## Original Paper

# Improved Protein-ligand Prediction Using Kernel Weighted Canonical Correlation Analysis

RAISSA RELATOR[1,a)]   TSUYOSHI KATO[1,b)]   RICHARD LEMENCE[2,c)]

**Abstract:** Protein-ligand interaction prediction plays an important role in drug design and discovery. However, wet lab procedures are inherently time consuming and expensive due to the vast number of candidate compounds and target genes. Hence, computational approaches became imperative and have become popular due to their promising results and practicality. Such methods require high accuracy and precision outputs for them to be useful, thus, the problem of devising such an algorithm remains very challenging. In this paper we propose an algorithm employing both support vector machines (SVM) and an extension of canonical correlation analysis (CCA). Following assumptions of recent chemogenomic approaches, we explore the effects of incorporating bias on similarity of compounds. We introduce kernel weighted CCA as a means of uncovering any underlying relationship between similarity of ligands and known ligands of target proteins. Experimental results indicate statistically significant improvement in the area under the ROC curve (AUC) and F-measure values obtained as opposed to those gathered when only SVM, or SVM with kernel CCA is employed, which translates to better quality of prediction.

**Keywords:** canonical correlation analysis, kernel methods, protein-ligand interaction, support vector machines

## 1. Introduction

Drug discovery is a multi-staged process which involves the determination of existing interactions between a compound and a protein. Many drugs are developed depending on the reaction they produce when coupled with the respective proteins acting during a biological process in the body. However, only a few existing interactions have actually been validated through experiments. Moreover, wet lab procedures are inherently time consuming and expensive due to the vast number of candidate compounds and target genes. Hence, computational approaches became imperative and have become popular due to their promising results and practicality.

The protein-ligand interaction prediction problem can be viewed as a task of filling up a protein-ligand matrix whose rows represent the candidate compounds and the columns represent the target proteins as shown in the example in **Fig. 1** (a). A matrix entry is +1 if there is interaction between the corresponding drug and target. Otherwise, −1. Only a few interactions have actually been verified and recorded which makes the protein-ligand matrix sparse. Termed as the 'chemogenomic approach' by Rognan [20], the ultimate goal of this task is to identify all the ligands of each target, thus, fully matching the ligand and target spaces [3].

Many in silico methods have already been developed to address this problem. We can classify these methods into two: the struc-

ture or docking approach and the ligand-based approach. Docking approaches make use of 3D structures of the chemical compounds or the proteins to find protein-ligand pairs which are more likely to bind [2], [4], [5]. On the other hand, ligand-based techniques usually employ machine learning algorithms in comparing known ligands and candidate ligands of a certain target even without any prior information regarding their structure [9], [14], [25]. In this study, we shall make use of the ligand-based approach.

There are two ways of approaching the task of interaction prediction: one is by using the global model [2], [17], and another one is via the local model [3], [14], [25]. The global model utilizes a large interaction matrix and imputation of missing values is done simultaneously. Each cell in the interaction matrix is considered as a sample to which statistical methods are applied. Descriptors of ligands in the form of a feature matrix and some information for target proteins are combined to generate a fused profile for each cell in the interaction matrix. An advantage is that interaction prediction for target proteins with few known interactions can still be formed. However, since the model aims
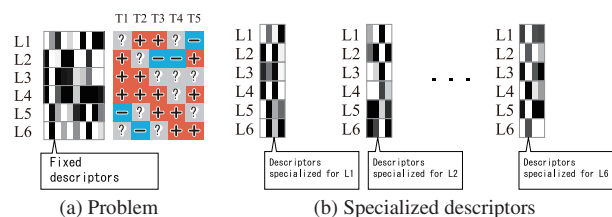


(a) Problem          (b) Specialized descriptors

**Fig. 1** Protein-ligand matrix and descriptors. In the example depicted in (a), the prediction task is to impute 11 missing entries in the 6×5 protein-ligand matrix using 10-dimensional raw descriptors of ligands. The problem can be divided into six sub-problems, each of which is to complete a row in the protein-ligand matrix. Our algorithm extracts compact descriptors specialized in each sub-problem.

---

1   Department of Computer Science, Graduate School of Engineering, Gunma University, Kiryu, Gunma 376–8515, Japan
2   Institute of Mathematics, College of Science, University of the Philippines, Diliman, Quezon City 1101, Philippines
a)   relator-raissa@kato-lab.cs.gunma-u.ac.jp
b)   katotsu@cs.gunma-u.ac.jp
c)   rslemence@math.upd.edu.ph

to exploit information from similar columns, some useful information for learning the rule for prediction may be corrupted by information from irrelevant columns.

Meanwhile, in the local model approach, prediction is made for each column of the protein-ligand table independently — the approach finds unknown chemical compounds which are similar to known ligands interacting with the target protein of interest. The local model often suffers from a small-sample problem. Many columns in the protein-ligand interaction matrix include few positive interactions, causing machine learning algorithms to be trained with few positive samples despite very high dimensionality of ligand descriptors.

The goal of determining interactions between targets and compounds is established under twofold assumptions [3], [20]: First is that compounds with similar properties tend to share targets. And, targets with similar ligands share similarities in structures such as binding sites. These have been verified by recent studies by considering drug side effects [7] and similarities among ligands [16]. Moreover, integrated approaches exploring both protein and compound similarities have also been investigated [6], [14], [26]. Thus, recent methodologies have allowed us to make predictions on interactions based on similarity measures for ligands and targets.

Motivated by the assumption that similar ligands tend to have similar target proteins [15], [23], our goal is to uncover any underlying relationship between a set of ligands and exploit this relationship, together with some known ligand-target interactions, to predict new interactions. We search for ligands with strong associations by finding correlations between them using their features.

In this paper, we present a weighted extension of canonical correlation analysis (WCCA) in the reproducing kernel Hilbert space (RKHS) in an attempt to introduce advantageous properties of local models to the global model approach. To estimate the missing entries in each row of the interaction matrix, we use kernel WCCA (KWCCA) to extract essential features which are specialized in imputation of the corresponding row. The extracted features are compact enough for local models to be trained with a small training set composed from the column. Through the experiments with data of GPCRs and odorant receptors, the prediction performance is shown to be improved when our algorithm is applied compared to several existing methods.

### 1.1   Related Works

A popular and useful technique in investigating relationships between sets of data is the so-called canonical correlation analysis (CCA) [12]. First introduced by Hotelling [13], CCA generally aims to find linear transformations which maximize the correlation between a pair of data. However, the common information extracted from the data sources may not be as useful if nonlinear correlations exist. For this reason, kernel CCA (KCCA) was introduced to offer an alternative solution via the kernel trick, where CCA is performed in a reproducing kernel Hilbert space (RKHS), typically a higher dimensional feature space [1].

Several variants of CCA have been developed and applied to different problem settings. For instance, Yu et al. [31] intro-

duced weights to CCA. Although we also introduce weighting in our proposed method, the authors' purpose and formulation are totally different from ours: they assumed more than two data sources and weight each source, whereas, in our formulation, we assume two data sources and each sample is weighted. On the other hand, in a biologically-related setting, Yamanishi et al. [29] employed multiple KCCA and integrated KCCA for gene cluster extraction. One is done by maximizing the sum of pairwise correlations and the other by maximizing correlation of combination of attributes.

For the problem of functional site prediction, Gonzalez et al. [11] incorporated KCCA to find amino acid pairs and protein functional classifications which are maximally correlated. This technique was motivated by the *Xdet* method [18] and CCA was employed as an alternative to computing Pearson correlation.

The indefinite kernel CCA (IKCCA) was developed by Samarov et al. [22] with a motivation similar to ours. They removed the similarity of samples outside the neighborhood to refine the analysis. The operation often yields an indefinite matrix. IKCCA finds a definite matrix close to the indefinite matrix to perform CCA on the definite matrix. However, their usage of employing CCA is different from ours: the inputs of their approach are positive pairs of ligands and proteins, whereas our approach applies CCA to two different types of ligand profiles. IKCCA is formulated with a saddle-point problem that is solved by minimizing a maximum, but the numerical algorithm to solve the problem has not been shown.

Another important variant is sparse CCA [27], [28] which uses *lasso* or *elastic net* techniques to encourage loading matrices to be sparse. This approach was also applied to a set of protein-ligand pairs with positive interactions in order to elucidate meaningful chemical descriptors in Ref. [28]. Another is the Supervised Regularized CCA [10] which allows integration of multimodal data. Such method can be very useful when involving non-image and image data samples.

## 2.   Materials and Methods

### 2.1   Data

The data used for this study was originally from Ref. [21]. The given interaction matrix consists of 62 mammalian odorant receptors (ORs) as target proteins and 63 odorants as candidate ligands. It is binary in form and contains 340 positive interactions. The number of known positive interactions for each target protein is at least one and at most thirty-seven, while the median is three. Some randomly selected protein-ligand pairs are assumed to be unknown to test prediction methods, and the values of the cells are set to zero. Each row in the interaction matrix provides an *interaction profile* of the ligand.

From the chemical IDs supplied, we searched PubChem [*1] for the chemical structures of the odorants to obtain the descriptors of the ligands. Frequent substructures are employed as descriptors of ligands. The frequent substructures are mined with a software named *gSpan* [30]. The software is applied to the 63 chemical structures, and the 60,311 binary descriptors are obtained as

---

[*1]   http://pubchem.ncbi.nlm.nih.gov/

*chemical profiles*.

## 2.2   Overview of the Algorithm

Our approach consists of two stages: First, we consider sub-problems, each of which involves imputation on a single row in the interaction matrix, and use weighted CCA to extract a compact vector representation for each sub-problem. Then, we apply SVM for prediction of each cell using the corresponding descriptor extracted in the previous stage. This technique is overviewed as follows.

*Chemical profiles* obtained from chemical structures contain numerous features that are not important for prediction. Extracting significant features from such chemical profiles is crucial for accurate prediction of protein-ligand interaction. To accomplish this, we have to find effective low-dimensional representations of the original chemical profiles lying in the extremely high-dimensional *chemical space*.

*Interaction profiles* describe the existence and the absence of interactions with several target proteins. More often than not, target proteins share similar properties. For this reason, interaction profiles approximately span a low-dimensional space, say $\mathbb{R}^m$, which we shall also extract from a high-dimensional *interaction space*, in a similar fashion as the chemical profiles.

Canonical correlation analysis uses a set of chemical profiles and interaction profiles to find two projection functions, $\phi_{\mathrm{ch}}$ and $\phi_{\mathrm{in}}$, simultaneously: The projection $\phi_{\mathrm{ch}}$ is from the chemical space to the low-dimensional canonical space $\mathbb{R}^m$, and $\phi_{\mathrm{in}}$ is from the interaction space to $\mathbb{R}^m$. The images of $\phi_{\mathrm{ch}}$ are used to approximate the images of $\phi_{\mathrm{in}}$. The projections obtained by CCA are shown mathematically to be the minimizer of the expected deviation of the image of $\phi_{\mathrm{ch}}$ from the image of $\phi_{\mathrm{in}}$.

**Figure 2** (a) is an illustration of how CCA works with chemical profiles and interaction profiles. In this figure, the shaded squares are data representations of the feature vector of each ligand in the chemical space. While the open circles are the data representation of the interaction vector of each ligand in the interaction space. The images under $\phi_{\mathrm{ch}}$ and $\phi_{\mathrm{in}}$ of these data points are plotted in the canonical space, and their corresponding images are linked with a dashed line. CCA finds the projections $\phi_{\mathrm{ch}}$ and $\phi_{\mathrm{in}}$ so that the average squared length of the dashed lines is minimized.

In application to protein-ligand interaction prediction, estimating the images for all ligands is not necessary; it is only for the ligand whose interactions we wish to predict that the image of the chemical compound is desired to be well approximated. To obtain a good approximation for a ligand of interest, it is sufficient to estimate projections so that only the images of similar ligands are approximated well. The precisions of the approximations for ligands dissimilar to the ligand of interest barely affect the accuracy of the solution. This consideration motivated us to assign weights to ligands according to their similarity to the ligand of interest, and to extend the classical CCA so that the weighted average deviation is minimized. The weighted CCA almost disregards ligands with small weights to find projections, achieving more accurate approximations for the ligand of interest. We refer to the extension of CCA as weighted CCA.

Figure 2 (b) illustrates the effects of weighted CCA when weights are added to similar ligands. In this context, we define similarity as the measure of affinity between features of compounds. This can be represented by the distance between the data representation of the ligands in the chemical space. In the given figure, the chemical profile for a ligand of interest is marked with a star, and profiles of similar ligands are colored red. In a similar manner, we interpret points of the same color as ligands sharing similarities in their chemical properties, hence their grouping in the chemical space. The two figures, (a) and (b), allow us to compare classical CCA with weighted CCA: the deviations for red points in (b) are smaller than those in (a). The deviations for other ligands are larger, which hardly worsen the performance of predicting the interaction of the protein of interest.

The final prediction result is obtained in the post-processing stage using SVM. The images of the projections are used for SVM learning. SVM is trained well if a good training set is given. Hence, ligands with poor approximations by CCA, which are noisy for SVM learning, are preferably excluded. The im-
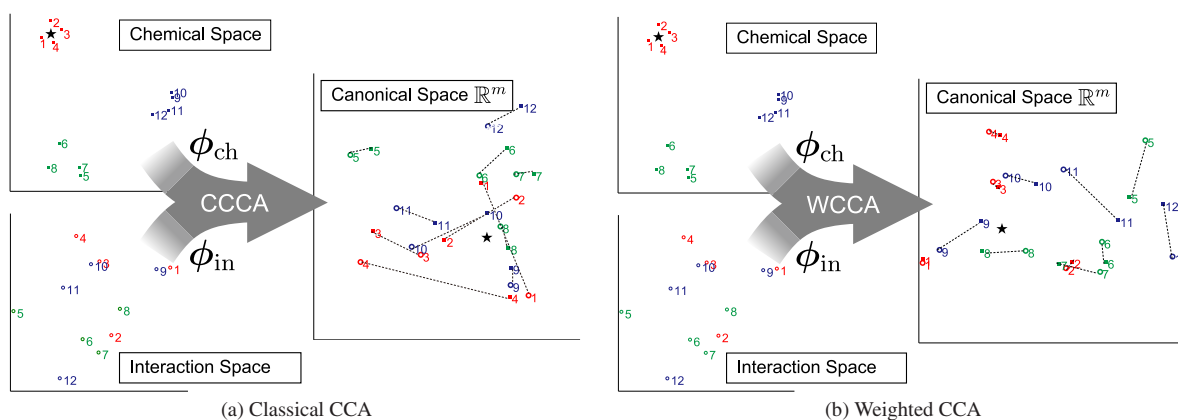


(a) Classical CCA                    (b) Weighted CCA

**Fig. 2**   Classical CCA and weighted CCA. Our approach projects chemical and interaction profiles into a low-dimensional canonical space so that the images are close to each other. The star point represents the ligand of interest, and red points are ligands sharing similarities with the ligand of interest. Although the classical CCA minimizes the average deviation over all the ligands, to achieve accurate prediction, it is sufficient that the deviations between the images of the target ligand and the ligands similar to it are small. The weighted CCA works with arbitrarily specified weights, which ensures small deviations for red points by giving them larger weights.

ages are already in a low-dimensional space in which SVM learning works well even with a small training set, encouraging us to assign smaller weights to ligands with poor approximations for SVM learning.

### 2.3 Weighted CCA

In this subsection we present the details of weighted CCA. We denote the chemical profile and the interaction profile, respectively, by a $p_{ch}$-dimensional vector $x^{ch}$ and a $p_{in}$-dimensional vector $x^{in}$. Assuming that the functions $\phi_{ch} : \mathbb{R}^{p_{ch}} \to \mathbb{R}^m$ and $\phi_{in} : \mathbb{R}^{p_{in}} \to \mathbb{R}^m$ are affine transformations allows us to express them as $\phi_{ch}(x^{ch}) = W_{ch}^\top(x^{ch} - \mu_{ch})$, $\phi_{in}(x^{in}) = W_{in}^\top(x^{in} - \mu_{in})$, where $W_{ch} \in \mathbb{R}^{p_{ch} \times m}$, $\mu_{ch} \in \mathbb{R}^{p_{ch}}$, $W_{in} \in \mathbb{R}^{p_{in} \times m}$, and $\mu_{in} \in \mathbb{R}^{p_{in}}$ are their respective parameters. We wish to find the pair of projection functions minimizing the expected deviation between the images given by $J(\phi_{ch}, \phi_{in}) \equiv \mathbb{E}[\|\phi_{ch}(x^{ch}) - \phi_{in}(x^{in})\|^2]$, where $\mathbb{E}$ is the expectation operator.

The expected deviation can be reduced arbitrarily by setting the projections so that the images are scaled down. A trivial solution is $W_{ch} = 0$ and $W_{in} = 0$ at which the expected deviation vanishes for any dataset. To avoid trivial solutions, the size of the images is adjusted by fixing the second moment matrices, $\mathbb{E}[\phi_{ch}(x^{ch})\phi_{ch}(x^{ch})^\top]$ and $\mathbb{E}[\phi_{in}(x^{in})\phi_{in}(x^{in})^\top]$, to identity matrices.

The expectation appearing in the derivation and the second moment matrices operates according to an empirical probabilistic distribution. Supposing $n$ ligands are given, the chemical profiles are denoted by $x_1^{ch}, \ldots, x_n^{ch}$, and the interaction profiles by $x_1^{in}, \ldots, x_n^{in}$. If we define an empirical distribution as $q(x^{ch}, x^{in}) = \sum_{j=1}^n v_j \delta(x^{ch} - x_j^{ch})\delta(x^{in} - x_j^{in})$, with weights $v_1, \ldots, v_n$ whose sum is one and $\delta(\cdot)$ is the Dirac delta function, then the expected deviation is reduced to the weighted average of deviation and can be expressed as

$$J(\phi_{ch}, \phi_{in}) = \sum_{j=1}^n v_j \|\phi_{ch}(x_j^{ch}) - \phi_{in}(x_j^{in})\|^2. \qquad (1)$$

This implies that approximations are refined locally by setting the weights so that ligands dissimilar from the target ligand are given smaller weights.

The optimal projections can be computed via the generalized eigen-decomposition, as given in Algorithm 1 in the Appendix. When setting $v_j = 1/n$, the algorithm is shown to be equivalent to the classical CCA. Hence, we can say that weighted CCA is an extension of the classical CCA.

Kernelization of weighted CCA is formulated with a similarity function of chemical profiles $K_{ch}(x_i^{ch}, x_j^{ch})$ and a similarity function of interaction profiles $K_{in}(x_i^{in}, x_j^{in})$ without using the vectors themselves explicitly. These similarity functions are said to be valid kernels guaranteeing the theory of the algorithms, which map the profiles non-linearly into other (typically high-dimensional) spaces $\mathcal{H}_{ch}$ and $\mathcal{H}_{in}$, respectively, called an RKHS. Kernelized weighted CCA finds affine-transforms from RKHS to a canonical space $\mathbb{R}^m$, so that the expected deviation between images in $\mathbb{R}^m$ is minimized. If we denote the composite mapping functions by $\psi_{ch}$ and $\psi_{in}$, respectively, the optimal solution is given by $\psi_{ch}(x^{ch}) = A_{ch}^\top D_v^{1/2} \bar{k}_{ch}(x^{ch})$, $\psi_{in}(x^{in}) =$

$A_{in}^\top D_v^{1/2} \bar{k}_{in}(x^{in})$. The algorithm for computing the two matrices, $A_{ch} \in \mathbb{R}^{m \times n}$ and $A_{in} \in \mathbb{R}^{m \times n}$, is presented in Algorithm 2 in the Appendix. The functions $\bar{k}_{ch}(\cdot)$ and $\bar{k}_{in}(\cdot)$ are called the empirical kernel mapping, and their definition is as given in Eq. (A.7) in the Appendix.

### 2.4 Weighted SVM

Prediction of the interaction between ligand $i$ and target $t$ is performed with the SVM score given by $f(x_i^{ch}; w_{(i,t)}, b_{(i,t)}) = w_{(i,t)}^\top \psi_{ch}(x_i^{ch}) + b_{(i,t)}$, where $x_i^{ch}$ is the chemical profile of ligand $i$. The SVM parameters, $w_{(i,t)}$ and $b_{(i,t)}$, are obtained beforehand by the SVM learning algorithm. This is performed only with ligands whose interaction with the target $t$ is known. This study employs the similarity of ligands as weights in the learning process, as in the algorithm presented in the Appendix.

### 2.5 Weighting Schemes

Ligands are given weights in both stages of the weighted CCA and the weighted SVM. These weights are dependent on the ligand to be predicted. Larger weights are given for ligands that are more similar to the ligand of interest. In predicting the interaction of the $i$th ligand, the weight of $j$th ligand is given by the normalization of $v_j' = \frac{1}{\|\bar{k}_{ch}(x_j^{ch}) - \bar{k}_{ch}(x_i^{ch})\| + \|\bar{k}_{in}(x_j^{in}) - \bar{k}_{in}(x_i^{in})\| + \epsilon}$, where $\epsilon$ is a positive constant and set to 10 in our analysis. Normalization is done by setting $v_j = \frac{v_j'}{\sum_{k=1}^n v_k'}$ so that the sum of the weights is one.

## 3. Results

### 3.1 Experimental Setting

To illustrate the effectiveness of the kernel weighted CCA (KWCCA), we carried out experiments on an interaction dataset of GPCRs and odorant receptors described in the previous section. For evaluation of prediction performance, we applied a 10-fold Monte-Carlo cross validation, where data is randomly divided into 2 disjoint sets of training and test data for 10 repetitions. Data was partitioned such that for each target protein, 50% of the positive and negative interactions are used for training, and the other half for testing. KCCA, KWCCA, and the weighted SVM were implemented in Matlab, and LIBSVM [8] was used for the classical SVM.

We also performed prediction using SVM in the global model setting for comparison. The kernel function for the global model here is defined as the product of the inner product among chemical profiles and the inner product among columns of the interaction matrix.

Parameters of the local models are determined by finding respective values where the test data perform best using SVM and KCCA. Namely, the regularization parameter $C$ and the kernel function for SVM are chosen so that SVM achieves the highest prediction performance, while the regularization parameters for CCA $\gamma_{ch}$ and $\gamma_{in}$, and the number of dimensions of the canonical space $m$, are determined via the performance of KCCA. As a result, the values of the parameters are set as $C = 1000$, $\gamma_{ch} = \gamma_{in} = 1$, and $m = 4$. The RBF kernel is applied and the kernel width is determined as the mean of the distance within sets. These mentioned parameters are then fed into the algorithm em-

**Table 1** Abbreviation of methods.

| Abbreviation | Description |
|---|---|
| WW | KWCCA + Weighted SVM |
| WU | KWCCA + Classical SVM |
| KW | Classical KCCA + Weighted SVM |
| KU | Classical KCCA + Classical SVM |
| S | SVM of local model |
| SGL | Linear SVM of global model |
| SGR | RBF SVM of global model |

ploying KWCCA. The parameters are not tuned specifically for KWCCA. Thus, it is believed that there is a chance of improvement in the perfomance of this algorithm if careful and suitable parameter selection is done.

For the global model, the kernel which achieves the best performance is the linear kernel. The regularization parameter is chosen as $C = 10$, achieving the best performance among other values. Results for the case of the RBF kernel with the best $C$ value obtained are also reported for comparison.

The methods based on KWCCA involve two stages upon implementation. First, we exploit KWCCA to extract a set of features for each compound. Second, we use them for training a machine learning algorithm employing SVMs before testing them to make predictions. In total, seven methods are implemented in the experiments: two using SVM in the global model setting, and the other five following the local model. One of the two global model methods uses RBF kernel for SVM, and the other uses the linear kernel. On the other hand, the methods used for the local models are as follows: SVM, KCCA with classical SVM, KCCA with weighted SVM, KWCCA with classical SVM, and KWCCA with weighted SVM. For simplicity of notation, we shall refer to each of the seven methods using the abbreviations in **Table 1**.

### 3.2 Performance Evaluation Criteria

The following criteria were used to compare the seven prediction methods:

( 1 ) Area under the ROC curve (AUC) – The receiver operating characteristic (ROC) curve is a plot of the true positive rate (TPR) versus the false positive rate (FPR) where TPR $= \frac{\text{TP}}{\text{TP+FN}}$, FPR $= \frac{\text{FP}}{\text{FP+TN}}$, and TP, FN, FP, and TN are the number of true positives, false negatives, false positives, and true negatives, respectively. For performance comparison using the ROC, the AUC value is further computed.

( 2 ) F-measure – A value which is given by the harmonic mean of precision and recall: $F = \frac{2\text{Prec} \times \text{Recall}}{\text{Prec+Recall}}$ where the precision Prec and the recall Recall are defined by Prec $= \frac{\text{TP}}{\text{TP+FP}}$, Recall $= \frac{\text{TP}}{\text{TP+FN}}$, respectively. Since the problem is presented as a binary classification problem, only the maximum value of the F-measure values for each target is considered. The scores obtained via SVM are used as confidence levels, thus, changing the threshold yields different predictions.

These values are calculated for each target protein and averaged over the ten data divisions. However, there are instances when the test set does not contain a true positive interaction, hence AUC and F-measures cannot be computed. Therefore, these values were disregarded and, out of 62 target proteins, AUC and F-measures were computed for 49 of them. The Wilcoxon signed test was used for the statistical significance of the differ-
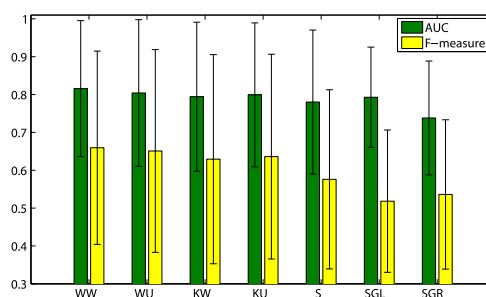


**Fig. 3** Average performance of the methods. Data was randomly split into training and test sets, and 10 training-testing data divisions were used for each method. Following the local model, AUC and F-measure were computed for each of the 62 targets. The bar plots represent the average AUC (green) and average F-measure (yellow) over the 10 cross validation sets and the 49 targets containing true positives. The two KWCCA-based methods, WW and WU, and the other methods were implemented for comparison. The difference of the performances of WW and WU from the other five methods showed to be statistically significant in terms of the P-values (by Wilcoxon signed rank test).

ence among the values of the evaluation measures.
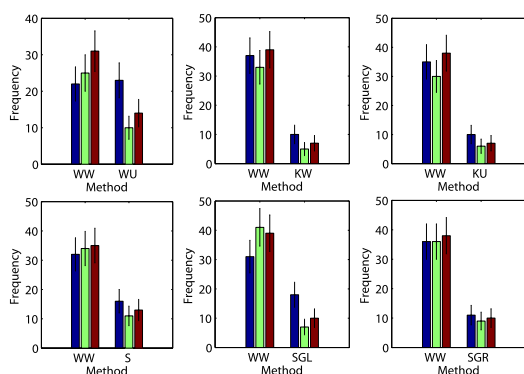
### 3.3 Effects of the Use of CCA

The average AUCs and F-measures are reported in **Fig. 3**. Error bars are also included to present standard deviations. In comparison with the local models, four CCA-based methods, WW, WU, KW, and KU, achieve remarkably better AUCs and F-measures compared to those of S: the differences between the AUCs and F-measures of KW, the worst among the four CCA-based methods, and S are 0.014 and 0.053, respectively, (P-values: $5.81 \times 10^{-7}$ and $9.49 \times 10^{-9}$ respectively). The AUC of the global model SGL is comparable to some of the local models, whereas the F-measure is not worse than that of S. A closer inspection on the results of SGL indicate that it has the lowest average number of true positives over all cross-validations among all models, around 161, which may be the reason behind a very small F-measure value.
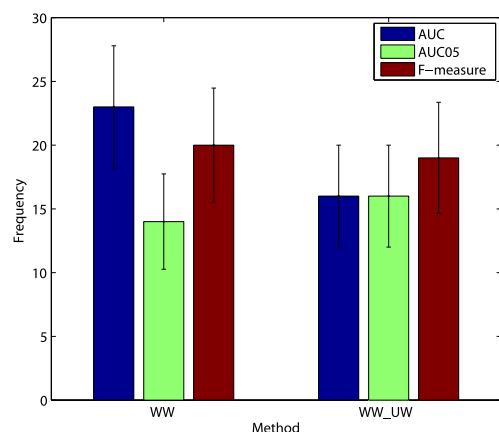
### 3.4 Improvement by Weighting

The effects of the weighted extension of CCA are manifested via comparison among four CCA-based methods. WW achieves significantly higher AUC and F-measures in average compared to KW and KU, where the P-values for the difference in the AUCs are $4.85 \times 10^{-11}$ and $6.91 \times 10^{-10}$, and the P-values for the F-measures are $3.59 \times 10^{-7}$ and $3.96 \times 10^{-6}$, respectively.

### 3.5 Histogram Comparison

The frequencies of WW besting the AUC or F-measure values of the other methods in predicting interactions for a certain target protein are shown in the histograms in **Fig. 4** (a). These values represent the number of target proteins such that the evaluated AUC and F-measure values for the method WW is better than the AUC and F-measure values of the other method in comparison. Instances when there are ties between the methods were unaccounted. For the evaluated AUC and F-measure values, WW outputs are more desirable than most of the others which indicates higher quality of prediction performance.

(a) Histogram comparisons of the proposed method WW vs. other methods.



(b) Histogram comparisons of Weighted SVM using weighting scheme (2) vs. Classical SVM.

**Fig. 4** Histogram comparisons between the proposed method WW and other methods. Frequencies when the AUC (blue), AUC between 0 and 0.05 (green), and F-measure (red) values of WW outperform the other methods, and vice-versa, are illustrated. It can be observed that AUC and F-measure values histograms for WW are more desirable than the rest.

### 3.6 Weighted and Classical SVM

WU yields interesting results in the histogram (Fig. 4 (a)): The frequency of WW yielding better AUCs are comparable to that of WU's, although frequency of better F-measures are relatively higher for WW than WU. To further investigate the comparison between WW and WU, we compute the area under the curve of the region of FPR between 0 and 0.05. This area, which we shall refer to as AUC05, allows us to evaluate the true positive rate with higher confidence. The histogram on AUC05 shows WW bests WU more frequently than WU does, which implies the use of weights in the SVM stage can find more true positives confidently than the classical SVM.

The motivation to endow the weights with training data in SVM learning is that the projections in the canonical space from chemical profiles with larger weights are expected to be better approximations of the projections from interaction profiles. It is possible to directly evaluate how good the approximations are by computing the distances among the projections. This motivation leads to another weighting scheme using the normalization of

$$v'_j = \frac{1}{\left\| \phi_{\mathrm{ch}}(x_j^{\mathrm{ch}}) - \phi_{\mathrm{in}}(x_j^{\mathrm{in}}) \right\| + \epsilon} \tag{2}$$

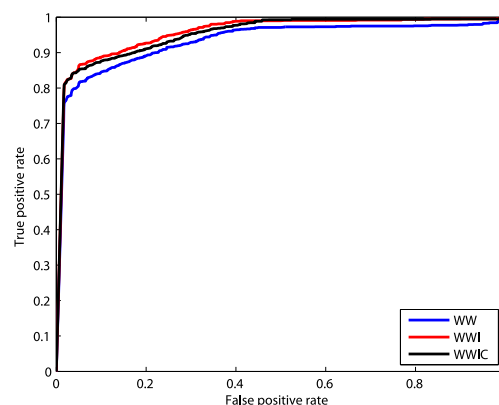instead of Eq. (1) in the SVM learning stage. We investigate the



**Fig. 5** ROC curves. WWI uses the projections from interaction profiles in the SVM stage, and WWIC uses the projections from both chemical and interaction profiles for SVM.

performance when the weighting scheme is changed to Eq. (2) in the SVM learning stage. We refer to this approach as WW$_{\mathrm{UW}}$ hereinafter. The average AUC and F-measure of WW$_{\mathrm{UW}}$ are 0.802 and 0.649, respectively, which are slightly worse than those of WW. The number of target proteins, for which the prediction performance of WW$_{\mathrm{UW}}$ is better than that of WW is not larger than the number of WW besting WW$_{\mathrm{UW}}$, as depicted in Fig. 4 (b). These facts imply that the changing weighting scheme in SVM learning does not achieve significant improvements.

### 3.7 Using Interaction Profiles

When a sufficient number of known positive and negative interactions are given for a certain ligand, the image of the interaction profile in the canonical image can provide good descriptors for predicting the remaining interactions. We further implemented two methods, herein referred to as WWI and WWIC, to investigate the performance of the interaction profile. WWI replaces the image of a chemical profile with the image of the interaction profile in the SVM stage, while WWIC concatenates the two images to feed them to the weighted SVM. The two methods achieved significant improvement. WWI achieved an average AUC of 0.857 and average F-measure of 0.699, while WWIC obtained a 0.835 average AUC and a 0.692 average F-measure. The P-values of the differences on AUC from WW are $5.27 \times 10^{-9}$ and 0.021, respectively, and P-values on F-measures are $1.05 \times 10^{-5}$ and $9.17 \times 10^{-7}$, respectively. **Figure 5** compares the average ROC curves of WW, WWI, and WWIC. The curves of WWI and WWIC are higher than that of WW, which supports the claim that introducing the interaction profiles improves the prediction performance.

## 4. Conclusions

A kernel version of weighted canonical correlation analysis is proposed, which is implemented using a derived form of the generalized eigenvalue problem. Similar to the linear CCA and its kernelized version, this can be applied to machine learning problems for dimension reduction and feature extraction. The paper presents an application to improving the prediction quality obtained in the protein-ligand interaction problem setting. By adding bias to more similar samples, better prediction can be

made which is evident on the higher AUC and F-measure values obtained.  Weighting scheme on SVM based on CCA outputs were also explored and are judged to be better than classical SVM.

Even in the field of computational biology, CCA for more than two data sources has been widely used [19], [24], [29] and their usual objectives involve maximizing the sum of correlations for every pair of data sources.  For future work, it could be worth exploring the extension of weighted CCA for analysis of multiple data sets in a biological setting.  It could also be interesting to investigate the effectiveness of applying the proposed method to other biological problems aside from protein-ligand interaction prediction.

## References

[1] Akaho, S: Kernel method for canonical correlation analysis, *Proc. International Meeting of Psychometric Society*, Osaka, Japan (2001).

[2] Antes, I., Siu, S.W. and Lengauer, T.: DynaPred: A structure and sequence based method for the prediction of MHC class I binding peptide sequences and conformations, *Bioinformatics*, Vol.22, No.14, pp.e16–24 (July 2006).

[3] Bajorath, J.: Computational analysis of ligand relationships within target families, *Curr. Opin. Chem. Biol.*, Vol.12, No.3, pp.352–358 (June 2008).

[4] Ballester, P.J. and Mitchell, J.B.: A machine learning approach to predicting protein-ligand binding affinity with applications to molecular docking, *Bioinformatics*, Vol.26, No.9, pp.1169–1175 (May 2010).

[5] Biniashvili, T., Schreiber, E. and Kliger, Y.: Improving classical substructure-based virtual screening to handle extrapolation challenges, *J. Chem. Inf. Model*, Vol.52, No.3, pp.678–685 (Mar. 2012).

[6] Bleakley, K. and Yamanishi, Y.: Supervised prediction of drug-target interactions using bipartite local models, *Bioinformatics*, Vol.25, No.18, pp.2397–2403 (Sep. 2009).

[7] Campillos, M., Kuhn, M., Gavin, A.-C., Jensen, L.J. and Bork, P.: Drug target identification using side-effect similarity, *Science*, Vol.321, pp.263–266 (2008).

[8] Chang, C.-C. and Lin, C.-J.: LIBSVM: A library for support vector machines, *ACM Trans. Intelligent Systems and Technology*, Vol.2, No.27, pp.1–27 (2011).

[9] Chen, B., Harrison, R.F., Papadatos, G., Willett, P., Wood, D.J., Lewell, X.Q., Greenidge, P. and Stiefl, N.: Evaluation of machine-learning methods for ligand-based virtual screening, *J. Comput. Aided Mol. Des.*, Vol.21, No.1-3, pp.53–62 (Jan.-Mar. 2007).

[10] Golugula, A., Lee, G., Master, S.R.. Feldman, M.D., Tomaszewski, J.E. and Madabhushi, A.: Supervised regularized canonical correlation analysis: Integrating histologic and proteomic data for predicting biochemical failures, *Conf. Proc. IEEE Eng. Med. Biol. Soc.*, pp.6434–6437 (2011).

[11] Gonzalez, A.J., Liao, L. and Wu, C.H.: Predicting ligand binding residues and functional sites using multi-positional correlations with graph theoretic clustering and kernel CCA, *IEEE/ACM Trans. Comput. Biol. Bioinform.* (Oct. 2011).

[12] Hardoon, D.R., Szedmak, S. and Shawe-Taylor, J.: Canonical correlation analysis: An overview with application to learning methods, *Neural Comput.*, Vol.16, No.12, pp.2639–2664 (Dec. 2004).

[13] Hotelling, H.: Relation between two sets of variates, *Biometrika*, Vol.28, No.3 (1936).

[14] Jacob, L. and Vert, J.P.: Protein-ligand interaction prediction: An improved chemogenomics approach, *Bioinformatics*, Vol.24, No.19, pp.2149–2156 (Oct. 2008).

[15] Klabunde, T.: Chemogenomic approaches to drug discovery: Similar receptors bind similar ligands, *Br. J. Pharmacol.*, Vol.152, pp.5–7 (2007).

[16] Martin, Y.C., Kofron, J.L. and Traphagen, L.M.: Do structurally similar molecules have similar biological activity, *J. Med. Chem.*, Vol.45, pp.4350–4358 (2002).

[17] Paolini, G.V., Shapland, R.H.B., van Hoorn, W.P., Mason, J.S. and Hopkins, A.L.: Global mapping of the pharmacological space, *Nat.*

*Biotechnol.*, Vol.24, pp.805–815 (2006).

[18] Pazos, F., Rausell, A. and Valencia, A.: Phylogeny-independent detection of functional residues, *Bioinformatics*, Vol.22, No.12, pp.1440–1448 (June 2006).

[19] Peng, Y., Zhang, D. and Zhang, J.: A new canonical correlation analysis algorithm with local discrimination, *Neural Process Lett.*, Vol.31, pp.1–15 (2010).

[20] Rognan, D.: Chemogenomic approaches to rational drug design, *Br. J. Pharmacol.*, Vol.152, No.1, pp.38–52 (Sep. 2007).

[21] Saito, H., Kubota, M., Roberts, R.W., Chi, Q. and Matsunami, H.: RTP family members induce functional expression of mammalian odorant receptors, *Cell*, Vol.119, No.5, pp.679–691 (Nov. 2004).

[22] Samarov, D., Marron, J.S., Liu, Y., Grulke, C. and Tropsha, A.: Local kernel canonical correlation analysis with application to virtual drug screening, *Ann. Appl. Stat.*, Vol.5, No.3, pp.2169–2196 (Sep. 2011).

[23] Schuffenhauer, A., Floersheim, P., Acklin, P. and Jacoby, E.: Similarity metrics for ligands reflecting the similarity of the target proteins, *J. Chem. Inf. Comput. Sci.*, Vol.43, pp.391–405 (2003).

[24] Tang, C.S. and Ferreira, M.A.: A gene-based test of association using canonical correlation analysis, *Bioinformatics*, Vol.28, No.6. pp.845–850 (Mar. 2012).

[25] van Laarhoven, T., Nabuurs, S.B. and Marchiori, E.: Gaussian interaction profile kernels for predicting drug-target interaction, *Bioinformatics*, Vol.27, No.21, pp.3036–3043 (Nov. 2011).

[26] Wassermann, A.M., Geppert, H. and Bajorath, J.: Ligand prediction for orphan targets using support vector machines and various target-ligand kernels is dominated by nearest neighbor effects, *J. Chem. Inf. Model.*, Vol.49, pp.2155–2167 (2009).

[27] Witten, D.M., Tibshirani, R. and Hastie, T.: A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis, *Biostatistics*, Vol.10, No.3, pp.515–534 (July 2009).

[28] Yamanishi, Y., Pauwels, E., Saigo, H. and Stoven, V.: Extracting sets of chemical substructures and protein domains governing drug-target interactions, *J. Chem. Inf. Model.* (May 2011).

[29] Yamanishi, Y., Vert, J.P., Nakaya, A. and Kanehisa, M.: Extraction of correlated gene clusters from multiple genomic data by generalized kernel canonical correlation analysis, *Bioinformatics*, Vol.19, Suppl.1, pp.i323–330 (2003).

[30] Yan, X. and Han, J.: gSpan: Graph-based substructure pattern mining, *Proc. 2002 Int'l Conf. Data Mining* (*ICDM '02*), pp.721–724 (2002).

[31] Yu, S., De Moor, B. and Moreau, Y.: Learning with heterogenous data sets by weighted multiple kernel canonical correlation analysis, *Proc. IEEE Machine Learning for Signal Processing* (*IEEE MLSP 2008*), pp.81–86, Thessaloniki, Greece (Aug. 2007).

# Appendix

Here we present the details of the supporting theories applied in this study.  Throughout the discussion, we shall use the following notations. We will denote vectors by boldface lower-case letters and matrices by boldface upper-case letters. The inverse of a matrix $A$ is denoted by $A^{-1}$ and its transpose is $A^{\top}$. The $n \times n$ identity matrix is given by $I_n$, while $D_v$ is defined as diag($v$), the matrix whose diagonal entries are from vector $v$. We let $\mathbb{S}^n$ be the set of all $n \times n$ symmetric matrices, $\mathbb{S}^n_+$ as the set of all $n \times n$ symmetric positive semi-definite matrices, $\mathbb{S}^n_{++}$, the set of all $n \times n$ symmetric positive definite matrices, and $\mathbb{O}^{m \times n}$ as the set of all $m \times n$ orthonormal matrices. We define $\mathbb{N}_n = \{i \in \mathbb{N} | i \leq n\}$, and $\langle \cdot, \cdot \rangle$ is used to denote the inner product between vectors.

## A.1   Generalized Eigendecomposition

The generalized eigendecomposition is defined as follows:
**Theorem 1.** (Generalized Eigendecomposition) Let $n \in \mathbb{N}$. For any $A \in \mathbb{S}^n$ and any $B \in \mathbb{S}^n_{++}$,

$$\exists U \in \mathbb{R}^{n \times n}, \exists \lambda \in \mathbb{R}^n, \text{ s.t. } \quad AU = BUD_\lambda,$$
$$U^{\top}BU = I_n.$$

The entries of $\lambda$ are called *generalized eigenvalues*, and the columns of $U$ are called *generalized eigenvectors*.

This section deals with the case when the two matrices $A \in \mathbb{S}^{n+m}$ and $B \in \mathbb{S}^{n+m}_{++}$ are of the form

$$A = \begin{bmatrix} \mathbf{0}_{m \times m} & C \\ C^\top & \mathbf{0}_{n \times n} \end{bmatrix}, \quad B = \begin{bmatrix} B_x & \mathbf{0}_{m \times n} \\ \mathbf{0}_{n \times m} & B_y \end{bmatrix}$$

where $C \in \mathbb{R}^{m \times n}$ with $r \equiv \text{rank}(C)$. Let us denote the generalized eigendecomposition of $(A, B)$ by

$$AU_{\text{all}} = BU_{\text{all}}\Lambda_{\text{all}},$$

where $U_{\text{all}}^\top B U_{\text{all}} = I_{m+n}$, and $\Lambda_{\text{all}} = \text{diag}(\lambda_{\text{all}})$ and $\lambda_{\text{all}} = [\lambda_1, \ldots, \lambda_{m+n}]^\top$, such that $\lambda_1 \geq \cdots \geq \lambda_{m+n}$. We denote the columns in $U_{\text{all}}$ by

$$U_{\text{all}} = [u_1, \ldots, u_{m+n}] = \begin{bmatrix} u_{x,1}, \cdots, u_{x,m+n} \\ u_{y,1}, \cdots, u_{y,m+n} \end{bmatrix}$$

where $(\forall i \in \mathbb{N}_{m+n})\, u_{x,i} \in \mathbb{R}^m$, $u_{y,i} \in \mathbb{R}^n$, and define

$$U_x \equiv [u_{x,1}, \ldots u_{x,r}], \quad U_y \equiv [u_{y,1}, \ldots, u_{y,r}].$$

The following theorem is the main result of this section.
**Theorem 2.** Consider the following optimization problem

$$\begin{aligned} \max \quad & \text{tr}(X^\top C Y) \\ \text{wrt} \quad & X \in \mathbb{R}^{m \times k}, Y \in \mathbb{R}^{n \times k}, \\ \text{subj to} \quad & X^\top B_x X = Y^\top B_y Y = I_k. \end{aligned}$$

An optimal solution is given by

$$X = \sqrt{2}[u_{x,1}, \ldots, u_{x,k}], \quad Y = \sqrt{2}[u_{y,1}, \ldots, u_{y,k}].$$

To prove Theorem 2, we will use the following lemma.
**Lemma 1.**

$$U_x^\top B_x U_x = U_y^\top B_y U_y = \frac{1}{2} I_r.$$

*Proof.* (of Lemma 1) Let

$$\Lambda \equiv \text{diag}\{\lambda_1, \ldots, \lambda_r\}.$$

We assume $\lambda_r > 0$ here. From this assumption,

$$CU_y = B_x U_x \Lambda, \quad C^\top U_x = B_y U_y \Lambda.$$

Pre-multiplying the former equation by $U_x^\top$ and post-multiplying the transpose of the latter equation by $U_y$ yield

$$U_x^\top B_x U_x \Lambda = U_x^\top C U_y = \Lambda U_y^\top B_y U_y.$$

For the diagonal entries of the above equality,

$$\forall i \in \mathbb{N}_r : \lambda_i u_{x,i}^\top B_x u_{x,i} = \lambda_i u_{y,i}^\top B_y u_{y,i},$$

resulting in

$$u_{x,i}^\top B_x u_{x,i} = u_{y,i}^\top B_y u_{y,i} = \frac{1}{2} \tag{A.1}$$

since $\lambda_i > 0$ and from the assumption that $U_{\text{all}}^\top B U_{\text{all}} = I_{m+n}$. For the off-diagonal entries,

$$u_{x,i}^\top B_x u_{x,j} - \frac{\lambda_i}{\lambda_j} u_{y,i}^\top B_y u_{y,j} = 0. \tag{A.2}$$

From the assumption $U_{\text{all}}^\top B U_{\text{all}} = I_{m+n}$,

$$\begin{bmatrix} u_{x,i} \\ u_{y,i} \end{bmatrix}^\top B \begin{bmatrix} u_{x,j} \\ u_{y,j} \end{bmatrix} = u_{x,i}^\top B_x u_{x,j} + u_{y,i}^\top B_y u_{y,j} = 0. \tag{A.3}$$

From the two equations (A.2) and (A.3), we have

$$u_{x,i}^\top B_x u_{x,j} = u_{y,i}^\top B_y u_{y,j} = 0. \tag{A.4}$$

Thus, Eqs. (A.1) and (A.4) establish Lemma 1. □

*Proof.* (of Theorem 2) Let

$$Z \equiv \frac{1}{\sqrt{2}} \begin{bmatrix} X \\ Y \end{bmatrix}.$$

There exists $R \in \mathbb{R}^{(m+n) \times k}$ such that $Z = U_{\text{all}} R$, and matrix $R$ is orthonormal, i.e., $R \in \mathbb{O}^{(m+n) \times k}$, since

$$\begin{aligned} R^\top R &= R^\top I_{(m+n)} R = R^\top U_{\text{all}}^\top B U_{\text{all}} R = Z^\top B Z \\ &= \frac{1}{2} X^\top B_x X + \frac{1}{2} Y^\top B_y Y = I_k. \end{aligned}$$

Let $r_i \in \mathbb{R}^k$ denote the $i$th row vector, i.e.,

$$R = \begin{bmatrix} r_1^\top \\ \vdots \\ r_{m+n}^\top \end{bmatrix}.$$

Then we have

$$\begin{aligned} \text{tr}(X^\top C Y) &= \frac{1}{2}\text{tr}(X^\top C Y) + \frac{1}{2}\text{tr}(Y^\top C^\top X) \\ &= \frac{1}{2}\text{tr}\left( [X^\top, Y^\top] \begin{bmatrix} \mathbf{0}_{m \times m} & C \\ C^\top & \mathbf{0}_{n \times n} \end{bmatrix} \begin{bmatrix} X \\ Y \end{bmatrix} \right) \\ &= \text{tr}(Z^\top A Z) = \text{tr}(R^\top U_{\text{all}}^\top A U_{\text{all}} R) \\ &= \text{tr}(R^\top U_{\text{all}}^\top B U_{\text{all}} D_{\lambda_{\text{all}}} R) = \text{tr}(R^\top D_{\lambda_{\text{all}}} R) \\ &= \sum_{i=1}^{m+n} \lambda_i \|r_i\|^2 = \langle w, \lambda_{\text{all}} \rangle, \tag{A.5} \end{aligned}$$

where $w \in \mathbb{R}_+^{m+n}$ is a nonnegative vector in which the $i$th entry is defined by $w_i = \|r_i\|^2$. To conclude the proof, we shall make use of the following two properties:

First, note that there exists $S \in \mathbb{O}^{(m+n) \times (m+n-k)}$ such that $R^\top S = \mathbf{0}_{k \times (m+n-k)}$. Then $[R, S] \in \mathbb{O}^{(m+n) \times (m+n)}$. Let $s_i \in \mathbb{R}^{(m+n-k)}$ denote the $i$th row vector, i.e.,

$$S = \begin{bmatrix} s_1^\top \\ \vdots \\ s_{m+n}^\top \end{bmatrix}.$$

Then, $\forall i \in \mathbb{N}_n$ we have

$$1 = \left\| \begin{bmatrix} r_i \\ s_i \end{bmatrix} \right\|^2 = \|r_i\|^2 + \|s_i\|^2 \geq \|r_i\|^2 = w_i,$$

where the first equality follows from the property of the square orthonormal matrix: $[R, S][R, S]^\top = I_n$. Second, observe that

$$\|w\|_1 = \sum_{i=1}^n w_i = \sum_{i=1}^n \|r_i\|^2 = \text{tr}(RR^\top) = \text{tr}(R^\top R) = \text{tr} I_k = k.$$

Now, the objective function $\text{tr}(X^\top C Y)$ is maximized when

$$w_i = \begin{cases} 1, & \text{if } 1 \le i \le k, \\ 0, & \text{if } k < i \le n, \end{cases}$$

and so the value of Eq. (A.5) is bounded above by

$$\langle w, \lambda \rangle \le \sum_{i=1}^{k} \lambda_i.$$

Finally, we show that equality holds when

$$R = \begin{bmatrix} I_k \\ 0_{(m+n-k) \times k} \end{bmatrix}.$$

Since $Z = UR = U_{\text{ma}}$, where $U_{\text{ma}} \equiv [u_1, \ldots, u_k]$, we have

$$\text{tr}(Z^\top A Z) = \text{tr}(U_{\text{ma}}^\top A U_{\text{ma}}) = \sum_{i=1}^{k} \lambda_i.$$

Thus, $Z = U_{\text{ma}}$ is an optimum, which implies

$$X = \sqrt{2}[u_{x,1}, \ldots, u_{x,k}], \quad Y = \sqrt{2}[u_{y,1}, \ldots, u_{y,k}].$$

$\square$

## A.2   Linear Weighted CCA

**Algorithm 1** (Linear Weighted CCA) Define two matrices,

$$X_{\text{ch}} \equiv [x_1^{\text{ch}}, \ldots, x_n^{\text{ch}}], \qquad X_{\text{in}} \equiv [x_1^{\text{in}}, \ldots, x_n^{\text{in}}],$$

and let $v = [v_1, \ldots, v_n]^\top$. Then the optimal offsets, $\mu_{\text{ch}}, \mu_{\text{in}}$, are computed as

$$\mu_{\text{ch}} = X_{\text{ch}} v, \qquad \mu_{\text{in}} = X_{\text{in}} v.$$

We use the optimal offsets to define $C_{\text{ch,ch}}, C_{\text{ch,in}}, C_{\text{in,ch}}, C_{\text{in,in}}$ as

$$C_{\text{ch,ch}} \equiv X_{\text{ch}} D_v X_{\text{ch}}^\top - \mu_{\text{ch}} \mu_{\text{ch}}^\top, \quad C_{\text{ch,in}} \equiv X_{\text{ch}} D_v X_{\text{in}}^\top - \mu_{\text{ch}} \mu_{\text{in}}^\top,$$
$$C_{\text{in,ch}} \equiv X_{\text{in}} D_v X_{\text{ch}}^\top - \mu_{\text{in}} \mu_{\text{ch}}^\top, \quad C_{\text{in,in}} \equiv X_{\text{in}} D_v X_{\text{in}}^\top - \mu_{\text{in}} \mu_{\text{in}}^\top,$$

and consider the following generalized eigen-decomposition problem:

$$\begin{bmatrix} 0_{p_{\text{ch}} \times p_{\text{ch}}} & C_{\text{ch,in}} \\ C_{\text{in,ch}} & 0_{p_{\text{in}} \times p_{\text{in}}} \end{bmatrix} \begin{bmatrix} w^{\text{ch}} \\ w^{\text{in}} \end{bmatrix} = \begin{bmatrix} C_{\text{in,in}} & 0_{p_{\text{in}} \times p_{\text{ch}}} \\ 0_{p_{\text{ch}} \times p_{\text{in}}} & C_{\text{ch,ch}} \end{bmatrix} \begin{bmatrix} w^{\text{ch}} \\ w^{\text{in}} \end{bmatrix}.$$

Denote the $h$th major eigen-vector by $\begin{bmatrix} w_h^{\text{ch}} \\ w_h^{\text{in}} \end{bmatrix}$. The optimal loading matrices $W_{\text{ch}}$ and $W_{\text{in}}$ are computed by setting the $h$th columns of $W_{\text{ch}}$ and $W_{\text{in}}$ to $w_h^{\text{ch}}$ and $w_h^{\text{in}}$, respectively.

**Theorem 3.** Algorithm 1 yields the parameters of the mapping functions $(\phi_{\text{ch}}, \phi_{\text{in}})$ which minimize the expected deviation subject to the scaling constraints that the second moment matrices are the identity matrix.

To prove this theorem, we employ the following result.

Let two affine mapping functions $\phi_{\text{ch}} : \mathbb{R}^{p_{\text{ch}}} \to \mathbb{R}^m$ and $\phi_{\text{in}} : \mathbb{R}^{p_{\text{in}}} \to \mathbb{R}^m$ be parametrized by

$$\phi_{\text{ch}}(x_{\text{ch}}) = W_{\text{ch}}^\top(x^{\text{ch}} - \mu_{\text{ch}}), \quad \text{and} \quad \phi_{\text{in}}(x^{\text{in}}) = W_{\text{in}}^\top(x^{\text{in}} - \mu_{\text{in}}),$$

respectively. Define the expected deviation by

$$J = \mathbb{E}_{p(x^{\text{ch}}, x^{\text{in}})}[\|\phi_{\text{ch}}(x^{\text{ch}}) - \phi_{\text{in}}(x^{\text{in}})\|^2],$$

where the expectation is according to a probabilistic distribution $p(x^{\text{ch}}, x^{\text{in}})$. If we consider the optimization problem for minimizing the expected deviation with respect to the parameters $(W_{\text{ch}}, \mu_{\text{ch}}, W_{\text{in}}, \mu_{\text{in}})$ subject to

$$\mathbb{E}_{p(x^{\text{ch}})}[\phi_{\text{ch}}(x^{\text{ch}})\phi_{\text{ch}}(x^{\text{ch}})^\top] = \mathbb{E}_{p(x^{\text{in}})}[\phi_{\text{in}}(x^{\text{in}})\phi_{\text{in}}(x^{\text{in}})^\top] = I_m,$$

a minimizer of the optimization problem is found by the following optimization problem.

Set the offset vectors as $\mu_{\text{ch}} = \mathbb{E}_{p(x^{\text{ch}})}[x^{\text{ch}}]$ and $\mu_{\text{in}} = \mathbb{E}_{p(x^{\text{in}})}[x^{\text{in}}]$. The optimal loading matrices $W_{\text{ch}}$ and $W_{\text{in}}$ are computed by setting the $h$th columns of $W_{\text{ch}}$ and $W_{\text{in}}$ to $w_h^{\text{ch}}$ and $w_h^{\text{in}}$, respectively, and denoting by $\begin{bmatrix} w_h^{\text{ch}} \\ w_h^{\text{in}} \end{bmatrix}$ the $h$th major eigenvector of the generalized eigendecomposition:

$$\begin{bmatrix} 0_{p_{\text{ch}} \times p_{\text{ch}}} & \check{C}_{\text{ch,in}} \\ \check{C}_{\text{in,ch}} & 0_{p_{\text{in}} \times p_{\text{in}}} \end{bmatrix} \begin{bmatrix} w^{\text{ch}} \\ w^{\text{in}} \end{bmatrix}$$
$$= \begin{bmatrix} \check{C}_{\text{in,in}} & 0_{p_{\text{in}} \times p_{\text{ch}}} \\ 0_{p_{\text{ch}} \times p_{\text{in}}} & \check{C}_{\text{ch,ch}} \end{bmatrix} \begin{bmatrix} w^{\text{ch}} \\ w^{\text{in}} \end{bmatrix}, \tag{A.6}$$

where the covariance matrices are given by

$$\check{C}_{\text{ch,ch}} \equiv \mathbb{E}_{p(x^{\text{ch}})}[(x^{\text{ch}} - \mu_{\text{ch}})(x^{\text{ch}} - \mu_{\text{ch}})^\top],$$
$$\check{C}_{\text{ch,in}} \equiv \mathbb{E}_{p(x^{\text{ch}}, x^{\text{in}})}[(x^{\text{ch}} - \mu_{\text{ch}})(x^{\text{in}} - \mu_{\text{in}})^\top],$$
$$\check{C}_{\text{in,ch}} \equiv \check{C}_{\text{ch,in}}^\top,$$
$$\check{C}_{\text{in,in}} \equiv \mathbb{E}_{p(x^{\text{in}})}[(x^{\text{in}} - \mu_{\text{in}})(x^{\text{in}} - \mu_{\text{in}})^\top].$$

To verify this, we let

$$x \equiv \begin{bmatrix} x^{\text{ch}} \\ x^{\text{in}} \end{bmatrix}, \qquad \mu_{\text{tot}} \equiv \begin{bmatrix} \mu_{\text{ch}} \\ \mu_{\text{in}} \end{bmatrix}.$$

Then the second order moment matrices are rewritten as

$$\mathbb{E}[\phi_{\text{ch}}(x^{\text{ch}})\phi_{\text{ch}}(x^{\text{ch}})^\top] = W_{\text{ch}}^\top \mathbb{E}[(\mu_{\text{ch}} - x^{\text{ch}})(\mu_{\text{ch}} - x^{\text{ch}})^\top]W_{\text{ch}},$$
$$\mathbb{E}[\phi_{\text{in}}(x^{\text{in}})\phi_{\text{in}}(x^{\text{in}})^\top] = W_{\text{in}}^\top \mathbb{E}[(\mu_{\text{in}} - x^{\text{in}})(\mu_{\text{in}} - x^{\text{in}})^\top]W_{\text{in}},$$

respectively, and the expected deviation is arranged as

$$\begin{aligned} J &= \mathbb{E}[\|\phi_{\text{ch}}(x^{\text{ch}}) - \phi_{\text{in}}(x^{\text{in}})\|^2] \\ &= \mathbb{E}[\|W_{\text{ch}}^\top(x^{\text{ch}} - \mu_{\text{ch}}) - W_{\text{in}}^\top(x^{\text{in}} - \mu_{\text{in}})\|^2] \\ &= \text{tr}(W_{\text{ch}}^\top \mathbb{E}[(x^{\text{ch}} - \mu_{\text{ch}})(x^{\text{ch}} - \mu_{\text{ch}})^\top]W_{\text{ch}}) \\ &\quad + \text{tr}(W_{\text{in}}^\top \mathbb{E}[(x^{\text{in}} - \mu_{\text{in}})(x^{\text{in}} - \mu_{\text{in}})^\top]W_{\text{in}}) \\ &\quad - 2\text{tr}(W_{\text{ch}}^\top \mathbb{E}[(x^{\text{ch}} - \mu_{\text{ch}})(x^{\text{in}} - \mu_{\text{in}})^\top]W_{\text{in}}) \\ &= 2m - 2\text{tr}(W_{\text{ch}}^\top \mathbb{E}[(x^{\text{ch}} - \mu_{\text{ch}})(x^{\text{in}} - \mu_{\text{in}})^\top]W_{\text{in}}). \end{aligned}$$

From here, we will first derive the optimal value of $\mu_{\text{tot}}$, and then give the algorithm to find the optimal $W_{\text{ch}}$ and $W_{\text{in}}$. Introducing the Lagrangian multipliers $\Lambda_{\text{ch}} \in \mathbb{S}^{p_{\text{ch}}}$ and $\Lambda_{\text{in}} \in \mathbb{S}^{p_{\text{in}}}$, the Lagrangian function is written as

$$\begin{aligned} \mathcal{L}_{\text{A}} &= 2m - 2\text{tr}(W_{\text{ch}}^\top \mathbb{E}[(\mu_{\text{ch}} - x^{\text{ch}})(\mu_{\text{in}} - x^{\text{in}})^\top]W_{\text{in}}) \\ &\quad - \langle \Lambda_{\text{ch}}, I_m - W_{\text{ch}}^\top \mathbb{E}[(\mu_{\text{ch}} - x^{\text{ch}})(\mu_{\text{ch}} - x^{\text{ch}})^\top]W_{\text{ch}} \rangle \\ &\quad - \langle \Lambda_{\text{in}}, I_m - W_{\text{in}}^\top \mathbb{E}[(\mu_{\text{in}} - x^{\text{in}})(\mu_{\text{in}} - x^{\text{in}})^\top]W_{\text{in}} \rangle. \end{aligned}$$

To obtain the values of $\mu_{\text{ch}}$ and $\mu_{\text{in}}$ at the saddle point, we set the derivative of the Lagrangian to zero:

$$\frac{\partial \mathcal{L}_A}{\partial \boldsymbol{\mu}_{\text{ch}}} = 2W_{\text{ch}}W_{\text{in}}^\top(\boldsymbol{\mu}_{\text{in}} - \mathbb{E}[\boldsymbol{x}^{\text{in}}])$$
$$+ 2W_{\text{ch}}\Lambda_{\text{ch}}W_{\text{ch}}^\top(\boldsymbol{\mu}_{\text{ch}} - \mathbb{E}[\boldsymbol{x}^{\text{ch}}]),$$
$$\frac{\partial \mathcal{L}_A}{\partial \boldsymbol{\mu}_{\text{in}}} = 2W_{\text{in}}W_{\text{ch}}^\top(\boldsymbol{\mu}_{\text{ch}} - \mathbb{E}[\boldsymbol{x}^{\text{ch}}])$$
$$+ 2W_{\text{in}}\Lambda_{\text{in}}W_{\text{in}}^\top(\boldsymbol{\mu}_{\text{in}} - \mathbb{E}[\boldsymbol{x}^{\text{in}}]).$$

The two derivatives vanish simultaneously when

$$\boldsymbol{\mu}_{\text{ch}} = \mathbb{E}[\boldsymbol{x}^{\text{ch}}], \qquad \text{and} \qquad \boldsymbol{\mu}_{\text{in}} = \mathbb{E}[\boldsymbol{x}^{\text{in}}].$$

Next, we derive the optimal $W_{\text{ch}}$ and $W_{\text{in}}$. Substituting the definitions of the covariance matrices into the expected deviation and the second moments of the affine transformations, we have

$$J = 2m - 2\text{tr}(W_{\text{ch}}^\top \breve{C}_{\text{ch,in}} W_{\text{in}})$$

and

$$\mathbb{E}_{p(\boldsymbol{x}^{\text{ch}})}[\boldsymbol{\phi}_{\text{ch}}(\boldsymbol{x}^{\text{ch}})\boldsymbol{\phi}_{\text{ch}}(\boldsymbol{x}^{\text{ch}})^\top] = W_{\text{ch}}^\top \breve{C}_{\text{ch,ch}} W_{\text{ch}},$$
$$\mathbb{E}_{p(\boldsymbol{x}^{\text{in}})}[\boldsymbol{\phi}_{\text{in}}(\boldsymbol{x}^{\text{in}})\boldsymbol{\phi}_{\text{in}}(\boldsymbol{x}^{\text{in}})^\top] = W_{\text{in}}^\top \breve{C}_{\text{in,in}} W_{\text{in}}.$$

Then, by omitting the constants, the problem for finding the optimal $W_{\text{ch}}$ and $W_{\text{in}}$ is reduced to the following optimization problem:

$$\begin{aligned} \max \quad & \text{tr}(W_{\text{ch}}^\top \breve{C}_{\text{ch,in}} W_{\text{in}}) \\ \text{wrt} \quad & W_{\text{ch}} \in \mathbb{R}^{p_{\text{ch}} \times m}, W_{\text{in}} \in \mathbb{R}^{p_{\text{in}} \times m} \\ \text{subj to} \quad & W_{\text{ch}}^\top \breve{C}_{\text{ch,ch}} W_{\text{ch}} = W_{\text{in}}^\top \breve{C}_{\text{in,in}} W_{\text{in}} = I_m. \end{aligned}$$

From Theorem 2, the optimization problem is solved by generalized eigendecomposition (A.6) previously described. Hence, the given algorithm finds a minimizer of the optimization problem.

Now we present the proof of Theorem 3.

*Proof.* Observe that if we substitute the probabilistic distribution $p(\boldsymbol{x}^{\text{ch}}, \boldsymbol{x}^{\text{in}})$ to the empirical distribution $q(\boldsymbol{x}^{\text{ch}}, \boldsymbol{x}^{\text{in}})$ defined in the main text, the first moments and the second moments are expressed as

$$\mathbb{E}_{p(\boldsymbol{x}^{\text{ch}})}[\boldsymbol{x}^{\text{ch}}] = X_{\text{ch}}\boldsymbol{v}, \qquad \mathbb{E}_{p(\boldsymbol{x}^{\text{in}})}[\boldsymbol{x}^{\text{in}}] = X_{\text{in}}\boldsymbol{v},$$
$$\breve{C}_{\text{ch,ch}} = X_{\text{ch}}D_{\boldsymbol{v}}X_{\text{ch}}^\top - \boldsymbol{\mu}_{\text{ch}}\boldsymbol{\mu}_{\text{ch}}^\top, \quad \breve{C}_{\text{ch,in}} = X_{\text{ch}}D_{\boldsymbol{v}}X_{\text{in}}^\top - \boldsymbol{\mu}_{\text{ch}}\boldsymbol{\mu}_{\text{in}}^\top,$$
$$\breve{C}_{\text{in,ch}} = \breve{C}_{\text{ch,in}}^\top, \qquad \breve{C}_{\text{in,in}} = X_{\text{in}}D_{\boldsymbol{v}}X_{\text{in}}^\top - \boldsymbol{\mu}_{\text{in}}\boldsymbol{\mu}_{\text{in}}^\top,$$

which implies that the optimization algorithm given in the result above is equivalent to Algorithm 1 in this case. Thus, Theorem 3 is established. □

## A.3   Kernel Weighted CCA

Suppose we are given $n$ drugs. The kernel matrices of the chemical and interaction kernels $K_{\text{ch}} \in \mathbb{S}^{n \times n}$ and $K_{\text{in}} \in \mathbb{S}^{n \times n}$ are defined so that their respective entries are the values of the kernel functions, as given by

$$\forall i, \forall j \in \mathbb{N}_n: \quad K_{i,j}^{\text{ch}} = K_{\text{ch}}(\boldsymbol{x}_i^{\text{ch}}, \boldsymbol{x}_j^{\text{ch}}), \quad K_{i,j}^{\text{in}} = K_{\text{in}}(\boldsymbol{x}_i^{\text{in}}, \boldsymbol{x}_j^{\text{in}}).$$

The kernel empirical mapping of the two information sources are defined as

$$\boldsymbol{k}_{\text{ch}}(\boldsymbol{x}^{\text{ch}}) \equiv [K_{\text{ch}}(\boldsymbol{x}_1^{\text{ch}}, \boldsymbol{x}^{\text{ch}}), \dots, K_{\text{ch}}(\boldsymbol{x}_n^{\text{ch}}, \boldsymbol{x}^{\text{ch}})]^\top,$$
$$\boldsymbol{k}_{\text{in}}(\boldsymbol{x}^{\text{in}}) \equiv [K_{\text{in}}(\boldsymbol{x}_1^{\text{in}}, \boldsymbol{x}^{\text{in}}), \dots, K_{\text{in}}(\boldsymbol{x}_n^{\text{in}}, \boldsymbol{x}^{\text{in}})]^\top.$$

KWCCA needs the kernel values among shifted vectors in the kernel Hilbert spaces. The shifted kernels for chemical profiles can be computed as

$$\bar{K}_{\text{ch}}(\boldsymbol{x}_i^{\text{ch}}, \boldsymbol{x}_j^{\text{ch}}) \equiv K_{\text{ch}}(\boldsymbol{x}_i^{\text{ch}}, \boldsymbol{x}_j^{\text{ch}}) - \boldsymbol{v}^\top(\boldsymbol{k}_{\text{ch}}(\boldsymbol{x}_i^{\text{ch}}) - \boldsymbol{k}_{\text{ch}}(\boldsymbol{x}_j^{\text{ch}}))$$
$$+ \boldsymbol{v}^\top K_{\text{ch}}\boldsymbol{v},$$
$$\bar{\boldsymbol{k}}_{\text{ch}}(\boldsymbol{x}^{\text{ch}}) \equiv (I_n - 1_n\boldsymbol{v}^\top)(\boldsymbol{k}_{\text{ch}}(\boldsymbol{x}^{\text{ch}}) - K_{\text{ch}}\boldsymbol{v}),$$
$$\bar{K}_{\text{ch}} \equiv (I_n - 1_n\boldsymbol{v}^\top)K_{\text{ch}}(I_n - \boldsymbol{v}1_n^\top),$$

(A.7)

and $\bar{K}_{\text{in}}(\boldsymbol{x}_i^{\text{in}}, \boldsymbol{x}_j^{\text{in}})$, $\bar{\boldsymbol{k}}_{\text{in}}(\boldsymbol{x}^{\text{in}})$, and $\bar{K}_{\text{in}}$ are defined similarly.

If we define $\tilde{K}_{\text{ch}} \equiv D_{\boldsymbol{v}}^{1/2}\bar{K}_{\text{ch}}D_{\boldsymbol{v}}^{1/2}$, $\tilde{K}_{\text{in}} \equiv D_{\boldsymbol{v}}^{1/2}\bar{K}_{\text{in}}D_{\boldsymbol{v}}^{1/2}$, the expected deviation between the images in $\mathbb{R}^m$ is expressed as

$$\mathbb{E}[\|\psi_{\text{ch}}(\boldsymbol{x}^{\text{ch}}) - \psi_{\text{in}}(\boldsymbol{x}^{\text{in}})\|^2] = 2m - 2\text{tr}(A_{\text{ch}}^\top \tilde{K}_{\text{ch}}\tilde{K}_{\text{in}} A_{\text{in}}).$$

(A.8)

The second moment matrices can be written as

$$\mathbb{E}[\psi_{\text{ch}}(\boldsymbol{x}^{\text{ch}})\psi_{\text{ch}}(\boldsymbol{x}^{\text{ch}})^\top] = A_{\text{ch}}^\top \tilde{K}_{\text{ch}}^2 A_{\text{ch}},$$
$$\mathbb{E}[\psi_{\text{in}}(\boldsymbol{x}^{\text{in}})\psi_{\text{in}}(\boldsymbol{x}^{\text{in}})^\top] = A_{\text{in}}^\top \tilde{K}_{\text{in}}^2 A_{\text{in}}.$$

Hence, the algorithm can be reduced to the following maximization problem:

$$\begin{aligned} \max \quad & \text{tr}(A_{\text{ch}}^\top \tilde{K}_{\text{ch}}\tilde{K}_{\text{in}} A_{\text{in}}) \\ \text{wrt} \quad & A_{\text{ch}}, A_{\text{in}} \in \mathbb{R}^{n \times m}, \\ \text{subj to} \quad & A_{\text{ch}}^\top \tilde{K}_{\text{ch}}^2 A_{\text{ch}} = A_{\text{in}}^\top \tilde{K}_{\text{in}}^2 A_{\text{in}} = I_m. \end{aligned}$$

Since the rank of $I_n - \boldsymbol{v}1_n^\top$ is $n-1$, the two shifted kernel matrices, $\tilde{K}_{\text{ch}}$ and $\tilde{K}_{\text{in}}$, are always singular. To avoid overfitting, we introduce two regularization terms $\gamma_{\text{ch}}A_{\text{ch}}^\top A_{\text{ch}}$ and $\gamma_{\text{in}}A_{\text{in}}^\top A_{\text{in}}$ to scale constraints as

$$A_{\text{ch}}^\top \tilde{K}_{\text{ch}}^2 A_{\text{ch}} + \gamma_{\text{ch}}A_{\text{ch}}^\top A_{\text{ch}} = A_{\text{in}}^\top \tilde{K}_{\text{in}}^2 A_{\text{in}} + \gamma_{\text{in}}A_{\text{in}}^\top A_{\text{in}} = I_m,$$

(A.9)

where $\gamma_{\text{ch}}$ and $\gamma_{\text{in}}$ are constants. The following algorithm finds $A_{\text{ch}}$ and $A_{\text{in}}$ minimizing the expected deviation subject to the regularized constraints.

**Algorithm 2** (Kernel Weighted CCA) Solve the following generalized eigen-decomposition:

$$\begin{bmatrix} 0_{n \times n} & \tilde{K}_{\text{ch}}\tilde{K}_{\text{in}} \\ \tilde{K}_{\text{in}}\tilde{K}_{\text{ch}} & 0_{n \times n} \end{bmatrix} \begin{bmatrix} \boldsymbol{\alpha}^{\text{ch}} \\ \boldsymbol{\alpha}^{\text{in}} \end{bmatrix}$$
$$= \begin{bmatrix} \tilde{K}_{\text{ch}}^2 + \gamma_{\text{ch}}I_n & 0_{n \times n} \\ 0_{n \times n} & \tilde{K}_{\text{in}}^2 + \gamma_{\text{in}}I_n \end{bmatrix} \begin{bmatrix} \boldsymbol{\alpha}^{\text{ch}} \\ \boldsymbol{\alpha}^{\text{in}} \end{bmatrix}$$

The matrix $\begin{bmatrix} A_{\text{ch}} \\ A_{\text{in}} \end{bmatrix}$ is set so that its $h$th column is set to the $h$th major generalized eigenvector.

**Theorem 4.** Algorithm 2 yields the parameters of the mapping functions $(\psi_{\text{ch}}, \psi_{\text{in}})$ which minimize the expected squared deviation (A.8) subject to the scaling constraints given in Eq. (A.9).

*Proof.* The problem of minimizing Eq. (A.8) subject to Eq. (A.9) is rewritten as

$$\max \quad \mathrm{tr}(A_{\mathrm{ch}}^{\top} \tilde{K}_{\mathrm{ch}} \tilde{K}_{\mathrm{in}} A_{\mathrm{in}})$$

$$\mathrm{wrt} \quad A_{\mathrm{ch}} \in \mathbb{R}^{n \times m}, A_{\mathrm{in}} \in \mathbb{R}^{n \times m},$$

$$\mathrm{subj\ to} \quad A_{\mathrm{ch}}^{\top} \tilde{K}_{\mathrm{ch}}^{2} A_{\mathrm{ch}} + \gamma_{\mathrm{ch}} A_{\mathrm{ch}}^{\top} A_{\mathrm{ch}}$$

$$= A_{\mathrm{in}}^{\top} \tilde{K}_{\mathrm{in}}^{2} A_{\mathrm{in}} + \gamma_{\mathrm{in}} A_{\mathrm{in}}^{\top} A_{\mathrm{in}} = I_{m}.$$

From Theorem 2, the optimization problem is solved by generalized eigendecomposition in Algorithm 2.    □

## A.4    Weighted SVM

Let us denote the protein-ligand interaction table by $Y \in \{\pm 1, 0\}^{n \times p_{\mathrm{in}}}$, where each row represents a ligand, and each column represents a target protein. Each cell in the table $Y$ takes one of three values, $\pm 1$ and 0: $+1$ and $-1$ indicate the existence and the absence of the interaction, respectively, and unknowns are 0.

SVM learning algorithm finds a classification boundary minimizing the violations for the constraints that training points are kept out of the margin. The weighted SVM employed in this study counts the violations with weights $v_j$ (See Sections 2.5 and 3.6) given to each training data.

When the interaction between a ligand $i$ and a target $t$ is predicted, ligands whose interactions with the target $t$ are known are selected as training data. If the index set of the ligands for training is denoted by $\mathcal{I}_i$, then the weighted SVM minimizes

$$\frac{1}{2}\|w\|^2 + C \sum_{j \in \mathcal{I}_i} v_j \max\left(0, 1 - Y_{j,t} f(x_j^{\mathrm{ch}}; w, b)\right),$$

where $C$ is constant. This algorithm is reduced to the classical SVM if all $v_j$'s are set to be equal in value. The dual problem of the weighted SVM learning algorithm can be derived as a quadratic program with box constraints and a single equality linear constraint, enabling fast learning with kernels.

**Tsuyoshi Kato** was born in 1975. He received his B.E., M.E., and Ph.D. degrees in engineering from Tohoku University, Sendai, Japan, in 1998, 2000, and 2003, respectively. From 2003 to 2005, he was with the National Institute of Advanced Industrial Science Technology (AIST) as a Postdoctoral Fellow in the Computational Biology Research Center (CBRC), Tokyo. From 2005 to 2008, he was an Assistant Professor at the Graduate School of Frontier Sciences, the University of Tokyo, Japan. He became an Associate Professor at the Center for Informational Biology, Ochanomizu University, Tokyo, Japan, in 2008. Then, he stayed in the University of Tokyo for seven months, and now he is an associate professor at the Graduate School of Engineering, Gunma University. His current scientific interests include bioinformatics and statistical pattern recognition. Dr. Kato is a member of The Institute of Electronics, Information and Communication Engineers (IEICE) and Japanese Society for Bioinformatics (JSBi).

**Richard Lemence** was born in 1977. He received his Ph.D. degree in 2005 from Niigata University, Niigata City, Japan under his adviser Dr. Kouie Sekigawa. In 2009–2012, he was a JST post-doctoral fellow at Ochanomizu University under the supervision of Dr. Jesper Jansson. Currently, he is an Associate Professor at the Institute of Mathematics, College of Science, University of the Philippines-Diliman. He is a member of the Mathematical Society of the Philippines and the American Mathematical Society.

(Communicated by  *Kengo Kinoshita*)

**Raissa Relator** was born in 1984. She received her B.S. and M.S. degrees in Mathematics from the University of the Philippines-Diliman in 2005 and 2008, respectively. She is currently pursuing a Ph.D. in Computer Science at the Graduate School of Engineering, Gunma University.