

# Relative Principal Components Analysis: Application to Analyzing Biomolecular Conformational Changes

*Mazen Ahmad<sup>a\*</sup>, Volkhard Helms<sup>b</sup>, Olga V. Kalinina<sup>a</sup> and Thomas Lengauer<sup>a</sup>*

<sup>a</sup>Computational Biology Research Group, Max Planck Institute for Informatics, Saarland Informatics Campus, Campus E1 4, 66123 Saarbrücken, Germany.

<sup>b</sup>Center for Bioinformatics, Saarland University, 66123 Saarbrücken, Germany.

\*Correspondence author: mahmad@mpi-inf.mpg.de

## Abstract

A new method termed “Relative Principal Components analysis” (RPCA) is introduced that extracts optimal relevant principal components to describe the change between two data samples representing two macroscopic states. The method is widely applicable in data-driven science. Calculating the components is based on a unified physical framework which introduces the objective function, namely the Kullback-Leibler divergence, appropriate for quantifying the change of the macroscopic state as it is effected by the microscopic features. To demonstrate the applicability of RPCA, we analyze the thermodynamically relevant conformational changes of the protein HIV-1 protease upon binding to different drug molecules. In this case, the RPCA method provides a sound thermodynamic foundation for the analysis of the binding process. The relevant collective (global) conformational changes can be reconstructed from the informative latent variables to exhibit both the enhanced and the restricted conformational fluctuations upon ligand association. Moreover, RPCA characterizes the locally relevant conformational changes which can be presented on the structure of the protein.

## 1 Introduction

Studying the transitions and differences between multiple states populated by a dynamic system is a central topic in different fields including chemistry, physics, biology, machine learning and all of data-driven science. A typical task is to uncover how macroscopic changes of the dynamic system are related to the features (variables) that describe its microscopic individuals (instances). Two examples of such microscopic features would be the genetic sequences of a virus taken from snapshots during the course of evolution or the spatial conformations of two biomolecules when they bind to each other. The relationship between the “microscopic” factors of a system and the change of its macroscopic states requires the definition of an appropriate objective function for quantifying the change of the “macroscopic” state of the system. Such a rigorous definition of changes of the macroscopic state of a system in terms of its microscopic features is available for physical systems whose thermodynamic quantities can be measured or computed. For example, the change of free energy (a scalar value) is a suitable quantity to characterize macroscopic changes in physical, chemical and biochemical systems. However, in other areas of data-driven science, such a rigorous definition and quantification of macroscopic changes generally does not exist. Instead, various heuristic objective functions are used in practice. Examples include divergence measures from information theory (Lin, 1991) and the wide variety of objective functions which are used for prediction and feature extraction in pattern recognition (Fukunaga, 1990). Mining the factors informative for the change between two samples is of high importance and of general interest in all areas of data-driven science and is generally performed in a high-dimensional feature space. In fact, mining informative features is the central theme in a large domain of machine learning and includes methods such as dimensionality reduction (Murphy, 2012; Theodoridis, 2015), feature extraction (Fukunaga, 1990), and latent variable models (Bishop, 1998). However, one needs to select an objective function that is appropriate for quantifying the change before applying a multivariate method to extract the informative features.

Analyzing the conformational changes taking place during biomolecular reactions is one of the most important tasks in structural biology. Unfortunately, analyzing and mechanistically understanding biochemical interactions is quite tricky due to the complex conformational dynamics in the high-dimensional space where the interactions take place. The macroscopic changes in biochemical systems, on the other hand, are quantified using the change of free energy (a scalar quantity). Molecular simulations are becoming a more and more attractive tool for analyzing conformational changes of biomolecules. Current interest in analyzing molecular conformational changes (e.g. using Markov state models (Husic and Pande, 2018)) focusses on characterizing the kinetic changes between representative conformations within one macroscopic state and the analysis is performed in a data space where the points are the (clustered) conformations. Methods such as PCA (Amadei et al., 1993; Kitao and Go, 1999) and partial least squares (Krivobokova et al., 2012) were used to analyze the conformational changes in the feature space (coordinates). However, most of these methods are only suited for studying the dynamic changes within one macroscopic state and not for studying transitions and changes between two states. On the other hand, methods adapted from multivariate analysis often lack thermodynamic insight.

In this work, we introduce a unified framework rooted in statistical information theory and statistical mechanics (Kullback, 1997a; Barndorff-Nielsen, 1978; Ahmad et al., 2015a, 2017) for the purpose of studying the change between two data sets representing two states. A new method termed “*Relative Principal Components Analysis*” is introduced to extract the optimal relevant principal components describing the change between two data samples (two states). This method is applicable in all areas of data-driven science and is introduced in its generality in Section 2 and 3. As special but important example, starting in Section 4, we apply RPCA for analyzing the energetically relevant conformational changes of a biomolecule (protease of the human immunodeficiency virus, HIV-1) upon binding to various ligands.

## 2 A general formalism for analyzing changes in dynamic systems based on information theory

Before going into the technical details of finding the directions in feature space that are informative of the change between two states, we first introduce a physical framework for defining and quantifying the change of dynamic systems in all areas of data-driven sciences and justify the objective function used for quantifying the macroscopic change.

Let  $\mathbf{x}$  be a random vector encoding the microscopic features (variables) which are necessary to identify the difference between the instances of a system of interest. By a *system* we mean a collection of individuals or *instances*, e.g. viruses or molecules, where each instance is defined via a set of (microscopic) *features*. A macroscopic state ( $a$ ) of the system is defined by the probability density distribution  $P(\mathbf{x}|a) = P_a(\mathbf{x})$  of all possible instances ( $\mathbf{x}$ ) when the system is at equilibrium in this macroscopic state  $a$ . The macroscopic state  $a$  can be changed into a new macroscopic state  $b$  by perturbing the probability distribution of the microscopic instances  $\mathbf{x}$  to  $P(\mathbf{x}|b) = P_b(\mathbf{x})$ . Here, a Bayesian approach is used whereby the macroscopic state is viewed as a random variable, and the conditional probability  $P(\mathbf{x}|i) = P_i(\mathbf{x})$  is naturally interpreted as the equilibrium distribution of state  $i$  that can be taken from a finite or discrete set. A relationship between the two macroscopic states can be obtained by applying Bayes’ theorem, yielding the following equation for the probability density ratio  $P_b(\mathbf{x})/P_a(\mathbf{x})$ :

$$\frac{P(\mathbf{x}|b)}{P(\mathbf{x}|a)} = \frac{P_b(\mathbf{x})}{P_a(\mathbf{x})} = \frac{P(a)}{P(b)} \frac{P(b|\mathbf{x})}{P(a|\mathbf{x})} \quad (2.1)$$

Here,  $P(b)$  and  $P(a)$  are conditional probabilities under the (implicit) condition of the perturbing factors. Averaging the logarithm of the density ratio equation over the probability distribution of state  $b$  ( $P_b(\mathbf{x})$ ) yields:

$$\underbrace{\ln \left[ \frac{P(a)}{P(b)} \right]}_{\Delta F} = \int P_b(\mathbf{x}) \underbrace{\ln \left[ \frac{P(a|\mathbf{x})}{P(b|\mathbf{x})} \right]}_{U_p(\mathbf{x})} d\mathbf{x} + \underbrace{\int P_b(\mathbf{x}) \ln \left[ \frac{P_b(\mathbf{x})}{P_a(\mathbf{x})} \right] d\mathbf{x}}_{D_{kl}} \quad (2.2)$$

Indeed, this is a general derivation of the Perturbation Divergence Formalism (PDF), which was previously derived for physical systems for the purpose of decomposing the change of free energy ( $\Delta F$ ) between two macroscopic states  $a$  and  $b$  (Ahmad et al., 2017, 2015a):

$$\ln \left[ \frac{P(a)}{P(b)} \right] = \Delta F = \langle U_p(\mathbf{x}) \rangle_b + D_{kl}(P_b(\mathbf{x}) \parallel P_a(\mathbf{x})) \quad (2.3)$$

For physical systems, we use here the natural unit of the energy  $(kT)^{-1} = 1$ .  $D_{kl}$  is the Kullback-Leibler (KL) divergence (Kullback and Leibler, 1951) (also termed the relative entropy (Cover and Thomas, 2006) or the discriminant information (Kullback, 1997a)) between the probability distributions of states  $a$  and  $b$ . Interestingly, equation (2.2) provides a purely probabilistic formalism of free energy change via a new probabilistic definition of the perturbation  $U_p$  of the microscopic instances as the logarithm of the ratio of the posterior probabilities of two macroscopic states given a particular microscopic configuration  $\mathbf{x}$ :

$$U_p(\mathbf{x}) = \ln \left[ \frac{P(a|\mathbf{x})}{P(b|\mathbf{x})} \right]. \quad (2.4)$$

**Bridging the macroscopic change and the microscopic features.** The concept of the perturbation of a microscopic instance was originally introduced in statistical mechanics as an energetic quantity in order to formalize the relationship between the microscopic (atomistic) description of a physical system and its macroscopic changes between states (free energy change)(Kirkwood, 1934; Zwanzig, 1954). When considering the change between two macroscopic states, the perturbation  $U_p(\mathbf{x})$  is a one-dimensional (microscopic) variable which has the same KL divergence (discriminant information) as the high dimensional feature  $\mathbf{x}$  (Ahmad et al., 2015a) :

$$D_{kl}(P_b(\mathbf{x}) \parallel P_a(\mathbf{x})) = D_{kl}(P_b(U_p(\mathbf{x})) \parallel P_a(U_p(\mathbf{x}))). \quad (2.5)$$

In statistical inference theory (Berger and Casella, 2001; Lehmann and Casella, 2003), such a variable is termed a sufficient statistic. A sound framework for the relationship between the change of macroscopic states of a dynamic system and its microscopic elements can be inherited from the formalism of exponential families (Barndorff-Nielsen, 1978; Brown, 1986) in statistical estimation theory.

Let us assume we have a dynamic system at equilibrium in a macroscopic reference state labeled by the so-called *natural parameters*  $\lambda$ . The system can populate a new macroscopic state  $P_\lambda(\mathbf{x})$  upon perturbing its original state  $P_{\lambda_0}(\mathbf{x})$ . The principle of minimum discrimination information (Kullback, 1997b) (equivalent to the principle of maximum entropy (Grevin et al., 2003)) is applied to find the new distribution  $P_\lambda(\mathbf{x})$  via minimizing the KL divergence from the reference distribution  $P_{\lambda_0}(\mathbf{x})$  to the new distribution  $P_\lambda(\mathbf{x})$  under the constraint of a finite expected value of the sufficient statistic  $\langle \mathbf{T}(\mathbf{x}) \rangle_\lambda$ . The probability density distribution  $P_\lambda(\mathbf{x})$  forms an exponential family of distributions in terms of the reference distribution:

$$P_\lambda(\mathbf{x}) = \exp[\lambda^T \mathbf{T}(\mathbf{x}) - \psi(\lambda)] P_{\lambda_0}(\mathbf{x}) \quad (2.6)$$

$$\psi(\boldsymbol{\lambda}) = \ln \left[ \int_{\boldsymbol{x}} d\boldsymbol{x} P_{\lambda_0}(\boldsymbol{x}) \exp(\boldsymbol{\lambda}^T \boldsymbol{T}(\boldsymbol{x})) \right]$$

$$D_{kl}(P_{\boldsymbol{\lambda}}(\boldsymbol{x}) \parallel P_{\lambda_0}(\boldsymbol{x}); \langle \boldsymbol{T}(\boldsymbol{x}) \rangle_{\boldsymbol{\lambda}}) = \boldsymbol{\lambda}^T \langle \boldsymbol{T}(\boldsymbol{x}) \rangle_{\boldsymbol{\lambda}} - \psi(\boldsymbol{\lambda}). \quad (2.7)$$

Here, the sufficient statistic  $\boldsymbol{T}(\boldsymbol{x})$  is a vector (or a scalar) function of the microscopic configuration  $\boldsymbol{x}$  of the system. The cumulant generating function  $\psi(\boldsymbol{\lambda})$  depends only on the natural parameters  $\boldsymbol{\lambda}$ . Besides being used to formulate a theoretical framework for studying the error bound of parameter estimation (Kullback, 1997b), the formalism of exponential families plays a central role in different fields of machine learning such as generalized linear models and variational inference (Murphy, 2012). In statistical thermodynamics, Kirkwood introduced his thermodynamic integration (TI) equation (Kirkwood, 1934) using the exponential family to “*alchemically*” interpolate two macroscopic states (Kirkwood, 1934). Indeed, the one-dimensional sufficient statistic, “the perturbation”, is the appropriate tool for interpolating between two macroscopic states in free energy calculations. Generally, a higher-dimensional sufficient statistic is required when studying multiple macroscopic states (Berger and Casella, 2001). The sufficient statistic and the corresponding cumulant generating function in equations (2.6) are not unique. Notably, comparing equations (2.2) and (2.7) shows that the sufficient statistic can be selected such that the corresponding cumulant generating function is a function of the change of free energy between the macroscopic states of dynamical systems. For example, when selecting the negative value of the perturbation ( $-U_p(\boldsymbol{x})$ ) as a sufficient statistic for studying the change between two macroscopic states, the corresponding cumulant generating function is the negative of the free energy change and equation (2.6) turns into the free energy perturbation equation that is well known in statistical thermodynamics (Zwanzig, 1954). Another example is the logarithm of the density ratio, which is a well-known sufficient statistic in statistical inference, and its corresponding cumulant generating function is the relative change of free energy (see Supplementary Material S4). An interesting known relationship of an exponential family is the functional dependence of the change of  $\psi(\boldsymbol{\lambda})$  (the function of the macroscopic state) on the microscopic features  $\boldsymbol{x}$  that is given by the average and the covariance of the sufficient statistic:

$$\nabla \psi(\boldsymbol{\lambda}) = \langle \boldsymbol{T}(\boldsymbol{x}) \rangle_{\boldsymbol{\lambda}}$$

$$\nabla^2 \psi(\boldsymbol{\lambda}) = \text{Cov}_{\boldsymbol{\lambda}}(\boldsymbol{T}(\boldsymbol{x})). \quad (2.8)$$

The free energy thermodynamic integration equation (Kirkwood, 1934) is a special case of equation (2.8). However, quantification of the macroscopic change is not of general interest in the field of data-driven sciences. The important task here is identifying the unknown sufficient statistic that explains the influence of the microscopic features of the macroscopic change. Unlike free energy change, KL divergence is a quantity familiar in machine learning and can be computed using parametric or nonparametric models (Sugiyama et al., 2012). Indeed, equation (2.7) shows that KL divergence is a Legendre transformation (Barndorff-Nielsen, 1978; Rockafellar, 1970) of the cumulant generating function and can be used to quantify the change of the macroscopic state. The PDF given by equation (2.2) is a special case of equation (2.7) where we use the perturbation as a sufficient statistic for studying the change between two

macroscopic states ( $\lambda = 0,1$ ). Due to the Legendre duality (Barndorff-Nielsen, 1978; Rockafellar, 1970) between the KL divergence and the change of the free energy, the relevance of the PDF goes beyond being a decomposition of the change of free energy. In fact, the terms of equation (2.2) are the perturbations (the features informative of the change) and their fingerprints (the configurational changes) which are quantified by the KL divergence. In this view, significant perturbations are reflected by significant changes of KL divergence.

### 3 Relative principal components analysis

This section presents a new method for analyzing the change between two states (data sets) using multivariate analysis methods (Theodoridis, 2015). Studying the change between multiple states is beyond the scope of this work and will be presented in a future publication. The newly introduced method termed “relative principal components analysis” (RPCA) computes collective canonical variables (linear combinations of the original features) termed the relative principle components (RPCAs) to which the KL divergence factorizes additively. Indeed, factorization of the KL divergence is equivalent to factorization of the logarithm of the logarithm of the density ratio which is a sufficient statistic of interest in machine learning (Sugiyama et al., 2012) (see above). Factorization of KL divergence was introduced for multivariate normal distributions in the seminal work of Kullback (Kullback, 1997a). However, the theoretical approach of factorizing KL divergence as introduced by Kullback was not accessible in practice. Specifically, a solution is needed around the singularity of the covariance matrices and the resulting features have to be optimal with respect to maximizing their KL divergences (see below).

Let  $\mathbf{x} = (x_1 \dots x_d)^T$  be a  $d$ -dimensional continuous random variable with two samples from two macroscopic states that will be labeled ( $a$ ) and ( $b$ ). RPCA aims at finding  $k$  latent canonical variables  $\mathbf{y} = (y_1, y_2 \dots y_k)^T = f(\mathbf{x})$  which satisfy the following two conditions: (i) their marginal distributions are independent in the states  $a$  and  $b$  meaning that their KL divergences are additive, and (ii) optimal in terms of maximizing the KL divergence, such that we can use  $m$  ( $m \ll d$ ) latent variables to represent the significant directions informative of the change. The KL divergence of a new variable  $\mathbf{y} = f(\mathbf{x})$  is always non-negative (Kullback, 1997a; Cover and Thomas, 2006) and is bounded from above by the KL divergence of the original variable  $\mathbf{x}$  (Ahmad et al., 2017; Kullback, 1997a):

$$D_{kl}(P_b(\mathbf{x}) \parallel P_a(\mathbf{x})) \geq D_{kl}(P_b(f(\mathbf{x})) \parallel P_a(f(\mathbf{x}))) \geq 0 \quad (3.1)$$

The new variable  $\mathbf{y} = f(\mathbf{x})$  is “sufficient” if equality holds in (3.1). When studying the change between two macroscopic states, a sufficient one-dimensional variable always exists (e.g. the perturbation (Ahmad et al., 2017) ) regardless of the dimensionality of  $\mathbf{x}$ . Although the existence of a one-dimensional sufficient feature appears promising, it is not practically useful for two reasons: (i) the analytical nature of the sufficient one-dimensional variable is generally unknown, and (ii) the complexity and nonlinearity of a sufficient one-dimensional variable –if it is known– will hinder its interpretability in terms of the original features  $\mathbf{x}$ . In fact, its simple interpretability and analytical traceability is one reason for the widespread use of the Gaussian linear parametric model in latent variable models (Bishop, 1998) (e.g. PCA).

We will adopt this model here as well. Then the latent variables are linear combinations of the original variable and are normally distributed:

$$\begin{aligned} \mathbf{y} &= f(\mathbf{x}) = \mathbf{G}^T \mathbf{x} \\ \mathbf{y}_l &\sim \mathcal{N}(\boldsymbol{\mu}_l, \boldsymbol{\Lambda}_l) ; l = a, b. \end{aligned} \quad (3.2)$$

Here,  $\mathbf{G}$  is a  $d \times k$  transformation matrix with columns  $\mathbf{g}_i$ . Thus  $y_i = \mathbf{g}_i^T \mathbf{x}$  and  $\boldsymbol{\mu}_l ; l = a, b$  are the averages of the two distributions. The covariance matrices  $\mathbf{S}_l$  of the original variables are related to the covariance matrices of the latent variables by  $\boldsymbol{\Lambda}_l = \mathbf{G}^T \mathbf{S}_l \mathbf{G} ; l = a, b$ . Under the model assumption of normality, the independence of the variables  $\mathbf{y}_l$  requires  $\boldsymbol{\Lambda}_l$  to be diagonal in both states  $a$  and  $b$ :

$$\begin{aligned} \mathbf{G}^T \mathbf{S}_a \mathbf{G} &= \mathbf{I} \\ \mathbf{G}^T \mathbf{S}_b \mathbf{G} &= \boldsymbol{\Lambda} \end{aligned} \quad (3.3)$$

Here, the diagonal matrix of state  $b$ ,  $\boldsymbol{\Lambda} = \text{diag}(\lambda_1 \dots \lambda_k)$ , contains the variances  $\lambda_1 \dots \lambda_k$  of the latent variables at state  $b$  while the covariance matrix of the latent variables at state  $a$  is arbitrarily selected to be an identity matrix ( $\mathbf{I}$ ). In case of a multivariate normal distribution ( $\mathbf{S}_a$  is nonsingular), Kullback (Kullback, 1997a) formulated the solution for  $\mathbf{G}$  as the generalized eigenvectors corresponding to the generalized eigenvalue problem  $|\mathbf{S}_b - \lambda \mathbf{S}_a|$  which can be solved using Wilkinson's algorithm (Martin and Wilkinson, 1971) involving the Cholesky decomposition of  $\mathbf{S}_a$  (Stewart, 2001, p. 229). Practically, the covariance matrices of real data are mostly singular or ill-conditioned and the generalized eigenproblem is not solvable. Singularity arises due to the fact that the real dimensionality of the probability distributions is smaller than the apparent dimensionality of  $\mathbf{x}$  (Mardia et al., 1979, p. 41).

**RPCA via simultaneous diagonalization of two matrices.** Here, we present a general algorithm for the simultaneous diagonalization of the matrices in equation (3.3) which can be used even if  $\mathbf{S}_a$  is singular.

A transformation matrix  $\mathbf{G} \in \mathbb{R}^{d \times k}$  that simultaneously diagonalizes the symmetric matrices  $\mathbf{S}_a \in \mathbb{R}^{d \times d}$  (of rank  $k$ ) and  $\mathbf{S}_b \in \mathbb{R}^{d \times d}$  (equation (3.3)) can be found by a combination of two transformation matrices:

$$\mathbf{G} = \mathbf{W} \boldsymbol{\Psi}. \quad (3.4)$$

- 1) The so-called *whitening transformation* (Fukunaga, 1990) matrix  $\mathbf{W} \in \mathbb{R}^{d \times k}$  of the matrix  $\mathbf{S}_a$  is computed from its eigendecomposition  $\mathbf{S}_a = \mathbf{U} \mathbf{D} \mathbf{U}^T = \mathbf{U}_k \mathbf{D}_k \mathbf{U}_k^T$ :

$$\mathbf{W} = \mathbf{U}_k \mathbf{D}_k^{-1/2}. \quad (3.5)$$

Here, the  $k$  eigenvectors in the columns of  $\mathbf{U}_k$  correspond to the  $k$  nonzero eigenvalues in the diagonal matrix  $\mathbf{D}_k$ . Clearly,  $\mathbf{W}$  reduces  $\mathbf{S}_a$  to an identity matrix  $\mathbf{W}^T \mathbf{S}_a \mathbf{W} = \mathbf{I} \in \mathbb{R}^{k \times k}$  corresponding to the covariance matrix of the whitened data ( $\mathbf{W}^T \mathbf{x}$ ). The algorithm above is well known in case  $\mathbf{S}_a$  is nonsingular ( $d = k$ ), (Fukunaga, 1990).

- 2) The matrix  $\boldsymbol{\Psi}$  is formed using the eigenvectors from the symmetric matrix  $\mathbf{W}^T \mathbf{S}_b \mathbf{W}$ :

$$\mathbf{W}^T \mathbf{S}_b \mathbf{W} = \mathbf{\Psi} \mathbf{\Lambda} \mathbf{\Psi}^T. \quad (3.6)$$

It is straightforward to see from (3.5) and (3.6) that  $\mathbf{G} = \mathbf{W} \mathbf{\Psi}$  simultaneously diagonalizes  $\mathbf{S}_a$  and  $\mathbf{S}_b$  and satisfies equation (3.3).

The  $k$  columns of  $\mathbf{G}$  are the generalized eigenvectors with the corresponding generalized eigenvalues in the diagonal matrix  $\mathbf{\Lambda}$ . The relative principle components  $y_i = \mathbf{g}_i^T \mathbf{x}$  can be reordered based on their KL divergences. The additive KL divergences of the independent variables  $y_i$  can be analytically computed based on the model assumption of normality (Kullback, 1997a):

$$D_{kl}(\mathbf{y}; a: b) = \sum_{i=1}^k D_{kl}(y_i = \mathbf{g}_i^T \mathbf{x}; a: b) = \sum_{i=1}^k \frac{1}{2} \left[ \underbrace{-\ln \lambda_i + \lambda_i - 1}_{\text{Variance change}} + \underbrace{\frac{\mathbf{g}_i^T \Delta \Delta^T \mathbf{g}_i}{\text{Average change}}}_{\text{Average change}} \right] \quad (3.7)$$

Here  $\Delta = \boldsymbol{\mu}_b - \boldsymbol{\mu}_a$  is the change of the average of the distributions of  $\mathbf{x}$  in states  $a$  and  $b$ . However, it is important to keep in mind that the value of  $D_{KL}$  of the components from equation (3.7) is computed based on the model assumption (normality). Fortunately, a significant KL divergence of a variable  $y_i$  has to be reflected in a significant change between its distributions in states  $a$  and  $b$  which can be used to assess the accuracy of the model-based KL divergences (see the example below).

**Optimal RPCAs via average-covariance sub-spacing.** Although the simultaneous diagonalization algorithm above returns independent latent variables into which the KL divergence factorizes additively, the obtained latent variables are not optimal in terms of maximizing the KL divergence. Optimal latent variables are required for reconstructing the most informative approximation (maximizing KL divergence) of the original variable; see below. The KL divergences of the relative principle components in equation (3.7) can be decomposed into the terms holding the change of the variances  $\frac{1}{2}[-\ln \lambda_i + \lambda_i - 1]$  and the terms holding the change of the average  $\frac{1}{2}(\mathbf{g}_i^T \Delta \Delta^T \mathbf{g}_i)$ . Indeed, the latent variables from equation (3.4) are optimal with respect to maximizing the KL divergences due to change of the variances (Kullback, 1997a). Unfortunately, this optimality is violated by the contributions to the KL divergences due to the change of the average  $\frac{1}{2}(\mathbf{g}_i^T \Delta \Delta^T \mathbf{g}_i)$ . Therefore, the following average-covariance sub-spacing algorithm is introduced to achieve the optimality of RPCAs with respect to maximizing their KL divergences. The detailed derivation is presented in Supplementary Material S1. The idea is to find a transformation matrix  $\mathbf{G} = [\mathbf{g}_\mu \quad \mathbf{g}_v]$  such that the KL divergence of the variable  $y_\mu = \mathbf{g}_\mu^T \mathbf{x}$  summarizes the total KL divergence due to changes of the averages while the KL divergences of the variables  $\mathbf{y}_v = \mathbf{G}_v^T \mathbf{x}$  are purely due to the change of covariance. The required solution for  $\mathbf{g}_\mu$  is given by:

$$\mathbf{g}_\mu = \frac{\mathbf{g}_\mu^*}{\sqrt{\mathbf{g}_\mu^{*T} \mathbf{S}_a \mathbf{g}_\mu^*}}; \quad \mathbf{g}_\mu^* = \mathbf{S}_a^{-1} \Delta \quad (3.8)$$

$$\mathbf{S}_a^- = \mathbf{U}_k \mathbf{D}_k^{-1} \mathbf{U}_k^T.$$

Here, the generalized pseudoinverse  $\mathbf{S}_a^-$  is used for the general case (e.g. singular  $\mathbf{S}_a$ ) and  $\mathbf{g}_\mu$  is normalized with respect to the covariance matrix  $\mathbf{S}_a$  such that  $\mathbf{g}_\mu^T \mathbf{S}_a \mathbf{g}_\mu = 1$ . The KL divergence of the variable  $y_\mu = \mathbf{g}_\mu^T \mathbf{x}$  includes the total KL divergence due to the change of the average, which in turn equals half of the Mahalanobis distance between the averages  $\frac{1}{2} \Delta^T \mathbf{S}_a^- \Delta$ :

$$D_{\text{KL}}(y_\mu = \mathbf{g}_\mu^T \mathbf{x}; a: b) = \frac{1}{2} \left[ \underbrace{-\ln \lambda_\mu + \lambda_\mu - 1}_{\text{Variance change}} + \frac{\Delta^T \mathbf{S}_a^- \Delta}{\mathbf{g}_\mu^T \Delta \Delta^T \mathbf{g}_\mu} \right] \quad (3.9)$$

$$\text{Var}(y_\mu | b) = \lambda_\mu = \mathbf{g}_\mu^T \mathbf{S}_b \mathbf{g}_\mu$$

Here, the KL divergence of the new variable  $y_\mu = \mathbf{g}_\mu^T \mathbf{x}$  also includes a contribution due to the change of its variance ( $\lambda_\mu$ ).

The remaining generalized eigenvectors  $\mathbf{G}_v$ , which do not contain contributions arising from the change of the average ( $\mathbf{g}_{i \neq \mu}^T \Delta = 0$ ), are computed after deflating the contributions to the KL divergence due to vector  $\mathbf{g}_\mu$  from the matrices  $\mathbf{S}_a$  and  $\mathbf{S}_b$ . Given the vector  $\mathbf{g}_\mu$ , the restricted simultaneous diagonalization problem can be presented in the equations:

$$\begin{aligned} \mathbf{G}^T \mathbf{S}_a \mathbf{G} &= \begin{bmatrix} \mathbf{g}_\mu^T \\ \mathbf{G}_v^T \end{bmatrix} \mathbf{S}_a [\mathbf{g}_\mu \quad \mathbf{G}_v] = \mathbf{I} \\ \mathbf{G}^T \mathbf{S}_b \mathbf{G} &= \begin{bmatrix} \mathbf{g}_\mu^T \\ \mathbf{G}_v^T \end{bmatrix} \mathbf{S}_b [\mathbf{g}_\mu \quad \mathbf{G}_v] = \begin{bmatrix} \lambda_\mu & \mathbf{0} \\ \mathbf{0} & \mathbf{\Lambda}_v \end{bmatrix} \end{aligned} \quad (3.10)$$

$$\text{subject to } \mathbf{G}_v^T \Delta = \mathbf{0} \quad (\mathbf{g}_{i \neq \mu}^T \Delta = 0 \text{ ; no KL divergence due to the average change}). \quad (3.11)$$

Here,  $\mathbf{\Lambda}_v$  is a diagonal matrix and  $\mathbf{0}$  denotes a matrix or a vector of zeros of suitable dimensionality. To fulfill the conditions in (3.10) and (3.11), the generalized eigenvectors  $\mathbf{G}_v$  can be constructed using a combination of two transformation matrices (for details see Supplementary Material S1):

$$\mathbf{G}_v = \mathbf{W}_v \mathbf{\Psi}_v \quad (3.12)$$

(i) The whitening transformation matrix  $\mathbf{W}_v$ , similarly to equation (3.5), is constructed here from the eigendecomposition of the covariance matrix of the projection of  $\mathbf{x}$  on the subspace, which is orthogonal to the vector  $(\mathbf{S}_a \mathbf{g}_\mu)$ . (ii) The second transformation  $\mathbf{\Psi}_v$  is obtained from the eigenvectors of covariance matrix of the projection of the whitened data ( $\mathbf{W}_v^T \mathbf{x}$ ) onto the subspace which is orthogonal to the vector  $\mathbf{W}_v^T \mathbf{S}_b \mathbf{g}_\mu$ . Here, the number of generalized eigenvectors from the optimal RPCA sub-spacing is one less than the number of generalized eigenvectors from the non-optimal RPCA algorithm of equation (3.4).

**Data reconstruction from the latent variable.** The influence of the  $m$ -dimensional ( $m < d$ ) latent variable  $\mathbf{y}_m = \mathbf{G}_m^T \mathbf{x}$  on the original  $d$ -dimensional variable  $\mathbf{x} \in \mathbb{R}^d$  can be presented via the

reconstruction (projection)  $\hat{\mathbf{x}} \in \mathbb{R}^d$  in the subspace which is spanned by the corresponding  $m$  generalized eigenvectors  $\mathbf{G}_m = [\mathbf{g}_1 \dots \mathbf{g}_m]$  and is given by the relationship (Harville, 2008):

$$\hat{\mathbf{x}} = \underbrace{\mathbf{G}_m (\mathbf{G}_m^T \mathbf{G}_m)^{-1} \mathbf{G}_m^T}_{\mathbf{P}_m} \mathbf{x} = \mathbf{P}_m \mathbf{x} \quad (3.13)$$

$$\hat{\mathbf{x}} = \underbrace{\mathbf{G}_m (\mathbf{G}_m^T \mathbf{G}_m)^{-1} \mathbf{G}_m^T}_{\mathbf{R}} \mathbf{x} = \mathbf{R} \mathbf{y}_m \quad (3.14)$$

Here,  $\mathbf{P}_m \in \mathbb{R}^{d \times d}$  is the projection matrix (Harville, 2008) on the  $d$ -dimensional subspace which is spanned by columns of  $\mathbf{G}_m$ .  $\mathbf{R} \in \mathbb{R}^{d \times m}$  is the reconstruction matrix facilitating the transformation from the latent variable. It is straightforward (see Supplementary Material S2) to show that the KL divergence between the distributions of the reconstructed variable  $\hat{\mathbf{x}}$  equals the KL divergence between the distributions of its corresponding latent variable  $\mathbf{y}_m$ . In other words, the most informative approximation (the one maximizing KL divergence) of the original variable  $\mathbf{x}$  using a restricted number ( $m < d$ ) of variables is obtained using the  $m$ -dimensional latent variables  $\mathbf{y}_m$  with the highest KL divergences.

**Change hotspots from RPCA.** Besides representing the changes collectively (incorporating all  $x_i$ ), RPCA provides the possibility to map the feature-wise (local) contributions to the change from the individual elements  $x_i$  of  $\mathbf{x}$  such that the elements of  $\mathbf{x}$  with larger contribution to the change (KL divergence) can be interpreted as the hotspots of the change. The contributions to the divergence in equation (3.7) can be approximated as a sum of two quadratic terms:

$$D_{\text{KL}}(\mathbf{y}_l; a: b) \approx \frac{1}{2} [\lambda_l - 1 + \mathbf{g}_l^T \Delta \Delta^T \mathbf{g}_l]$$

$$\lambda_l = \mathbf{g}_l^T \mathbf{S}_b \mathbf{g}_l = \underbrace{\sum_i g_{li}(\mathbf{S}_b)_{ii} g_{li}}_{\text{local}} + \underbrace{\sum_{i \neq j} g_{li}(\mathbf{S}_b)_{ij} g_{lj}}_{\text{cooperative}} \quad (3.15)$$

$$\mathbf{g}_l^T \Delta \Delta^T \mathbf{g}_l = \underbrace{\sum_i g_{li}(\Delta \Delta^T)_{ii} g_{li}}_{\text{local}} + \underbrace{\sum_{i \neq j} g_{li}(\Delta \Delta^T)_{ij} g_{lj}}_{\text{cooperative}}$$

Here,  $(\mathbf{A})_{ij}$  denotes element  $(i, j)$  of matrix  $\mathbf{A}$ . The contributions to the quadratic terms are due to the local contributions from the elements and their cooperative (cross) terms taking into account that each element of a generalized eigenvector ( $g_{li}$ ) corresponds to an element  $x_i$ . However, it is clear that the contributions to the quadratic terms in equation (3.15) can be collected via arbitrary grouping of the elements into subgroups. The computation of the (local) group-wise contributions and their cooperative contributions is equivalent to the computation of the quadratic terms using the corresponding submatrices of  $\mathbf{S}_b$  and  $\Delta \Delta^T$ .

**Asymmetric nature of RPCA.** Multivariate analysis methods can be grouped into methods handling either one state or multiple states. The methods which study one state include methods handling one multivariate variable, such as principal component analysis, factor analysis and independent components analysis, and methods handling two multivariate variables with concurrent measurement (a joint

distribution) such as canonical correlation, regression and partial least squares. Discriminant analysis (classification) methods, on the other hand, aim at finding the variation between several states (classes). Although RPCA is similar to feature extraction for classification (Fukunaga, 1990), there are fundamental differences between pattern classification and the physical definition of the change, which is introduced above. For example, the discriminant features in Fisher discriminant analysis (FDA) are related to the change of the averages of distributions and the information from the covariance matrices is whitened using a unified pooled (average) covariance matrix (Fukunaga, 1990). Therefore, the usage of FDA for dimensionality reduction (Sugiyama et al., 2010) is known to be restricted by a limit on the number of dimensions which is equal to the number of classes minus one. While the objective functions in discriminant analysis are required to be symmetric (Fukunaga, 1990), the objective function (KL divergence) for RPCA is asymmetric, which reflects the directed character of changes (Feng and Crooks, 2008). When studying the backward change from state  $b$  to state  $a$ , the generalized eigenvectors  $\mathbf{G}_R$  pertaining to the reverse direction are the scaled eigenvectors pertaining to the forward direction ( $\mathbf{G}$ ) and the generalized eigenvalues, which present the change of the variance, are inverted ( $1/\lambda_i$ ):

$$\mathbf{G}_R = \mathbf{G}\mathbf{\Lambda}^{-1/2} = \left[ \frac{\mathbf{1}}{\sqrt{\lambda_1}} \mathbf{g}_1 \dots \frac{\mathbf{1}}{\sqrt{\lambda_k}} \mathbf{g}_k \right] \quad (3.16)$$

$$\begin{aligned} \mathbf{G}_R^T \mathbf{S}_b \mathbf{G}_R &= \mathbf{\Lambda}^{-1/2} \mathbf{G}^T \mathbf{S}_b \mathbf{G} \mathbf{\Lambda}^{-1/2} = \mathbf{\Lambda}^{-1/2} \mathbf{\Lambda} \mathbf{\Lambda}^{-1/2} = \mathbf{I} \\ \mathbf{G}_R^T \mathbf{S}_a \mathbf{G}_R &= \mathbf{\Lambda}^{-1/2} \mathbf{G}^T \mathbf{S}_a \mathbf{G} \mathbf{\Lambda}^{-1/2} = \mathbf{\Lambda}^{-1/2} \mathbf{I} \mathbf{\Lambda}^{-1/2} = \mathbf{\Lambda}^{-1} = \mathbf{\Lambda}_R \end{aligned} \quad (3.17)$$

Taking into account that the KL divergence due to the change of the variance  $1/2(-\ln \lambda_i + \lambda_i - 1)$  in equation (3.7) is a convex function with a minimum of 0 at  $\lambda_i = 1$ , we can divide the relative principal components into two groups. The first group includes the components with  $\lambda_i > 1$  where the variance (quantified by  $\lambda_i$ ) along the direction  $\mathbf{g}_i$  increases as we transit from state  $a$  to  $b$ . The second group includes the components with  $\lambda_i < 1$  since the respective variance decreases as we transit from state  $a$  to  $b$ . When the backward change is taking place from  $b$  to  $a$ , the roles of the components of these groups (quantified by  $1/\lambda_i$ ) are exchanged. The same role of these groups is also observed when taking the KL divergence due to the change of the average where the directions  $\mathbf{g}_i$  with  $\lambda_i > 1$  have diminished contributions ( $\frac{1}{2\lambda_i} \mathbf{g}_i^T \Delta \Delta^T \mathbf{g}_i$ ) to the backward direction from  $b$  to  $a$  in comparison with their contribution ( $\frac{1}{2} \mathbf{g}_i^T \Delta \Delta^T \mathbf{g}_i$ ) to the transition from state  $a$  to  $b$  and vice versa.

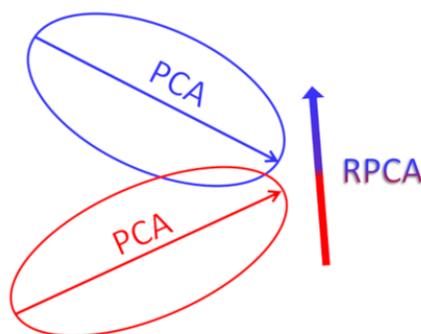
**RPCA and the distance metric.** Even though the task of RPCA (analyzing the change between states) differs from the task of PCA (finding the variation within one state; see Figure 1), a similarity exists regarding the applied distance metric. However, it is important to notice that the generalized eigenvectors are orthonormal with respect to the matrix (Harville, 2008)  $\mathbf{S}_a$  ( $\mathbf{g}_i^T \mathbf{S}_a \mathbf{g}_j = \delta_{ij}$ ), orthogonal with respect to  $\mathbf{S}_b$  ( $\mathbf{g}_i^T \mathbf{S}_b \mathbf{g}_i = \lambda_i$ ;  $\mathbf{g}_i^T \mathbf{S}_b \mathbf{g}_{j \neq i} = 0$ ) and not necessarily orthonormal to each other ( $\mathbf{g}_i^T \mathbf{g}_j \neq \delta_{ij}$ ). RPCA analysis can be interpreted as using the distance metric (Dryden and Mardia, 2016) of state  $a$  to analyze state  $b$ . Indeed, applying the whitening transformation (Fukunaga, 1990, p. 31) in equation (3.5) removes the “information” within state  $a$  ( $\mathbf{W}^T \mathbf{S}_a \mathbf{W} = \mathbf{I}$ ). The eigendecomposition of  $\mathbf{W}^T \mathbf{S}_b \mathbf{W}$  in (3.6) is

performed in the whitened space in which the squared Euclidean distance between two points  $\mathbf{x}_1$  and  $\mathbf{x}_2$  equals their Mahalanobis distance in the original space and the projections on the generalized eigenvectors factorize the Mahalanobis distance (see Supplementary Material S3):

$$\begin{aligned} (\mathbf{x}_1 - \mathbf{x}_2)^T \mathbf{S}_a^{-1} (\mathbf{x}_1 - \mathbf{x}_2) &= (\mathbf{G}^T (\mathbf{x}_1 - \mathbf{x}_2))^T (\mathbf{G}^T (\mathbf{x}_1 - \mathbf{x}_2)) = \sum_{i=1}^k [\mathbf{g}_i^T (\mathbf{x}_1 - \mathbf{x}_2)]^2 \\ &= (\mathbf{W}^T (\mathbf{x}_1 - \mathbf{x}_2))^T (\mathbf{W}^T (\mathbf{x}_1 - \mathbf{x}_2)) \end{aligned} \quad (3.18)$$

Moreover, the average Mahalanobis distance of points in state  $b$  to the average of state  $a$  can be written as the sum of the generalized eigenvalues and the Mahalanobis distance between the averages of the states (see supplementary material S3):

$$\langle (\mathbf{x} - \boldsymbol{\mu}_a)^T \mathbf{S}_a^{-1} (\mathbf{x} - \boldsymbol{\mu}_a) \rangle_b = \underbrace{\text{trace}(\mathbf{W}^T \mathbf{S}_b \mathbf{W})}_{\sum \lambda_i} + \underbrace{\Delta^T \mathbf{S}_a^{-1} \Delta}_{\text{Mahalanobis distance}} \quad (3.19)$$



**Figure 1.** Schematic representation of the conceptual difference between PCA which finds the largest variation within each state and RPCA which finds the change from the initial to the final state.

## 4 Application of RPCA to reveal the thermodynamically important molecular conformational changes

Although our RPCA method is not limited to biomolecular data, the starting point for its development was based on the new insight stemming from our recently introduced perturbation-divergence formalism (PDF) (Ahmad et al., 2017, 2015a). PDF reveals that the KL divergence ( $D_{kl}$ ) equals the work which is spent to change the conformational ensemble and hereby is the thermodynamically relevant objective function for analyzing conformational changes (Ahmad et al., 2015a). Indeed, PDF provides a flexible framework for studying conformational changes involving either all atoms of the structure, a part of them, or aspects of the structure based on phenotypical features (e.g. open/closed, pocket volume, domain movement, etc.); see the discussion in (Ahmad et al., 2017). An important benefit of using the KL divergence for quantifying the energetics related to conformational changes is avoiding the misleading

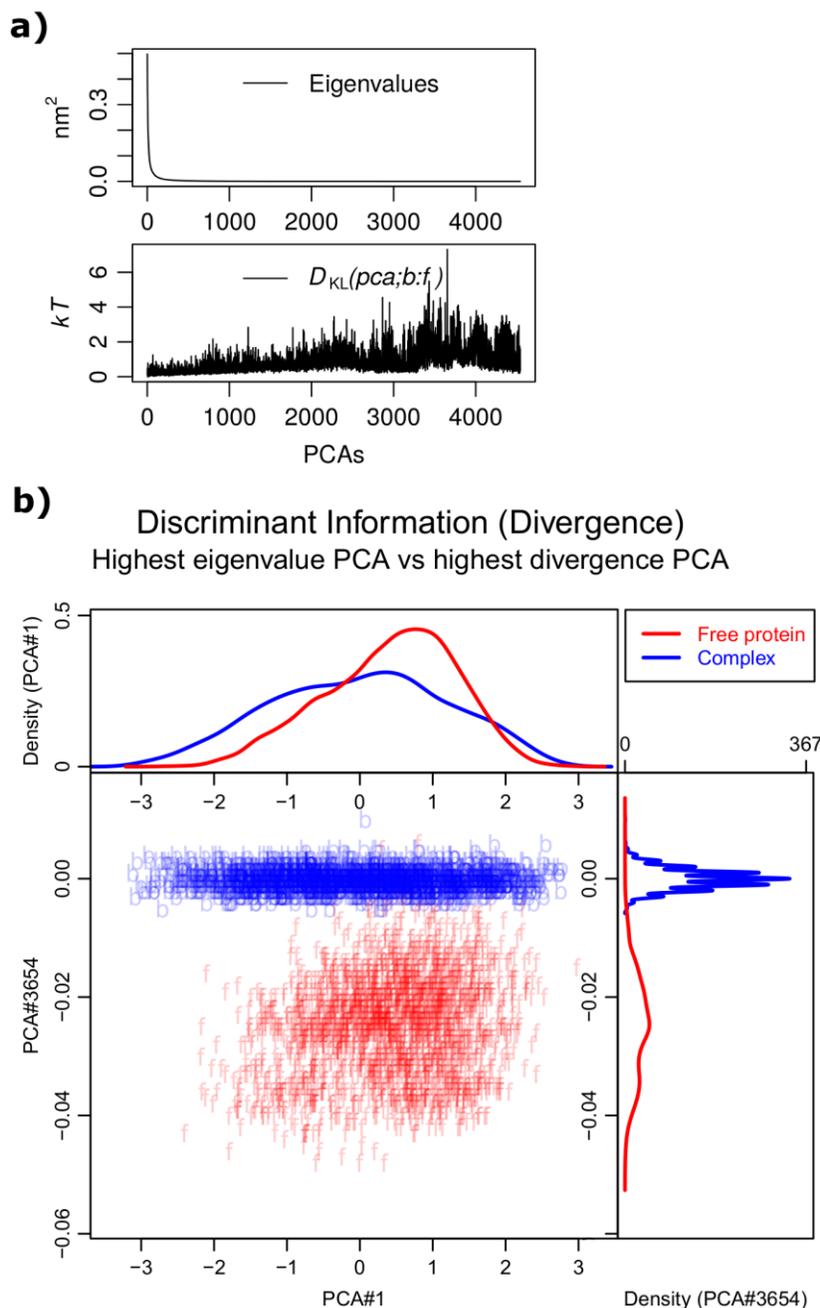
entropic terms which are subject to enthalpy-entropy compensation when using the changes in conformational entropy to estimate the importance of conformational changes, as was previously shown (Ahmad et al., 2016, 2015b, 2017).

In the following, we present use cases for applying RPCA to analyze the molecular conformational changes of the protein HIV-1 protease upon binding to several inhibitor molecules. The conformations of the protein are sampled via molecular dynamics simulations for both the initial (free, unbound) state and the final (bound) state. Technical details on these simulations are provided in Supplementary Material S5. Figure 2 shows an assessment of the relationship between the relevant components of the dynamic change within one state, which are analyzed using traditional PCA of the data points of the final state, and the thermodynamic importance of the corresponding conformational changes. The thermodynamic importance of the conformational changes along a principal component to the association process is quantified by the KL divergence of the distributions of projections of both the free and the bound state conformations on the component (Ahmad et al., 2016). Figure 2a shows that the eigenvalues of the principal components are not related to their thermodynamic importance for the association process. To illustrate this more clearly, Figure 2b shows that the principal component 3654 has more thermodynamic importance (KL divergence) than the first principal component, which is mostly irrelevant for the association process.

Figure 3 shows the RPCA of the conformational changes of wild-type HIV-1 protease upon binding to a drug molecule that has high affinity (Tiplranavir). Presented are both the non-optimal RPCA (Figure 3a) and the optimal RPCA calculated with the sub-spacing method (Figure 3b). The scores (KL divergence) of the importance of the components show that the first few components account for most of the KL divergence. The data points (conformations) of both the free, unbound state and the bound state are projected on selected components to illustrate the correspondence between the score (model based) and the real divergence which can be extracted from the difference between the distributions of the projections of the states. There are clear benefits of the optimal sub-spacing RPCA (Figure 3b). Namely, the first component has a significant divergence because it collects the total contribution to the KL divergence arising from the change of the average. In contrast, one obtains identical averages when projecting the conformations sampled in the two states onto the remaining components. Thus, the remaining components do not contribute to the KL divergence that is due to the change of the averages. Figure 4 shows a comparison of the KL divergence (discriminant information) of the top-scoring component versus the lowest-scoring and the principal component with largest eigenvalue from PCA of the bound state.

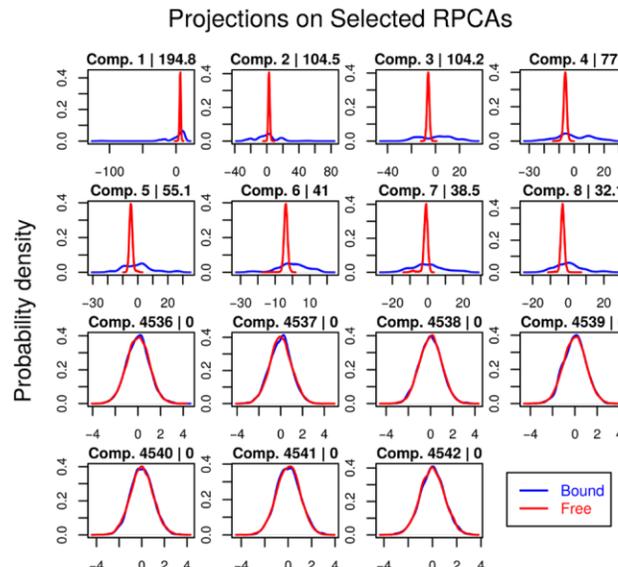
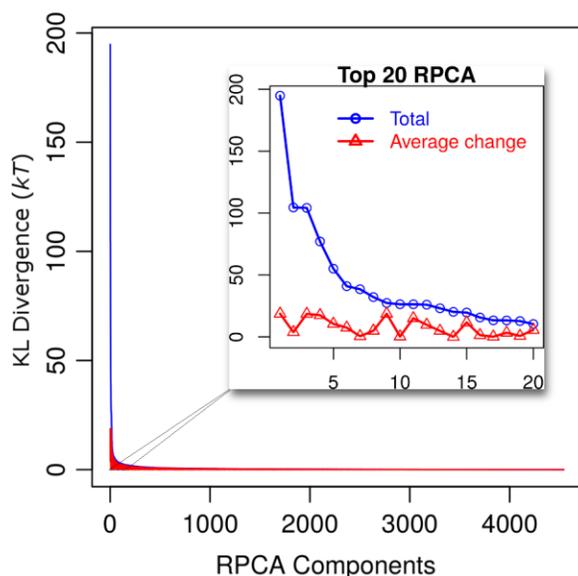
RPCA facilitates studying the conformational changes from both the collective and the local point of view. The relevant collective conformational changes (with respect to their KL divergence) can be presented via reconstructing the conformations in the subspace, which is spanned by the relevant generalized eigenvectors; see equation (3.13). Collective conformational changes can also be represented via reconstructing the conformations from the (normally distributed) latent variable  $\mathbf{y}_m$  using equation (3.14). Interestingly, RPCA provides a clear distinction between the directions (generalized eigenvectors)

along which the fluctuation of the conformations increases upon the change (corresponding to  $\lambda_i > 1$ ) and the directions along which the fluctuation of the conformations decreases (corresponding to  $\lambda_i < 1$ ). Figure 5a shows the collective conformational changes around the average conformation of the bound state along the 33 eigenvectors with the highest generalized eigenvalues ( $\lambda_i > 10$ ) which represent the enhanced motions of the wild-type HIV-1 protease upon binding Tipranavir. Figure 5b, on the other hand, shows the collective conformational changes around the average conformation of the free state along the 33 eigenvectors with the smallest generalized eigenvalues ( $\lambda_i < 0.009$ ) which, in turn, represent the most strongly restricted motions of the wild-type HIV-1 protease upon binding Tipranavir. It should be stressed that the importance of the local conformational changes (e.g. at residues) cannot be inferred from the collective conformational changes. In other words, a large fluctuation in a site, when presenting the collective conformational changes, does not imply larger thermodynamic importance than the associated less strongly fluctuating sites. Alternatively, hotspot analysis from RPCA in equation (3.15) can be used to rank the thermodynamic contribution of the local conformational changes (conformational hot spots) of the biological building blocks (residues) to the association process and the importance of the cooperative (correlation) interactions between them. Figure 6a shows a matrix plot representation (interaction map) of the contributions of the residues to the divergence. In Figure 5b, the relative local residue contributions to the divergence are mapped on the structure of the protein and the importance of the conformational changes is indicated by radii of varying size in the cartoon representation, in addition to using the color code. Interestingly, most of the marked hotspots are residues known to affect the binding affinity upon mutating them. Examples of the defined conformational hot spots are the residues in the active pocket (e.g. D25, V82) and the residues of the flap region (e.g. I50, I54). Figure 7a shows an analysis of the conformational hotspots from RPCA of the binding of a ligand (Saquinavir) to the HIV-1 protease mutant with a resistance-related mutation (I50V) which is located on the flap region outside the binding pocket. The conformational hotspots of the mutant are located at the flap region around the mutation V50. The same flap region does not show important conformational changes when applying the same analysis to the association of Saquinavir to the wild-type (I50) in Figure 7b.

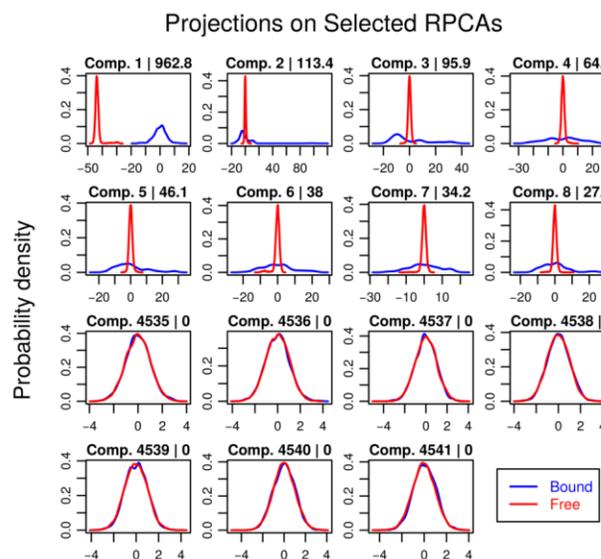
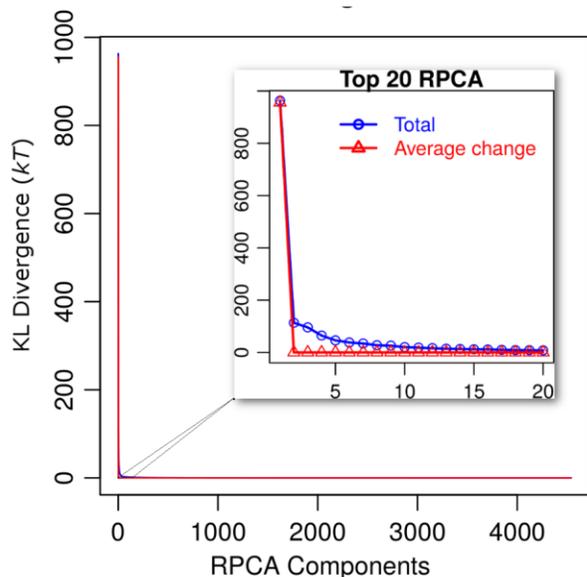


**Figure 2. Assessing the thermodynamic contribution of the conformational changes along the directions of the principal components from PCA to the association process.** **a)** Top panel: Shown are PCA eigenvalues derived from the covariance matrix of the conformational fluctuation (heavy atoms) in an MD simulation of HIV-1 protease bound to its inhibitor Tipranavir. Bottom panel: Shown is the importance (KL divergence) along PCAs computed from the projections of the data points (conformations) of the final state (bound) and the initial state (free). **b)** The divergence of the largest eigenvalue PCA#1 ( $D_{kl} \approx 0.1$ ) is compared to that of the PCA with highest divergence #3654 ( $D_{kl} \approx 7.3$ ). 2000 data points are projected on the PCA vectors of the final (bound; blue) and the initial (free; red) states. The marginal probability distribution densities are shown on the side panels. The PCA analysis is performed on the conformations from the bound state using the heavy atoms of the protein. The superposition of the conformations of both the bound and the free state ensembles and the superposition of the two ensembles to each other are performed in a way similar to RPCA (see below). The KL divergences between the distributions of the projections of the two states are computed using the KLIEP method (Sugiyama et al., 2008).

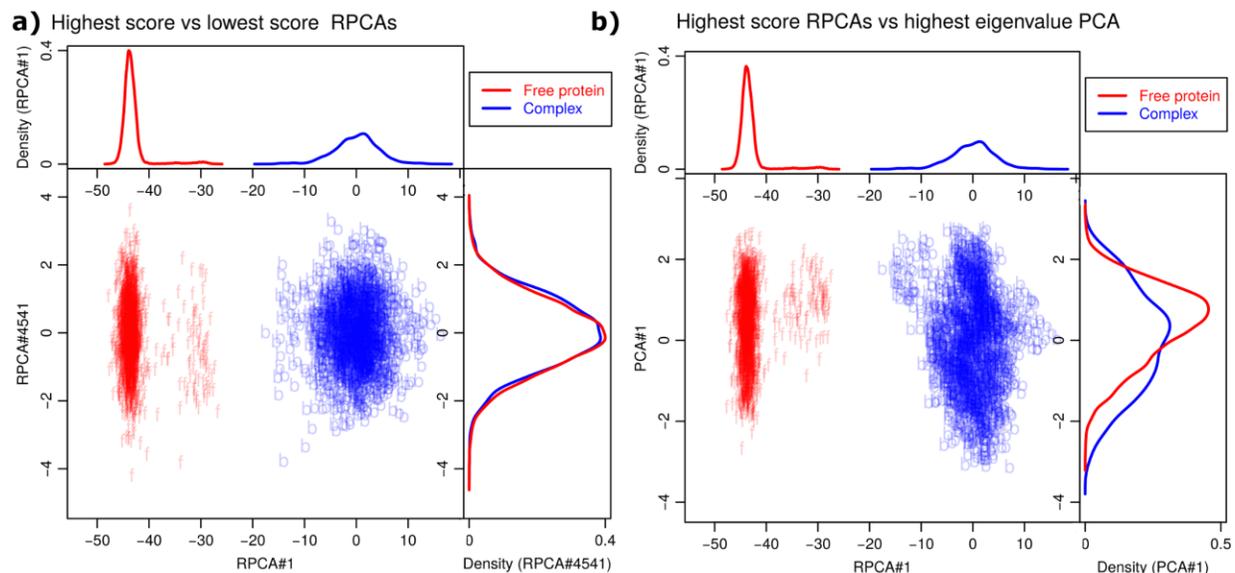
### a) RPCA (without sub-spacing)



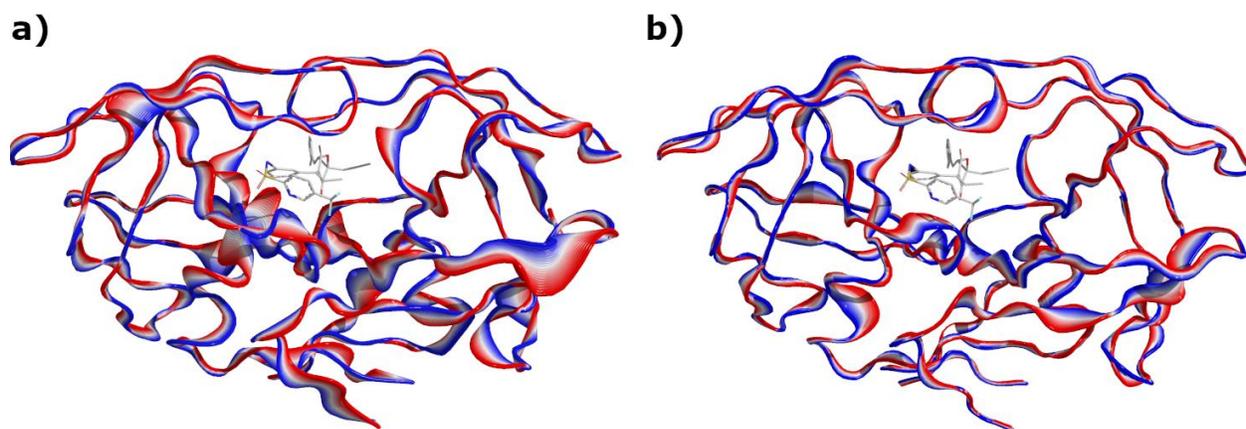
### b) RPCA with sub-spacing



**Figure 3. RPCA of the conformational changes of HIV-1 protease upon its association with Tipranavir.** Shown are the non-optimal RPCA (a) from equation (3.4) and the optimal RPCA with sub-spacing (b) from equation (3.12). The right panels show the KL divergences of the components (blue colored) and their corresponding contributions due to the change of the average (red colored). Left panels show projections of the data points (conformations) of both the initial (free, unbound) and the final bound state on selected components. The scores (KL divergences) of the components are displayed after their corresponding number. The projections show that the components (components 1-8) with the highest rank (KL divergence) distinguish the change between the free and bound states while the components with the lowest rank do not distinguish the change (similar projections). The analysis is performed using the heavy atoms of the protein. The plots are generated using R (*R*, 2012) and the densities are smoothed using the kernel density estimation.

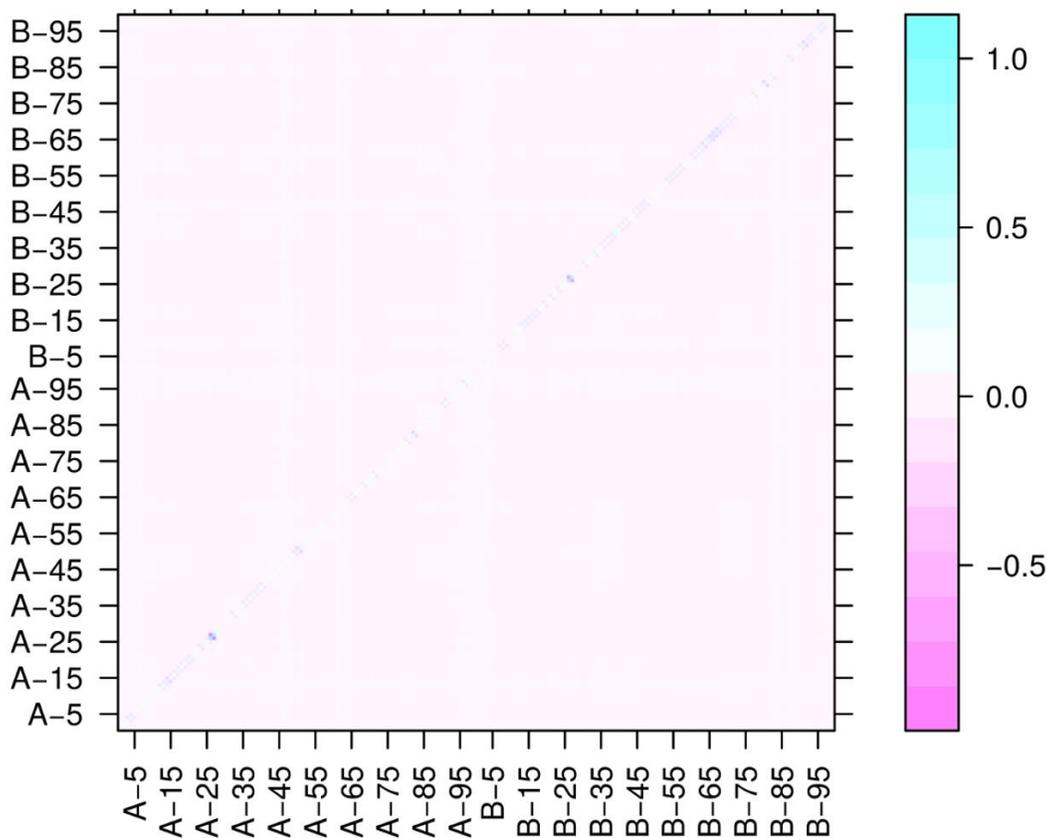


**Figure 4.** KL divergences of the highest scored relative principal component. The initial state (red) and the bound state (blue) are mapped on the top-scoring RPCA#1 against the lowest-scoring RPCA#4541 (a) and the top-scoring RPCA#1 against the top-scoring PCA#1 (highest eigenvalue) (b). RPCA is successful in extracting the components holding the most significant information on the change (binding process from free to bound state). The marginal probability distribution densities are shown on the side panels. Samples of 2000 data points were used for the projections.

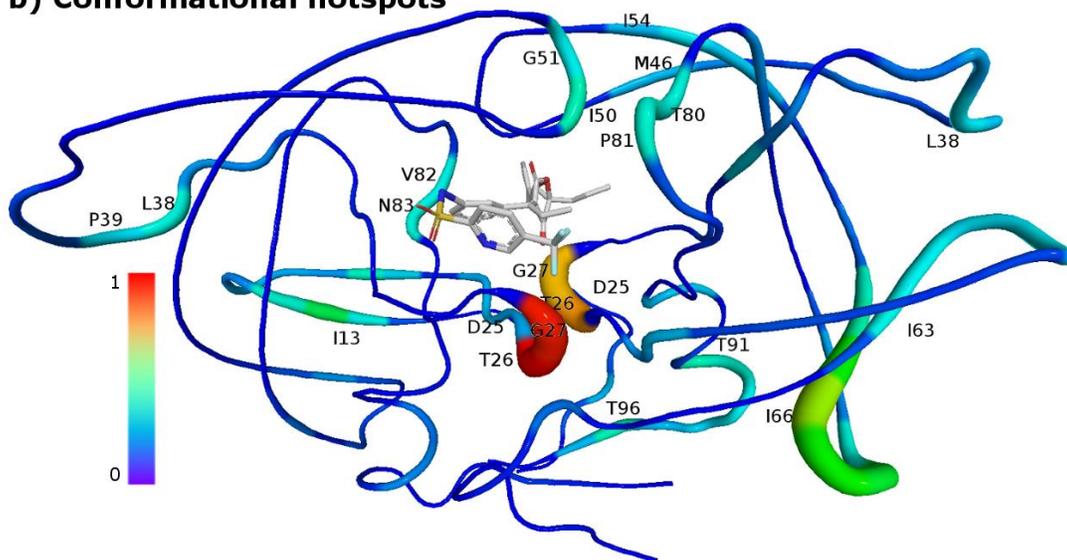


**Figure 5. Representation of the enhanced and restricted conformational fluctuations of HIV-1 protease upon binding the inhibitor Tipranavir.** Conformations around the average conformation are reconstructed from the latent variable after interpolation around its average along selected generalized eigenvectors; see equation (3.14). **a)** Enhanced conformational fluctuations around the average structure of the bound state along the 33 eigenvectors with the highest generalized eigenvalues ( $\lambda_i > 10$ ). These conformational fluctuations increase the affinity of binding by optimizing the local conformations in the ligand-receptor interface. **b)** Conformational fluctuations around the average structure of the free state along the 33 eigenvectors with the smallest generalized eigenvalues. These latter fluctuations are highly restricted upon the association ( $\lambda_i < 0.009$ ) because they decrease the affinity of binding via adverse local movements in the binding pocket and opening of the flap regions; see the supplementary movies. The cartoon representation is generated using Pymol (Schrödinger, LLC, 2018). The conformation of the ligand is taken from the experimental structure. Movies of these conformational changes are presented in the Supplementary Material S6.

**a) Residues interactions map (WT HIV-1 protease with Tipranavir)**

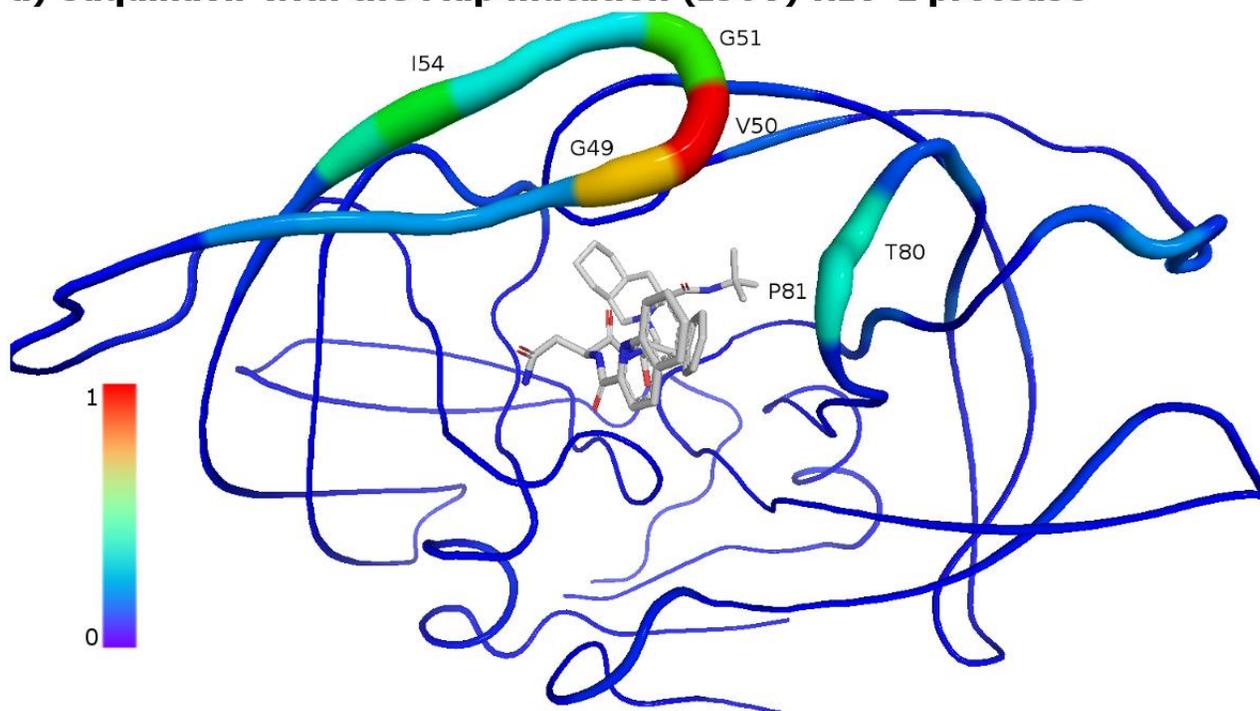


**b) Conformational hotspots**

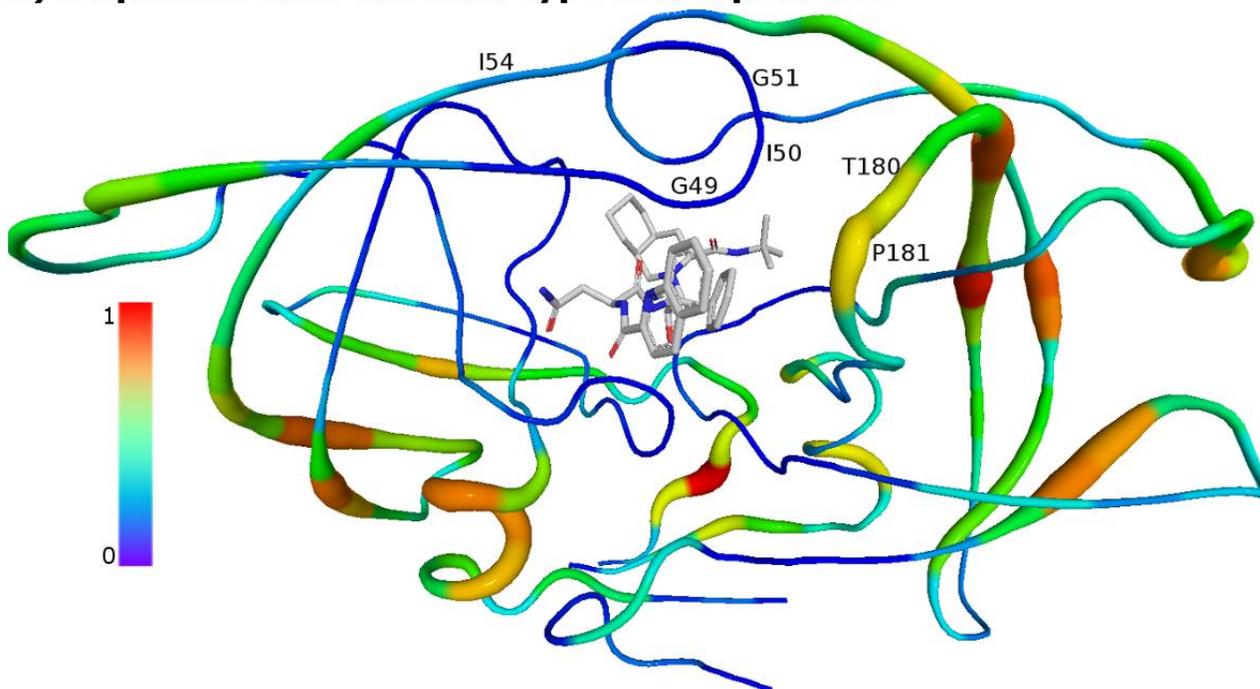


**Figure 6. a) Conformational interaction map** between the residues of HIV-1 protease upon binding its inhibitor Tipranavir. The contributions are concentrated around the diagonal elements indicating the important local conformational changes and the propagation of the conformational changes through the neighbor residues. **b) Conformational hotspots** mapped on the structure of the protein. The radius and the color of the cartoon indicate the importance of the conformational changes. The importance is normalized relative to the highest value that is set to 1.

**a) Saquinavir with the Flap mutation (I50V) HIV-1 protease**



**b) Saquinavir with the Wild-type HIV-1 protease**



**Figure 7.** Conformational hotspots from RPCA analysis recognize the importance of the resistance-related mutation (I50V) of HIV-1 protease when bound to Saquinavir. The conformational hotspots of the mutant (a) are located at the flap region around the mutation V50. The corresponding flap residues of the wild-type (I50) in figure (b) do not show energetically important (expensive) conformational changes. The radius and the color of the cartoon indicate the importance of the conformational changes. The importance is normalized relative to the highest value.

## Technical aspects of using RPCA of molecular conformational changes

**Optimal superposition of conformations of the ensembles.** The first step in analyzing an ensemble of conformations is removing the six external rotational and translational degrees of freedom via superimposing the conformations. Although the internal coordinates are not altered due to these similarity transformations (rotations and translations), the average conformation and the covariance matrix are highly affected by the way of superimposing the conformations to remove these external degrees of freedom. Traditionally, the conformations are superimposed on an arbitrary reference structure (e.g. the starting structure of the MD simulation). However, the resulting ensemble and the average structure are dependent on the reference structure. The solution of this problem is well known in statistical shape analysis as the Generalized Procrustes Analysis (Dryden and Mardia, 2016; Gapsys and de Groot, 2013) (GPA). GPA returns a compact ensemble of the conformations via minimizing the sum of the Euclidean distances between the conformations (which is equal to their Euclidean distances to the average conformation). GPA is performed via an iterative algorithm having two steps: (1) compute the average structure, (2) superimpose the conformations on the new average structure. Repeating these steps ensures the convergence of the average structure. Fortunately, few iterations are usually enough for convergence.

**Superposition of two ensembles.** The Mahalanobis distance term ( $\Delta^T \mathbf{S}_a^{-1} \Delta$ ) of the KL divergence in equation (3.9) is affected by the fashion in which we superimpose the average conformations ( $\mu_b, \mu_a$ ). However, superimposition via minimizing the Euclidean distance may overestimate the KL divergence via an artificial contribution to the Mahalanobis distance  $\Delta^T \mathbf{S}_a^{-1} \Delta$ ;  $\Delta = \mu_b - \mu_a$ . Therefore, superimposition of the average conformations should aim at minimization of the Mahalanobis distance. In contrast to the minimization of the Euclidean distance, there is no analytical solution known for minimizing the Mahalanobis distance. This type of nonlinear minimization is known as the Covariance Weighted Procrustes Analysis (Brignell et al., 2016) (CWPA) and numerical methods can be used to find the optimal superposition (rotations and translations) for minimizing the quadratic term  $\Delta^T \mathbf{S}_a^{-1} \Delta$ .

### Steps for performing a RPCA analysis of two molecular states:

- (1) GPA fitting of the conformations sampled in the simulation of the first state. The covariance matrix of this state is also computed at this step. A successful superimposition of the ensemble will lead to a singular covariance matrix with at least six eigenvalues of zero value accounting for removing the external degrees of freedom.
- (2) GPA fitting of the conformations sampled in the simulation of the second state to obtain the average conformation.
- (3) Covariance weighted fitting of the average conformation of the second state on the average conformation of the first state via minimizing their Mahalanobis distance. This unconstrained nonlinear optimization is numerically performed using the line-search algorithm and the BFGS factored method to update the Hessian (Dennis and Schnabel, 1996, p. 356).
- (4) The new average conformation of the second state is used as a reference to refit the conformations of the second state and to compute the covariance matrix of the second state.

- (5) Simultaneous diagonalization of the covariance matrices is performed. Optionally, the sub-spacing optimal algorithm can be used.
- (6) KL divergences of the relative principal components are computed and the components are reordered based on their scores (KL divergences).

We have developed efficient computational tools to perform these steps. The tools are written in the C programming language and the numerical linear algebra operations are performed using the BLAS and LAPACK routines (Anderson et al., 1999). The limit of the zero value of the eigenvalues is defined using the machine precision multiplied by the largest eigenvalue. The covariance-weighted superimposition is performed using the nonlinear minimization algorithm by Dennis and Schnabel (Dennis and Schnabel, 1996, p. 356).

## 5 Conclusion

Here, we introduced the RPCA method, which extracts the relevant principal components describing the change between two macroscopic states of a dynamic system represented by two data sets. The definition of the macroscopic change of a dynamic system and its quantification are based on previous work where we derived a generalized quantification of the change of a physical system based on statistical mechanics. We presented use cases for conformational changes taking place upon ligand binding to HIV-1 protease that clearly illustrate the power of RPCA to characterize the relevant changes between two ensembles in a high-dimensional space. Moreover, software solutions were introduced to ensure the removal of the similarity transformations (rotations and translations) via superposing the conformations using GPA and CWPA. These procedures may also be beneficial for preparing the conformations for other analysis methods.

Although RPCA is currently limited to handling only continuous variables and two macroscopic states, the introduced framework for quantifying changes of dynamic systems using the exponential family of distributions is flexible regarding the nature of probability distributions and the nature of the microscopic variables (continuous and categorical). Therefore, the presented theoretical formalism opens the door for developing improved new methods for mining the factors underlying changes in dynamic systems in the directions of (i) handling both continuous and categorical data (e.g. the effect of sequence changes (mutations) on the binding affinity) and of (ii) handling multiple macroscopic states (e.g. study the binding of a series of ligands to a receptor).

## 6 References

- Ahmad M, Helms V, Kalinina OV, Lengauer T. 2017. Elucidating the energetic contributions to the binding free energy. *J Chem Phys* **146**:014105.
- Ahmad M, Helms V, Kalinina OV, Lengauer T. 2016. The Role of Conformational Changes in Molecular Recognition. *J Phys Chem B* **120**:2138–2144. doi:10.1021/acs.jpcc.5b11593
- Ahmad M, Helms V, Lengauer T, Kalinina OV. 2015a. How Molecular Conformational Changes Affect Changes in Free Energy. *J Chem Theory Comput* **11**:2945–2957. doi:10.1021/acs.jctc.5b00235
- Ahmad M, Helms V, Lengauer T, Kalinina OV. 2015b. Enthalpy–Entropy Compensation upon Molecular Conformational Changes. *J Chem Theory Comput* **11**:1410–1418. doi:10.1021/ct501161t
- Amadei A, Linssen ABM, Berendsen HJC. 1993. Essential dynamics of proteins. *Proteins Struct Funct Bioinforma* **17**:412–425. doi:10.1002/prot.340170408
- Anderson E, Bai Z, Bischof C, Blackford LS, Demmel J, Dongarra J, Du Croz J, Greenbaum A, Hammarling S, McKenney A. 1999. LAPACK Users' guide. SIAM.
- Barndorff-Nielsen OE. 1978. Information and Exponential Families in Statistical Theory. Chichester: John Wiley & Sons Ltd.
- Berger R, Casella G. 2001. Statistical Inference, 2 ed edition. ed. Australia ; Pacific Grove, CA: Cengage Learning, Inc.
- Bishop CM. 1998. Latent Variable Models In: Jordan MI, editor. Learning in Graphical Models, NATO ASI Series. Springer Netherlands. pp. 371–403. doi:10.1007/978-94-011-5014-9\_13
- Brignell CJ, Dryden IL, Browne WJ. 2016. Covariance Weighted Procrustes Analysis In: Turaga PK, Srivastava A, editors. Riemannian Computing in Computer Vision. Springer International Publishing. pp. 189–209. doi:10.1007/978-3-319-22957-7\_9
- Brown LD. 1986. Fundamentals of Statistical Exponential Families with Applications in Statistical Decision Theory. *Lect Notes-Monogr Ser* **9**:i–279.
- Cover TM, Thomas JA. 2006. Elements of Information Theory, 2nd ed. Hoboken, N.J: Wiley-Interscience.
- Dennis J, Schnabel R. 1996. Numerical Methods for Unconstrained Optimization and Nonlinear Equations, Classics in Applied Mathematics. Society for Industrial and Applied Mathematics. doi:10.1137/1.9781611971200
- Dryden IL, Mardia KV. 2016. Statistical Shape Analysis: With Applications in R. John Wiley & Sons.
- Feng E, Crooks G. 2008. Length of Time's Arrow. *Phys Rev Lett* **101**. doi:10.1103/PhysRevLett.101.090602
- Fukunaga K. 1990. Introduction to Statistical Pattern Recognition, Second Edition, 2 edition. ed. Boston: Academic Press.
- Gapsys V, de Groot BL. 2013. Optimal Superpositioning of Flexible Molecule Ensembles. *Biophys J* **104**:196–207. doi:10.1016/j.bpj.2012.11.003
- Greven A, Warnecke G, Keller G, editors. 2003. Entropy. Princeton: Princeton University Press.
- Harville DA. 2008. Matrix Algebra From a Statistician's Perspective. Springer Science & Business Media.
- Husic BE, Pande VS. 2018. Markov State Models: From an Art to a Science. *J Am Chem Soc.* doi:10.1021/jacs.7b12191

Kirkwood JG. 1934. On the Theory of Strong Electrolyte Solutions. *J Chem Phys* **2**:767–781. doi:doi:10.1063/1.1749393

Kitao A, Go N. 1999. Investigating protein dynamics in collective coordinate space. *Curr Opin Struct Biol* **9**:164–169. doi:10.1016/S0959-440X(99)80023-2

Krivobokova T, Briones R, Hub JS, Munk A, de Groot BL. 2012. Partial Least-Squares Functional Mode Analysis: Application to the Membrane Proteins AQP1, Aqy1, and CLC-ec1. *Biophys J* **103**:786–796. doi:10.1016/j.bpj.2012.07.022

Kullback S. 1997a. Information Theory and Statistics, Revised edition. ed. Mineola, N.Y: Dover Publications.

Kullback S. 1997b. Inequalities of Information Theory Information Theory and Statistics. Dover Publications. pp. 36–59.

Kullback S, Leibler RA. 1951. On Information and Sufficiency. *Ann Math Stat* **22**:79–86. doi:10.1214/aoms/1177729694

Lehmann EL, Casella G. 2003. Theory of Point Estimation, 2nd ed. 1998. Corr. 4th printing 2003 edition. ed. New York: Springer.

Lin J. 1991. Divergence measures based on the Shannon entropy. *IEEE Trans Inf Theory* **37**:145–151.

Mardia KV, Kent JT, Bibby JM. 1979. Multivariate analysis. Academic Press.

Martin RS, Wilkinson JH. 1971. Reduction of the Symmetric Eigenproblem  $Ax = \lambda Bx$  and Related Problems to Standard Form In: Bauer PDFL, editor. Linear Algebra, Handbook for Automatic Computation. Springer Berlin Heidelberg. pp. 303–314. doi:10.1007/978-3-662-39778-7\_21

Murphy KP. 2012. Machine Learning: A Probabilistic Perspective. MIT Press.

R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria, 2012. 2012.

Rockafellar RT. 1970. Convex Analysis. Princeton University Press.

Schrödinger, LLC. 2018. The PyMOL Molecular Graphics System, Version 2.1.

Stewart GW. 2001. Matrix Algorithms Volume 2: Eigensystems. SIAM.

Sugiyama M, Kawanabe M, Chui PL. 2010. Dimensionality reduction for density ratio estimation in high-dimensional spaces. *Neural Netw* **23**:44–59.

Sugiyama M, Suzuki T, Kanamori T. 2012. Density Ratio Estimation in Machine Learning. Cambridge University Press.

Sugiyama M, Suzuki T, Nakajima S, Kashima H, Bünau P von, Kawanabe M. 2008. Direct importance estimation for covariate shift adaptation. *Ann Inst Stat Math* **60**:699–746. doi:10.1007/s10463-008-0197-x

Theodoridis S. 2015. Machine Learning: A Bayesian and Optimization Perspective. Academic Press.

Zwanzig RW. 1954. High-Temperature Equation of State by a Perturbation Method. I. Nonpolar Gases. *J Chem Phys* **22**:1420–1426. doi:doi:10.1063/1.1740409