show absolute mean value and standard deviation of VCC of five C_{α} pairs respectively as a function of time. Panels b and e show mean value and standard deviation of GCC of the same pairs, respectively. Panels c and f show mean value and standard deviation of DiCC of the same pairs, respectively.

ASSOCIATED CONTENT

S Supporting Information

Detailed discussion of relation between generalized correlation coefficient and canonical correlation coefficient, correlation between components of position vectors calculated form MD simulation of Src SH2, and comparison of VCC and DiCC, VCC and GCC of pairs of C_{α} atoms calculated from the MD simulations. This material is available free of charge via the Internet at http://pubs.acs.org/.

AUTHOR INFORMATION

Corresponding Author

*E-mail: amitroy@purdue.edu.

Notes

The authors declare no competing financial interest.

ACKNOWLEDGMENTS

We thank the Rosen Center for Advanced Computing at Purdue University for providing computing resources. We thank Dr. Joshua M. Ward for allowing us to analyze Src SH2 complex trajectories from MD simulation. We also thank Jyotishka Dutta for his insightful discussion about R^2 . This work was supported by National Institutes of Health Grant No. AI039639.

Article

REFERENCES

(1) Berendsen, H. J. C.; Hayward, S. Collective protein dynamics in relation to function. *Curr. Opin. Struct. Biol.* **2000**, *10*, 165–169.

(2) Roy, A.; Post, C. B. Long-distance correlations of rhinovirus capsid dynamics contribute to uncoating and antiviral activity. *Proc. Natl. Acad. Sci. U.S.A.* **2012**, *109*, 5271–5276.

(3) Wereszczynski, J.; McCammon, J. A. Simulations of the p97 complex suggest novel conformational states of hydrolysis intermediates. *Protein Sci.* **2012**, *109*, 475–486.

(4) Tan, Y. S.; Fuentes, G.; Verma, C. A comparison of the dynamics of pantothenate synthetase from *M. tuberculosis* and *E. coli*: Computational studies. *Proteins* **2011**, *79*, 1715–1727.

(5) Mishra, S.; Caflisch, A. Dynamics in the active site of β -secretase: A network analysis of atomistic simulations. *Biochemistry* **2011**, *50*, 9328–9339.

Journal of Chemical Theory and Computation

(6) Ichiye, T.; Karplus, M. Collective motions in proteins: A covariance analysis of atomic fluctuations in molecular dynamics and normal mode simulations. *Proteins* **1991**, *11*, 205–217.

(7) Lange, F. O.; Grubmüller, H. Generalized correlation for biomolecular dynamics. *Proteins* **2006**, *62*, 1053–1061.

(8) Olbrich, C.; Strümpfer, J.; Schulten, K.; Kleinekathöfer, U. Quest for spatially correlated fluctuations in the FMO light-harvesting complex. J. Phys. Chem. B 2011, 115, 758–764.

(9) Kamberaj, H.; van der Vaart, A. Correlated motions and interactions at the onset of the DNA-induced partial unfolding of ets-1. *Biophys. J.* **2009**, *97*, 1747–1755.

(10) Székely, G. J.; Rizzo, M. L.; Bakirov, N. K. Measuring and testing dependence by correlation of distances. *Ann. Stat.* **2007**, *35*, 2769–2794.

(11) Székely, G. J.; Rizzo, M. L. Brownian distance covariance. Ann. Appl. Stat. 2009, 4, 1236–1265.

(12) Kraskov, A.; Stogbauer, H.; Grassberger, P. Estimating mutual information. *Phys. Rev. E* **2004**, *69*, 066138.

(13) MIxnyn. http://www.klab.caltech.edu/ kraskov/MILCA/ (accessed June 13, 2012).

(14) Rao, C. R. Linear Statistical Inference and Its Application; Wiley & Sons: New York, 1965; p 220.

(15) Nagelkerke, N. \hat{J} . D. A note on a general definition of the coefficient of determination. *Biometrika* **1991**, 78, 691–692.

(16) Reshef, D. N.; Reshef, Y. A.; Finucane, H. K.; Grossman, S. R.; McVean, G.; Turnbaugh, P. J.; Lander, E. S.; Mitzenmacher, M.; Sabeti, P. C. Detecting novel associations in large data sets. *Science* **2011**, *334*, 1518–1524.

(17) Joe, H. Relative entropy measures of multivariate dependence. J. Am. Stat. Assoc. 1989, 84, 157–164.

(18) Davidson, J. P.; Lubman, O.; Rose, T.; Waksman, G.; Martin, S. F. Calorimetric and structural studies of 1, 2, 3-trisubstituted cyclopropanes as conformationally constrained peptide inhibitors of Src SH2 domain binding. J. Am. Chem. Soc. 2002, 124, 205–215.

(19) Ward, J. M.; Gorenstein, N. M.; Tian, J.; Martin, S. F.; Post, C. B. Constraining binding hot spots: NMR and molecular dynamics simulations provide a structural explanation for enthalpy–entropy compensation in SH2–ligand binding. *J. Am. Chem. Soc.* **2010**, *132*, 11058–11070.

Supporting information for: Detection of long-range concerted motions in protein by a distance covariance

Amitava Roy* and Carol Beth Post

Medicinal Chemistry and Molecular Pharmacology, Purdue University, West Lafayette, USA

E-mail: amitroy@purdue.edu

1 Generalized correlation coefficient and canonical correlation

Let $\{A\}$ and $\{B\}$ be two *d*-dimensional Gaussian random vector series with zero mean. The total covariance matrix

$$\mathbf{C} = \begin{bmatrix} \mathbf{C}_{\mathbf{A}\mathbf{A}} & \mathbf{C}_{\mathbf{A}\mathbf{B}} \\ \mathbf{C}_{\mathbf{B}\mathbf{A}} & \mathbf{C}_{\mathbf{B}\mathbf{B}} \end{bmatrix}$$
(1)

is a block matrix where C_{AA} and C_{BB} are within-vector covariance matrices of $\{A\}$ and $\{B\}$ respectively and $C_{AB} = C_{BA}^{\dagger}$ is between-vector covariance matrix.

The mutual information (MI) between \mathbf{A} and \mathbf{B} is ¹

$$MI(\mathbf{A}, \mathbf{B}) = H(\mathbf{A}) + H(\mathbf{B}) - H(\mathbf{A}, \mathbf{B})$$
(2)

^{*}To whom correspondence should be addressed

where H is entropy and for a Gaussian distribution the entropy can be written as¹

$$H(\mathbf{A}) = \frac{1}{2} \ln(2\pi e)^{d} \parallel \mathbf{C}_{\mathbf{A}\mathbf{A}} \parallel$$
$$H(\mathbf{B}) = \frac{1}{2} \ln(2\pi e)^{d} \parallel \mathbf{C}_{\mathbf{B}\mathbf{B}} \parallel$$
$$H(\mathbf{A}, \mathbf{B}) = \frac{1}{2} \ln(2\pi e)^{2d} \parallel \mathbf{C} \parallel$$
(3)

with $\| . \|$ denoting the determinant. From Eq. (3) the MI between **A** and **B** is

$$MI(\mathbf{A}, \mathbf{B}) = -\frac{1}{2} \ln \left(\frac{\|\mathbf{C}\|}{\|\mathbf{C}_{\mathbf{A}\mathbf{A}}\| \| \mathbf{C}_{\mathbf{B}\mathbf{B}} \|} \right)$$
(4)

The generalized correlation coefficient, GCC,² of **A** and **B** is then

$$GCC(\mathbf{A}, \mathbf{B}) = \sqrt{1 - e^{\frac{-2Ml}{d}}}$$
$$= \sqrt{1 - \left(\frac{\|\mathbf{C}\|}{\|\mathbf{C}_{\mathbf{A}\mathbf{A}}\|\|\|\mathbf{C}_{\mathbf{B}\mathbf{B}}\|}\right)^{\frac{1}{d}}}$$
(5)

One can write

$$\frac{\|\mathbf{C}\|}{\|\mathbf{C}_{AA}\|\|\mathbf{C}_{BB}\|}$$

$$= \frac{\|\mathbf{C}_{AA}\|\|\mathbf{C}_{BB} - \mathbf{C}_{BA}\mathbf{C}_{AA}^{-1}\mathbf{C}_{AB}\|}{\|\mathbf{C}_{AA}\|\|\mathbf{C}_{BB}\|}$$

$$= \|\mathbf{I} - \mathbf{C}_{BB}^{-1}\mathbf{C}_{BA}\mathbf{C}_{AA}^{-1}\mathbf{C}_{AB}\|$$

$$= \prod_{i} (1 - \lambda_{i}^{2}).$$
Hence $GCC(\mathbf{A}, \mathbf{B}) = \sqrt{1 - \left(\prod_{i} (1 - \lambda_{i}^{2})\right)^{\frac{1}{d}}},$ (6)

where the λ_i^2 are eigenvalues of $\mathbf{C_{BB}}^{-1}\mathbf{C_{BA}}\mathbf{C_{AA}}^{-1}\mathbf{C_{AB}}$ and λ_i are called canonical correlations³ between **A** and **B**. λ_i^2 have values between zero and one. See Ref.⁴ for a short review on canonical correlations. The values λ_i^2 's are also the eigenvalues of $\mathbf{C_{AA}}^{-1}\mathbf{C_{AB}}\mathbf{C_{BB}}^{-1}\mathbf{C_{BA}}$. The eigenvectors

of the former matrix are the basis vectors for **B** and of the later matrix the basis vectors for **A**. For reference if **C**_{BB}, **C**_{AA} and **C**_{AB} are all diagonal then **C**_{BB}⁻¹**C**_{BA}**C**_{AA}⁻¹**C**_{AB} becomes **R**² matrix introduced in the main article. Let us assume λ_i^2 values are ordered, \mathbf{V}_i^A , \mathbf{V}_i^B are the corresponding eigenvectors and λ_1^2 is the largest eigenvalue. Then \mathbf{V}_1^A and \mathbf{V}_1^B are linear combinations of $A_1, ..., A_d$ and $B_1, ..., B_d$ such that the correlation between them is maximum among all possible combinations of $A_1, ..., A_d$ and $B_1, ..., B_d$ and λ_1 is the Pearson's correlation coefficient (PCC) between them. Similarly \mathbf{V}_2^A and \mathbf{V}_2^B are linear combinations of $A_1, ..., A_d$ and $B_1, ..., B_d$, which have a PCC of zero with \mathbf{V}_1^A and \mathbf{V}_1^B respectively, with the second largest PCC λ_2 between them. While canonical correlations have a definite physical meaning, the value $\sqrt{1 - (\prod_i (1 - \lambda_i^2))^{\frac{1}{d}}}$ is always dominated by the largest λ_i^2 and cannot be considered a proper correlation coefficient (CC) between position vectors. Consider the case where λ_1^2 is close to 1 while other λ^{2^2} s are close to zero. Then the product term will be close to zero and GCC will be close to 1. In such a case, although **A** and **B** are highly correlated in one direction, they are not correlated at all in the other d - 1 directions, yet the GCC value is still 1.

For example, if we take the series $\{A\}$ and $\{B\}$ generated from the model in Figure 1 with $\Delta \theta = \pi/6$ we have

$$\mathbf{C}_{\mathbf{A}\mathbf{A}} = \begin{bmatrix} 18.4 & 17.3 \\ 17.3 & 18.4 \end{bmatrix}, \mathbf{C}_{\mathbf{A}\mathbf{B}} = \begin{bmatrix} 7.4 & 24.2 \\ 5.7 & 24.7 \end{bmatrix}, \\ \mathbf{C}_{\mathbf{B}\mathbf{B}} = \begin{bmatrix} 5.0 & 12.6 \\ 12.6 & 48.5 \end{bmatrix}.$$
(7)

And we have $\lambda_1^2 = 0.90$ and $\lambda_2^2 = 0.68$, so that

$$\mathbf{V}_{1}^{A} = \begin{bmatrix} -0.64\\ 0.77 \end{bmatrix}, \mathbf{V}_{2}^{A} = \begin{bmatrix} -0.77\\ -0.64 \end{bmatrix},$$
$$\mathbf{V}_{1}^{B} = \begin{bmatrix} -0.91\\ 0.41 \end{bmatrix}, \mathbf{V}_{2}^{B} = \begin{bmatrix} 0.68\\ -0.41 \end{bmatrix}.$$
(8)

 \mathbf{V}_1^A and \mathbf{V}_1^B are almost perpendicular to **A** and **B** respectively. So λ_1^2 is reflecting $R^2(B_{\theta}, A_{\theta})$ which is 1 in the model in Figure 1. Similarly \mathbf{V}_2^A and \mathbf{V}_2^B are almost parallel to **A** and **B** respectively and λ_2^2 reflects $R^2(B_r, A_r)$ which is 0.69. And GCC of (**B**, **A**) becomes 0.91. GCC calculated by methods developed by Kraskov et al.⁵ gave a value of 0.96 as plotted in Figure 1, which is not an accurate reflection of the correlation in {**A**} and {**B**} or B_r and A_r .

2 Correlation between components of vector



Figure S1: Pearson's correlation coefficient between \hat{x} , \hat{y} and \hat{z} of the position vectors of C_{α} atoms calculated from the molecular dynamics simulation of Src SH2 domain. About 40% of the calculated PCC values are more than 0.3 and about 10% are more than 0.6.

We calculated Pearson's correlation coefficient (PCC) between \hat{x} , \hat{y} and \hat{z} components of the position vectors of C_{α} atoms in the 80,000 conformations saved during the molecular dynamics simulation of Src SH2 domain described in the main article. Figure S1 shows the probability

density of these PCC values. About 40% of the calculated PCC values are more than 0.3 and about 10% are more than 0.6.



3 Comparison among correlation coefficients

Figure S2: (a) DiCC and ||VCC|| of pairs of C_{α} atoms calculated from 40x2 ns long trajectory of Src SH2 domain in complex with the ligand pYEEI. (b) GCC and ||VCC|| of pairs of C_{α} atoms.

References

- Cover T.; Thomas J. Elements of Information Theory; Wiley & Sons: New York, 1991; p 230.
- (2) Lange F.O.; Grubmüller H. Generalized Correlation for Biomolecular Dynamics. *Proteins* 2006, 62, 1053–1061.
- (3) Hotelling H. The most predictable criterion. J. Edu. Psych. 1935, 26, 139–142.
- (4) Kettenring J.R. Canonical analysis of several sets of variables. *Biometrika* 1971, 58, 433–451.

(5) Kraskov A.; Stogbauer H.; Grassberger P. Estimating mutual information. *Phys. Rev. E* 2004, 69, 066138.