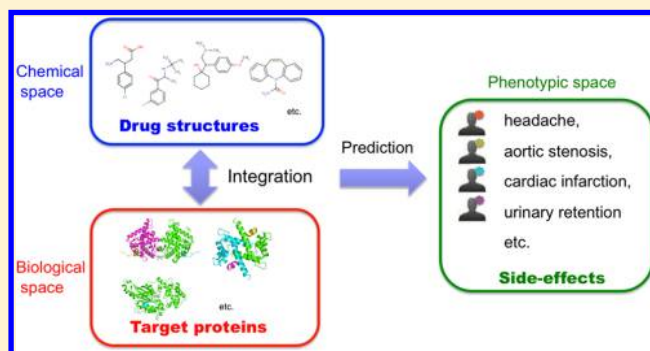


Drug Side-Effect Prediction Based on the Integration of Chemical and Biological Spaces

Yoshihiro Yamanishi,^{*,†} Edouard Pauwels,^{§,‡,⊥} and Masaaki Kotera^{||}[†]Division of System Cohort, Multi-scale Research Center for Medical Science, Medical Institute of Bioregulation, Kyushu University, 3-1-1 Maidashi, Higashi-ku, Fukuoka, Fukuoka 812-8582, Japan[‡]Mines ParisTech, Centre for Computational Biology, 35 rue Saint-Honore, F-77305 Fontainebleau Cedex, France[§]Institut Curie, F-75248, Paris, France[⊥]INSERM U900, F-75248, Paris, France^{||}Bioinformatics Center, Institute for Chemical Research, Kyoto University, Gokasho, Uji, Kyoto 611-0011, Japan

Supporting Information

ABSTRACT: Drug side-effects, or adverse drug reactions, have become a major public health concern and remain one of the main causes of drug failure and of drug withdrawal once they have reached the market. Therefore, the identification of potential severe side-effects is a challenging issue. In this paper, we develop a new method to predict potential side-effect profiles of drug candidate molecules based on their chemical structures and target protein information on a large scale. We propose several extensions of kernel regression model for multiple responses to deal with heterogeneous data sources. The originality lies in the integration of the chemical space of drug chemical structures and the biological space of drug target proteins in a unified framework. As a result, we demonstrate the usefulness of the proposed method on the simultaneous prediction of 969 side-effects for approved drugs from their chemical substructure and target protein profiles and show that the prediction accuracy consistently improves owing to the proposed regression model and integration of chemical and biological information. We also conduct a comprehensive side-effect prediction for uncharacterized drug molecules stored in DrugBank and confirm interesting predictions using independent information sources. The proposed method is expected to be useful at many stages of the drug development process.



■ INTRODUCTION

Drug side-effects, or adverse drug reactions, have become a major public health concern. It is one of the main causes of failure in the process of drug development and of drug withdrawal once a drug has reached the market. As an illustration to the extent of this problem, serious drug side-effects are estimated to be the fourth leading cause of death in the United States, resulting in 100 000 deaths per year.¹ The identification of potential severe adverse side-effects is a challenging issue at many stages of the drug development process.

From the viewpoint of systems biology, drugs can be regarded as molecules that induce perturbations to biological systems (consisting of various molecular interactions such as protein–protein interactions, chemical reactions, and signal transductions), leading to observed side-effects.² Actually, the body's response to a drug not only reflects the expected favorable effects due to the interaction with its target but also integrates the overall impact of off-target interactions. Indeed, even if a drug has a strong affinity for its target, it also often binds to other protein pockets with varying affinities, leading to

potential side-effects. It has been suggested that drugs with similar side-effects tend to share similar protein targets, and the use of side-effect similarity was proposed to identify unknown drug targets for marketed drugs.³

A useful experimental approach for predicting side-effects is preclinical in vitro safety profiling which tests compounds with biochemical and cellular assays, but experimental detection of drug side-effects remains very challenging in terms of cost and efficiency.⁴ Therefore, in silico prediction of potential side-effects early in the drug discovery process, before reaching the clinical stages, is of great interest to improve this long and expensive process and to provide new, efficient, and safe therapies for patients.

Recently, several computational methods for analyzing or predicting drug side-effects have been proposed, and the previous methods can be categorized into pathway-based approaches and chemical structure-based approaches. The principle of pathway-based approaches is to relate drug side-

Received: November 21, 2011

effects to perturbed biological pathways or subpathways because these pathways involve proteins targeted by the drug. A link between drug side-effects and biological pathways has been suggested by comparing biological pathways affected by toxic compounds and those affected by nontoxic compounds.⁵ A docking-based method has been proposed to identify off-targets of a drug by docking the drug into a protein's binding pocket similar to that of its primary target, and to map them onto known biological pathways, which allows to suggest potential side-effects.⁶ However, the method depends heavily on the availability of protein 3D structures, which limits its large-scale applicability. This limitation is critical for membrane proteins such as GPCRs or ion channels, which are major therapeutic targets, but whose 3D structure determination is difficult.

The principle of chemical structure-based approaches is to relate drug chemical structures with drug side-effects, following the spirit of QSAR (quantitative structure–activity relationship) and QSPR (quantitative structure–property relationship).^{7–9} A graph-based method for identifying chemical substructures associated with side-effects was proposed, but this method does not provide any framework for predicting side-effects for new drug candidate molecules.¹⁰ A method for predicting pharmacological information using chemical structures was proposed, but the method cannot be applied to predict high-dimensional side-effect profiles directly.¹¹ An algorithmic framework for predicting side-effect profiles from chemical structure data was established using canonical correlation analysis (CCA), which is a pioneering work in terms of simultaneous prediction of many side-effects.¹² A chemical fragment-based approach was proposed to relate drug chemical substructures with side-effects by sparse canonical correlation analysis (SCCA), and the extracted chemical substructures were used for predicting side-effect profiles.¹³

All previous works on the problem of drug side-effect prediction have focused on the use of information about either drug targeted proteins only⁶ or drug chemical structures only.^{12,13} However, it is more logical to consider both target protein and chemical structure information simultaneously when predicting drug side-effects, rather than using each of them independently.

In this paper, we develop a new method to predict potential side-effect profiles of drug candidate molecules based on their chemical structures and target protein information on a large scale. We propose several extensions of kernel regression model for multiple responses to deal with heterogeneous data sources. The originality lies in the integration of the chemical space of drug chemical structures and the biological space of drug target proteins in a unified framework. Figure 1 shows an illustration of the proposed method. We demonstrate the usefulness of the proposed method on the simultaneous prediction of 969 side-effects for 658 approved drugs from their chemical substructure and target protein profiles and show that the prediction accuracy consistently improves owing to the proposed regression model and integration of chemical and biological information. We also conduct a comprehensive side-effect prediction for many uncharacterized drug molecules stored in DrugBank and confirm interesting predictions using independent information sources.

MATERIALS

Drug Side-Effect Profiles. Side-effect terms were obtained from the SIDER database which contains information about

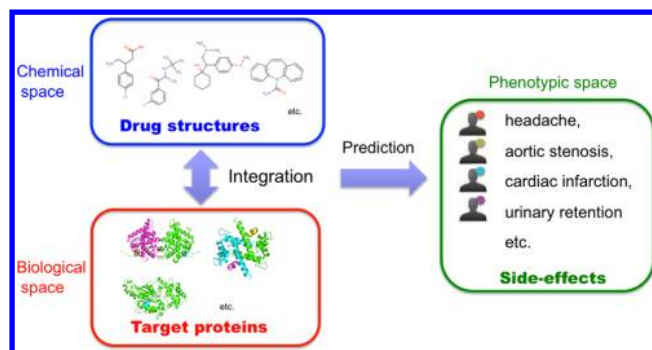


Figure 1. Illustration of the proposed method.

marketed medicines and their recorded side-effects.¹⁴ In this study, we focused on side-effects associated with drugs which are annotated as “small molecules” in the DrugBank database.¹⁵ There are some side-effects which are associated with almost all drugs (e.g., nausea, dizziness, vomiting, and rash), while there are some side-effects associated with very few drugs (e.g., agnosia, variant angina, aspergillosis, and gouty arthritis). Therefore, we removed side-effects which lie at the top 10% in terms of frequency (which are associated with more than 131 drugs), and we also removed side-effects which are associated with only one drug (singletons). This produced a data set consisting of 658 drugs, 969 side-effect terms, and 23 061 associations between drugs and side-effects, and each drug had 35.1 side-effects on average. The left panel in Figure 2 shows

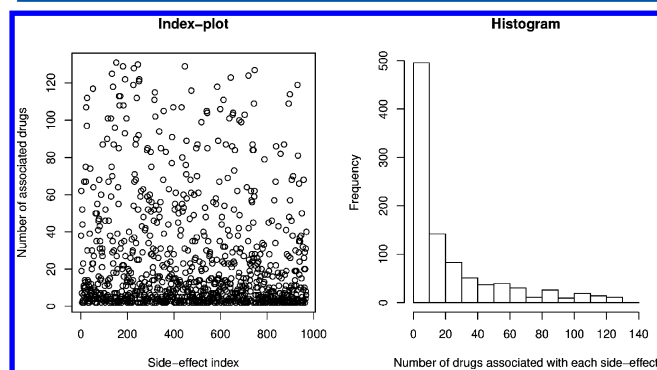


Figure 2. Characteristics of side-effect data. The left panel shows the index-plot of the number of associated drugs for each side-effect, and the right panel shows the histogram of the number of associated drugs for each side-effect.

the index-plot of the number of associated drugs for each side-effect, and the right panel in Figure 2 shows the histogram of the number of associated drugs for each side-effect.

This data set is used as gold standard data, and drugs in the set are referred to as reference drugs throughout this study. Each drug was represented by a 969 dimensional feature vector $y = (y_1, \dots, y_q)^T$, where each element encodes for the presence or absence of each of the side-effect terms by 1 or 0, respectively, and q is the number of side-effects. The feature vector is referred to as a side-effect profile.

Chemical Profiles. To encode the chemical structure of drugs, we used a fingerprint corresponding to the 881 chemical substructures defined in the PubChem database.¹⁶ Each drug was represented by an 881 dimensional feature vector $x^{(\text{chem})} = (x_1, \dots, x_{p_1})^T$, where each element encodes for the presence or absence of each PubChem substructure by 1 or 0, respectively,

and p_1 is the number of chemical substructures. The feature vector is referred to as chemical profile.

Biological Profiles. Drug–protein interactions were extracted from the DrugBank database¹⁵ and the Matador database.¹⁷ We used proteins that are indicated as direct interactions and removed proteins whose annotations were “Multidrug”, “Cytochrome”, “ATP-binding cassette”, “Glutathione S-transferases”, or “Flavin containing monooxygenase” in the UniProt database,¹⁸ because these proteins are involved in drug metabolism not drug mechanism. As a result, we extracted 5074 drug–protein interactions involving 1368 unique target proteins. Each drug is represented by a 1368 dimensional feature vector $\mathbf{x}^{(\text{bio})} = (x_1, \dots, x_{p_2})^T$, where each element encodes for the presence or absence of each target protein by 1 or 0, respectively, and p_2 is the number of target proteins. The feature vector is referred to as biological profile.

METHODS

Suppose that we have a set of n drugs with chemical profiles of p_1 chemical substructures, with biological profiles of p_2 target proteins, and with side-effect profiles of q side-effects. Each drug is represented by a chemical feature vector $\mathbf{x}^{(\text{chem})} = (x_1, \dots, x_{p_1})^T$, by a biological feature vector $\mathbf{x}^{(\text{bio})} = (x_1, \dots, x_{p_2})^T$, and by a side-effect feature vector $\mathbf{y} = (y_1, \dots, y_q)^T$.

We consider the situation where we are given a new drug candidate molecule with chemical profile $\mathbf{x}_{\text{new}}^{(\text{chem})}$ and biological profile $\mathbf{x}_{\text{new}}^{(\text{bio})}$, and we want to predict its potential side-effect profile \mathbf{y}_{new} based on the chemical and biological profiles.

Kernel Regression (KR). We consider a regression model to predict q -dimensional feature vector \mathbf{y} (response variables) from p -dimensional feature vector \mathbf{x} (explanatory variables). Here we propose to apply a kernel regression (KR) model for multiple responses, formulated as follows:

$$\mathbf{y}_{\text{new}} = f(\mathbf{x}_{\text{new}}) = \sum_{i=1}^n k(\mathbf{x}_{\text{new}}, \mathbf{x}_i) \mathbf{w}_i + \varepsilon \quad (1)$$

where \mathbf{x}_{new} is the feature vector of a new object, \mathbf{y}_{new} is the response vector of the new object, f is the projection $f: \mathcal{X} \rightarrow \mathbf{R}^q$, \mathcal{X} is a set of objects, n is the number of objects in a training set, \mathbf{x}_i is a feature vector of the i th object in the training set, $\mathbf{w}_i \in \mathbf{R}^q$ is a weight vector, $k(\cdot, \cdot)$ is a positive definite kernel, that is, a symmetric function $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbf{R}$ satisfying $\sum_{i,j=1}^n a_i a_j k(\mathbf{x}_i, \mathbf{x}_j) \geq 0$ for any $a_i, a_j \in \mathbf{R}$, and ε is a noise vector. In this study we use Gaussian RBF (Radial Basis Function) kernel defined as $k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / 2\sigma^2)$, where σ is the width parameter.

The fitting of the model can be done by finding \mathbf{w}_i which minimizes the following penalized loss function:

$$\begin{aligned} L &= \|\mathbf{Y} - \mathbf{KW}\|_{\text{F}}^2 + \lambda \|\mathbf{W}\|_{\text{F}}^2 \\ &= \text{trace}\{(\mathbf{Y} - \mathbf{KW})(\mathbf{Y} - \mathbf{KW})^T\} + \lambda \text{trace}\{\mathbf{W}\mathbf{W}^T\} \end{aligned} \quad (2)$$

where \mathbf{K} is an $n \times n$ kernel similarity matrix $(K)_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$, $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_n)^T$, $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)^T$, λ is a regularization parameter, and $\|\cdot\|_{\text{F}}$ is the Frobenius norm. Note that we put the above penalty term to avoid overfitting.

Taking the differential of L with respect to \mathbf{W} and setting to zero, the weight matrix \mathbf{W} can be analytically obtained by solving the following equation:

$$(\mathbf{K}^2 + \lambda \mathbf{I})\mathbf{W} = \mathbf{KY} \quad (3)$$

Once we have trained our model, which is computing \mathbf{W} , we can apply the model to unseen objects (drugs in this case).

For the side-effect prediction problem in this study, we consider four predictive models. The side-effect prediction for a new drug candidate molecule with chemical profile $\mathbf{x}_{\text{new}}^{(\text{chem})}$ and biological profile $\mathbf{x}_{\text{new}}^{(\text{bio})}$ is performed as follows:

- **KR with chemical information (KR chem).** The prediction is performed based only on $\mathbf{x}^{(\text{chem})}$ as follows:

$$\mathbf{y}_{\text{new}} = \sum_{i=1}^n k(\mathbf{x}_{\text{new}}^{(\text{chem})}, \mathbf{x}_i^{(\text{chem})}) \mathbf{w}_i \quad (4)$$

- **KR with biological information (KR bio).** The prediction is performed based only on $\mathbf{x}^{(\text{bio})}$ as follows:

$$\mathbf{y}_{\text{new}} = \sum_{i=1}^n k(\mathbf{x}_{\text{new}}^{(\text{bio})}, \mathbf{x}_i^{(\text{bio})}) \mathbf{w}_i \quad (5)$$

- **KR with the kernel integration of chemical and biological information (KR Kchem + Kbio).** The prediction is performed based on the sum of the two kernel functions for chemical and biological profiles as follows.

$$\begin{aligned} \mathbf{y}_{\text{new}} &= \sum_{i=1}^n \{c_1 k_1(\mathbf{x}_{\text{new}}^{(\text{chem})}, \mathbf{x}_i^{(\text{chem})}) + c_2 k_2(\mathbf{x}_{\text{new}}^{(\text{bio})}, \mathbf{x}_i^{(\text{bio})})\} \mathbf{w}_i \end{aligned} \quad (6)$$

where k_1 and k_2 are kernel functions for chemical profiles and biological profiles, respectively, and c_1 and c_2 are the weights for k_1 and k_2 in the data integration. A possible useful solution for the weights in practice would be to assign weight values according to the prediction accuracies of individual data sources. For example, we could use weights that are proportional to AUC – 0.5 or AUPR – a random AUPR, where AUC means the area under the rate of change (ROC) curve and AUPR means the area under the precision–recall curve. In fact, the usefulness of this procedure was already shown in the context of protein network inference.¹⁹

Note that \mathbf{W} is computed based on the training set for each model. In any predictive models, if the j th element in \mathbf{y}_{new} has a high score, the new molecule is predicted to have the j th side-effect ($j = 1, 2, \dots, q$).

Multiple Kernel Regression (MKR). In practice the same objects are represented by multiple different data sources. Suppose that each object is represented by m heterogeneous feature vectors as $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$ ($m = 2$ in our case, since drugs are represented by chemical and biological profiles).

To deal with such a situation, we propose an extension of the kernel regression model, which we call multiple kernel regression (MKR), formulated as follows:

$$\mathbf{y}_{\text{new}} = f(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}) = \sum_{l=1}^m \sum_{i=1}^n k_l(\mathbf{x}_{\text{new}}^{(l)}, \mathbf{x}_i^{(l)}) \mathbf{w}_i^{(l)} + \varepsilon \quad (7)$$

where $\mathbf{x}_{\text{new}}^{(l)}$ is the feature vector of a new object for the l th data source, n is the number of objects in a training set, $\mathbf{x}_i^{(l)}$ is a feature vector of the i th element for the l th data source in the training set, $\mathbf{w}_i^{(l)} \in \mathbf{R}^q$ is a weight vector, $k_l(\cdot, \cdot)$ is a kernel

similarity function for $\mathbf{x}^{(l)}$, and ε is a noise vector. In this study we use Gaussian RBF kernel for each k_i .

The fitting of the model can be done by finding $\mathbf{w}_i^{(1)}, \dots, \mathbf{w}_i^{(m)}$ which minimizes the following penalized loss function:

$$\begin{aligned} L &= \|Y - \sum_{l=1}^m K_l W_l\|_F^2 + \sum_{l=1}^m \lambda_l \|W_l\|_F^2 \\ &= \text{trace}\{(Y - \sum_{l=1}^m K_l W_l)(Y - \sum_{l=1}^m K_l W_l)^T\} \\ &\quad + \sum_{l=1}^m \lambda_l \text{trace}\{W_l W_l^T\} \end{aligned} \quad (8)$$

where K_l is an $n \times n$ kernel similarity matrix $(K_l)_{ij} = k_i(\mathbf{x}_i, \mathbf{x}_j)$, $W_l = (\mathbf{w}_1^{(l)}, \dots, \mathbf{w}_n^{(l)})^T$, $Y = (y_1, \dots, y_n)^T$, λ_l is a regularization parameter, and $\|\cdot\|_F$ is the Frobenius norm.

Taking the differential of L with respect to W_1, \dots, W_m and setting to zero, the weight matrices W_1, \dots, W_m can be analytically obtained by solving the following equation:

$$\begin{pmatrix} K_1 K_1 + \lambda_1 I & K_1 K_2 & \dots & K_1 K_m \\ K_2 K_1 & K_2 K_2 + \lambda_2 I & \dots & K_2 K_m \\ \dots & \dots & \ddots & \dots \\ K_m K_1 & K_m K_2 & \dots & K_m K_m + \lambda_m I \end{pmatrix} \begin{pmatrix} W_1 \\ W_2 \\ \dots \\ W_m \end{pmatrix} = \begin{pmatrix} K_1 Y \\ K_2 Y \\ \dots \\ K_m Y \end{pmatrix} \quad (9)$$

Once we have trained our model, which is computing W_1, \dots, W_m , we can apply the model to unseen objects (drugs in this case).

For the side-effect prediction problem in this study, we consider a predictive model for a new drug candidate molecule with chemical profile $\mathbf{x}_{\text{new}}^{(\text{chem})}$ and biological profile $\mathbf{x}_{\text{new}}^{(\text{bio})}$, as follows:

- **MKR based on chemical and biological information (MKR chem + bio).**

$$\begin{aligned} y_{\text{new}} &= \sum_{i=1}^n k_1(\mathbf{x}_{\text{new}}^{(\text{chem})}, \mathbf{x}_i^{(\text{chem})}) \mathbf{w}_i^{(1)} \\ &\quad + \sum_{i=1}^n k_2(\mathbf{x}_{\text{new}}^{(\text{bio})}, \mathbf{x}_i^{(\text{bio})}) \mathbf{w}_i^{(2)} \end{aligned} \quad (10)$$

where k_1 and k_2 are kernel functions for chemical profiles and biological profiles, respectively.

Note that W_1, \dots, W_m are computed based on the training set. If the j th element in y_{new} has a high score, the new molecule is predicted to have the j th side-effect ($j = 1, 2, \dots, q$).

Other Possible Prediction Methods. Canonical Correlation Analysis (CCA). The use of canonical correlation analysis (CCA)²⁰ was proposed to predict drug side-effect profiles from chemical substructure profile.¹² We make a brief review of the CCA-based method.

Suppose that we have a set of n drugs $\{\mathbf{x}_i\}_{i=1}^n$ represented by chemical profiles and $\{y_i\}_{i=1}^n$ represented by side-effect profiles, where each drug is represented by a chemical substructure feature vector $\mathbf{x} = (x_1, \dots, x_p)^T$, and by a side-effect feature vector $\mathbf{y} = (y_1, \dots, y_q)^T$. Consider two linear combinations for chemical substructures and side-effects as $u_i = \alpha^T \mathbf{x}_i$ and $v_i = \beta^T \mathbf{y}_i$ ($i = 1, 2, \dots, n$), where $\alpha = (\alpha_1, \dots, \alpha_p)^T$ and $\beta = (\beta_1, \dots, \beta_q)^T$ are weight vectors. The goal of ordinary CCA is to find weight vectors α and β which maximize the following canonical correlation coefficient:

$$\rho = \text{corr}(u, v) = \frac{\sum_{i=1}^n \alpha^T \mathbf{x}_i \beta^T \mathbf{y}_i}{\sqrt{\sum_{i=1}^n (\alpha^T \mathbf{x}_i)^2} \sqrt{\sum_{i=1}^n (\beta^T \mathbf{y}_i)^2}} \quad (11)$$

where u and v are assumed to be centered. Multiple canonical components (m canonical components, $k = 1, 2, \dots, m$) can be obtained by the following manipulation: $(\alpha_k, \beta_k) = \arg \max \rho = \text{corr}(u, v)$ under the following orthogonality constraints: $\alpha \perp \alpha_1, \dots, \alpha_{k-1}$ and $\beta \perp \beta_1, \dots, \beta_{k-1}$.

For the prediction of the side-effect profile for a new drug \mathbf{x}_{new} , we compute the following prediction score:

$$y_{\text{new}} = B^{-T} A^T \mathbf{x}_{\text{new}} \quad (12)$$

where $A = [\alpha_1, \dots, \alpha_m]$, $B = [\beta_1, \dots, \beta_m]$, m is the number of components, and B^{-T} is the pseudo-inverse matrix of B^T .

For the side-effect prediction problem in this study, we consider three predictive models for a new drug candidate molecule with chemical profile $\mathbf{x}_{\text{new}}^{(\text{chem})}$ and biological profile $\mathbf{x}_{\text{new}}^{(\text{bio})}$, as follows:

- **CCA with chemical information (CCA chem).** The prediction is performed based only on $\mathbf{x}_{\text{new}}^{(\text{chem})}$.
- **CCA with biological information (CCA bio).** The prediction is performed based only on $\mathbf{x}_{\text{new}}^{(\text{bio})}$.
- **CCA with chemical and biological information (CCA chembio).** The prediction is performed based on the concatenated profile $\mathbf{x}^{(\text{chembio})} = (\mathbf{x}^{(\text{chem})T}, \mathbf{x}^{(\text{bio})T})^T$.

RESULTS

Performance Evaluation. We applied the proposed methods to predict drug side-effect profiles from chemical profiles (chemical substructures) and biological profiles (target protein patterns). We tested seven approaches: (1) “CCA chem”, (2) “CCA bio”, (3) “CCA chembio”, (4) “KR chem”, (5) “KR bio”, (6) “KR Kchem + Kbio”, and (7) “MKR chem + bio” on their abilities to predict known side-effect profiles using 658 reference drugs (see the Methods section for a description of each approach). Note that the CCA chem approach corresponds to a previous method.¹² For simplicity, the same weight is used in the kernel integration of the KR Kchem + Kbio method ($c_1 = c_2 = 0.5$) in this study. We performed the following fivefold cross-validation: Drugs in the gold standard set were split into five subsets of roughly equal size, each subset was then taken in turn as a test set, and we performed the training on the remaining four sets. For accurate comparison, we kept the same experimental conditions, where the same training drugs and test drugs are used across the different methods in each cross-validation fold.

We evaluated the performance of each method by the precision–recall curve (PR curve), because the PR curve is known to be a practical performance measure²¹ and the PR curve was also featured in a previous study on drug side-effect prediction.¹² The performance for predicting all side-effects can be summarized by the area under the PR curve (AUPR). All parameters in each method (e.g., regularization parameters, number of components, kernel parameters) were optimized using grid search with the AUPR score as an objective function. To obtain a robust result, we repeated the overall cross-validation procedure three times and computed the average and standard deviation (SD) of the AUPR score.

We also performed the same analysis with the receiver operating curve (ROC curve), where all parameters in each method were optimized using grid search with the AUC score

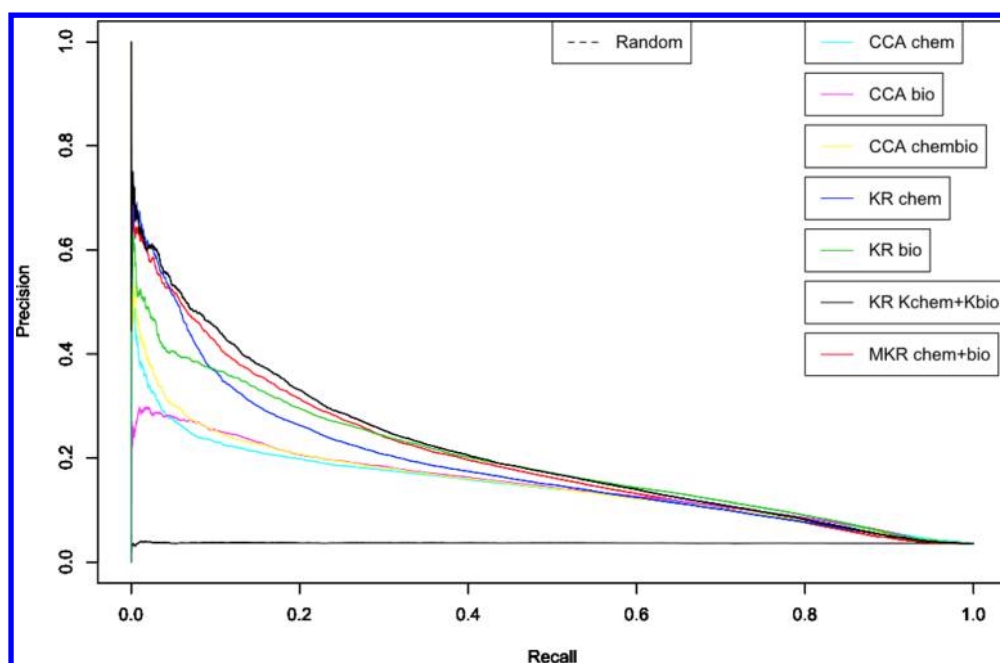


Figure 3. Precision-recall curves based on the fivefold cross-validation experiment. The PR curve is the plot of precision (positive predictive value) as a function of recall (sensitivity) based on various thresholds, where the recall is defined as $TP/(TP + FN)$ and the precision is defined as $TP/(TP + FP)$ where TP, FP, TN, and FN are the number of true positives, false positives, true negatives, and false negatives, respectively.

as an objective function. The ROC curve and the AUC score are shown in Supplemental Figures 1 and 2 in the Supporting Information. The result shows similar tendencies as those observed in the precision–recall curves and AUPR scores. The PR curve is more informative than the receiver operating curve (ROC curve) when the number of positive examples is much lower than that of negative examples.²¹ Therefore, we focused on the PR curve and the AUPR score below.

Figure 3 shows the PR curves for the seven different approaches based on the cross-validation experiment, where the prediction scores for all side-effects were merged and a global PR curve was drawn for each method. It seems that KR-based methods outperform CCA-based methods using any data sources, which suggests that the proposed method works better than the previous method.¹² In addition, the KR-based methods are computationally efficient compared with the CCA-based method, which is shown in Supplemental Figure 6 in the Supporting Information. The KR Kchem + Kbio and MKR chem + bio methods seem to work better than the KR chem and KR bio methods, which suggests that the integration of chemical and biological information is meaningful. We examined the effect of weighted data integration in the KR Kchem + Kbio method, and the result is shown in Supplemental Figure 3 in the Supporting Information. It is observed that the weighted method and the unweighted method ($c_1 = c_2 = 0.5$) are comparable when using our data.

We investigated the prediction accuracy of predicted side-effects for each drug with a high level of confidence. We computed the number of drugs having a known side-effect which is ranked highest in the prediction score, and the number of drugs having at least one known side-effect which is ranked among the five highest ranking side-effects. In fact these statistics were used in the previous study.¹² The resulting numbers are shown in Table 1. The proposed KR-based or MKR-based methods seem to outperform the previous CCA-based method. For example, the proposed KR Kchem + Kbio

Table 1. Performance Statistics Based on Fivefold Cross-Validation^a

method	AUPR \pm SD	Top1 \pm SD	Top5 \pm SD
random	0.0364 \pm 0.0000	27 \pm 2.64	99 \pm 5.56
CCA chem	0.1493 \pm 0.0015	203 \pm 6.00	398 \pm 8.50
CCA bio	0.1479 \pm 0.0010	221 \pm 9.84	418 \pm 4.50
CCA chembio	0.1527 \pm 0.0009	218 \pm 5.68	412 \pm 14.22
KR chem	0.1817 \pm 0.0025	250 \pm 3.78	431 \pm 1.52
KR bio	0.1933 \pm 0.0011	259 \pm 9.45	451 \pm 10.96
KR Kchem + Kbio	0.2089 \pm 0.0024	280 \pm 14.97	461 \pm 6.02
MKR chem + bio	0.2008 \pm 0.0020	272 \pm 5.29	461 \pm 8.18

^aAUPR is the total area under the precision–recall curve. Top1 indicates the number of drugs having a known side-effect which is ranked highest in the prediction score. Top5 indicates the number of drugs having at least one known side-effect which is ranked among the top five high scoring side-effects. SD indicates the standard deviation over the repetition of the cross-validation. The best result in each column is shown in bold.

method ranked first one of the known side-effects of 41.7% (275) of the 658 reference drugs and ranked a correct side-effect among the top five scoring side-effects for 70.0% (461) of the 658 reference drugs. This result suggests that predicted side-effects with high scores for drugs are more accurate in practice and heterogeneous data integration is useful in drug side-effect prediction. The performance of the KR Kchem + Kbio method is slightly better than that of MKR chem + bio method. One explanation for the better performance of KR Kchem + Kbio over MKR chem + bio is that MKR requires too many parameters to be estimated (MKR chem + bio, W1 and W2, compared with KR Kchem + Kbio, W) and the predictive model MKR chem + bio may be overfitting.

We investigated the overall accuracy of the predicted side-effect profile for each drug. The misclassification rate of each drug differs from side-effect to side-effect, so we evaluated the accuracy of each drug by computing the AUPR for each, based

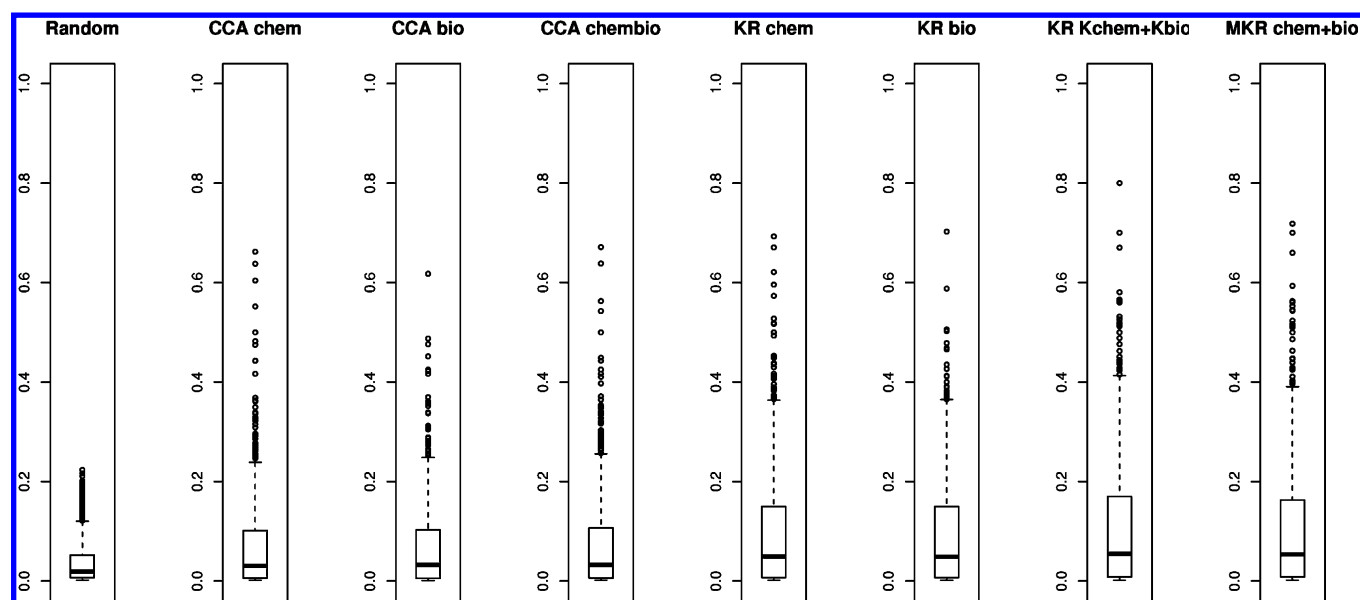


Figure 4. Boxplots of the AUPR (area under the precision–recall curve) scores for individual side-effects.

Table 2. Side-effect Prediction for Betamethasone Based on Chemical, Biological, and Both Information^a

keyword	freq	MKR C + B score	KR C + B score	KR bio score	KR chem score	MKR C + B rank	KR C + B rank	KR bio rank	KR chem rank	confirmation
cataract	103	0.651	0.959	0.886	0.853	3	3	7	48	yes
striae	24	0.655	0.882	0.785	0.818	2	8	66	77	yes
glaucoma	71	0.533	0.837	0.691	0.791	21	16	210	122	yes
hypopigmentation	31	0.573	0.83	0.868	0.699	10	18	12	392	yes
ecchymosis	108	0.383	0.701	0.518	0.646	360	97	1135	699	yes
hypertrichosis	31	0.516	0.7	0.688	0.612	33	101	216	985	yes
miliaria	20	0.521	0.695	0.7	0.61	30	109	186	1007	no
skin atrophy	23	0.546	0.695	0.701	0.623	16	110	180	877	yes
osteoporosis	47	0.414	0.69	0.544	0.635	200	118	895	778	yes
folliculitis	36	0.53	0.683	0.7	0.61	24	129	183	1003	no
pseudotumor cerebri	20	0.356	0.658	0.496	0.609	553	189	1332	1016	yes
exophthalmos	14	0.322	0.656	0.445	0.625	1041	194	1978	855	yes
aseptic necrosis	8	0.324	0.654	0.444	0.622	995	197	2017	882	yes
petechiae	70	0.316	0.653	0.441	0.623	1149	202	2095	874	yes
peptic ulcer	60	0.318	0.653	0.446	0.62	1123	203	1959	901	yes

^aThe frequency means the number of drugs with the predicted side-effects in the reference set. The score means the prediction score for each method (KR or MKR). “C + B” means the integration of chemical and biological data. The ranks are given according to the scores of all predicted keyword-molecule associations for each method. The confirmation column indicates whether the prediction validity were confirmed or not in independent sources.

on its predicted score vector and the true side-effect profile in each cross-validation fold. Examples of drug categories with bad AUPR scores are GABA agents, gastrointestinal agents, antiacne agents, antidyskinetics, analgesics, and adjuvants. The whole results are shown in Supplemental Table 1 in the Supporting Information.

We examined the prediction accuracy for individual side-effects. We calculated the PR curve for each side-effect and, then, computed the AUPR score (area under the PR curve) for each. Figure 4 shows the boxplot representing the distribution of the resulting AUPR scores for 969 side-effects using each method. It also seemed that the KR-based methods outperformed the CCA-based methods; the integration of chemical and biological information produced the best results. Examples of side-effects with high predictive accuracy are “chronic active hepatitis”, “pemphigus”, “pseudomembranous

colitis”, “Aseptic necrosis”, “interstitial nephritis”, and “gynecomastia”. Examples of side-effects with low predictive accuracy are “diabetic neuropathy”, “nephrogenic diabetes insipidus”, “bladder tumors”, “bone neoplasm”, “narcolepsy”, and “parotid gland enlargement”. The whole results are shown in Supplemental Table 2 in the Supporting Information. We also investigated the dependence of performance of the proposed method on the prevalence of the side-effect. The scatter-plot of the individual AUPR scores against the side-effect frequencies (the number of the associated drugs in the reference data) is shown in Supplemental Figure 4 in the Supporting Information. The AUC scores seemed to be positively correlated with side-effect frequencies to some extent, but there were many side-effects which had low frequency and high AUPR scores. These results demonstrate

Table 3. Side-Effect Prediction for a Different Drug Class^a

ATC	efficacy	drug name	side-effect	score	confirmation
M	muscle relaxants, centrally acting agents	baclofen	hypoventilation	0.218	yes
A	antiobesity	sibutramine	testicular swelling	0.154	no
A	antiobesity	sibutramine	abdominal cramps	0.149	no
M	muscle relaxants, centrally acting agents	baclofen	sneezing	0.144	yes
N	psychoanaleptics	atomoxetine	memory loss	0.139	no
N	anesthetics	desflurane	encephalopathy	0.137	yes
L	antineoplastic agents	aminolevulinic acid	hypocalcemia	0.133	no
L	anticancer	Ifosfamide	hypocalcemia	0.128	no
M	muscle relaxants, centrally acting agents	baclofen	respiratory arrest	0.127	yes
M	affecting bone structure and numeralization	pamidronate	hyperpigmentation	0.125	no
L	cytotoxic antibiotics	valrubicin	pseudomembranous colitis	0.124	no
L	antineoplastic agents	hydroxyurea	hypocalcemia	0.122	no
M	affecting bone structure and numeralization	alendronate	hyperpigmentation	0.119	yes

^aThe ATC column indicates 1 of the 14 ATC classes. The score column indicates the prediction score by the proposed method. The confirmation column indicates whether the prediction validities were confirmed or not in independent sources.

the high-performance predictive power of the proposed KR method for side-effect prediction in practical applications.

Comprehensive Side-Effect Prediction for Uncharacterized Drugs. We then made a comprehensive prediction of side-effects for uncharacterized drugs. We considered molecules in DrugBank for which side-effect information was not available in the SIDER database. We focused on 730 drugs labeled as small molecules in DrugBank for which chemical structure and target protein information is available. We predicted the side-effect profiles of the 730 drugs using all 658 reference drugs as a training set, based on the chemical profiles, the biological profiles, and the integration of both profiles. We used the proposed “KR Kchem + Kbio” method. All the prediction results can be viewed in the Supporting Information.

The adverse drug reactions of unmarketed drugs are usually not available. Additionally, publications regarding adverse drug effects deal with reported effects; there is no information about negative results where a given drug does not induce a certain effect. Therefore, we tried to confirm some true positive side-effect predictions using independent sources. We provide two examples of side-effect prediction confirmed by independent sources.

Lvonorgestrel (DB00367) is a synthetic progestational hormone used for contraception, control of menstrual disorders, and treatment of endometriosis. The top scoring keywords predicted for this drug are relatively frequent side-effects. It is more challenging to predict rare side-effects rather than frequent side-effects. Thus, we examined infrequent side-effects (e.g., appearing less than 20 times in our reference data). Among them the three first predicted side-effects for Lvonorgestrel are “ovarian cyst”, “breast tenderness”, and “melasma” which were confirmed in the literature.^{22,23}

To gain insights into the differences between the KR and MKR methods and between the data sources (chemical profiles, biological profile, and the integration of both), we took the example of betamethasone (DB00443), a glucocorticoid. Table 2 shows examples of predicted side-effects with high scores for betamethasone, where only top 15 side-effects predicted by the proposed method are shown due to space limitation. We were able to validate the prediction for the side-effect keywords based on independent sources.²⁴ The validated side-effects are marked as “confirmed” in the last column in Table 2. Thirteen out of the 15 top predicted side-effects for this drug were confirmed based on the literature. This example illustrates that

predicting side-effects based on both chemical and biological information gives higher score and lower rank, compared to the prediction based using separate information. It is observed that the MKR method and the KR method have a similar tendency in high scoring predictions. The use of both chemical and biological information improved prediction accuracy, as suggested in the Performance Evaluation section.

Both drugs considered here are part of large classes (progesterones and glucocorticoids), within which the side-effects are mostly similar. Therefore, the high prediction accuracy for these drugs is not surprising. We further surveyed the result for different families of drugs for which side-effect information was not available from SIDER database and found some side-effects were written in other data sources.

Sparfloxacin (DB01208) is a fluoroquinolone antibiotic used in the treatment of bacterial infections. Among the 12 first predicted side-effects, 9 could be confirmed (pseudomembranous colitis, myasthenia gravis, dysphasia, intestinal perforation, cerebral thrombosis, phobia, anosmia, and increased amylase keratoconjunctivitis) using the package insert from the FDA Web site (US Food and Drug Administration Home Page).²⁵

Escitalopram (DB01175) belongs to a class of antidepressant agents known as selective serotonin-reuptake inhibitors. Eight of the 20 first predicted side-effects for this drug could be confirmed (serotonin syndrome, drug dependence, hyperprolactinemia, choreoathetosis, sensory disturbance, suicidal tendency, and eye abnormality).²⁶

Cytarabine (DB00987) is an antimetabolite antineoplastic agent that inhibits the synthesis of DNA. It is used mainly in the treatment of leukemia. Five of the 12 first predicted side-effects for this drug could be confirmed (acute respiratory distress syndrome, megaloblastic anemia, polyps, peritonitis, and paraplegia).²⁷

Codeine is an opioid analgesic known to induce relatively strong physical dependence. The first three predicted side-effects for this drug contain drug dependence and drug withdrawal which could be confirmed.

These confirmations suggest that the method could be used to provide insights regarding the side-effect profiles of uncharacterized drugs.

Side-Effect Prediction for Drugs Belonging to a Different Class. Finally, we tested the potential of our method to predict side-effects for drugs belonging to a different

class. We assume the situation where the class of a given drug is different from those of drugs in a training set. In this study drugs were grouped into different classes based on the anatomical therapeutic chemical (ATC) classification of the World Health Organization (WHO). The ATC classification has the following 14 classes: (A) alimentary tract and metabolism, (B) blood and blood forming organs, (C) cardiovascular system, (D) dermatologicals, (G) genito urinary system and sex hormones, (H) systemic hormonal preparations, excluding sex hormones and insulins, (J) antiinfectives for systemic use, (L) antineoplastic and immunomodulating agents, (M) musculo-skeletal system, (N) nervous system, (P) antiparasitic products, insecticides, and repellents, (R) respiratory system, (S) sensory organs, and (V) various.

We performed the following procedure: (1) A class of drugs in our data was taken as a test set, and the other classes of drugs (belonging to the remaining 13 ATC classes) were used as a training set. (2) The predictive model of the KR method with chemical and biological data was learned using the training set only. (3) Then, the side-effects were predicted for drugs in the test set. (4) These operations were performed in turn for each of 14 ATC classes. Among the prediction results, we picked up some pairs of drugs and side-effects with high prediction scores, and we investigated whether the relations between drugs and side-effects can be found in an independent source that we did not use in the learning process. Note that we used drugs belonging to only one class, which removed the effect of drugs belonging to multiple classes.

Table 3 shows some examples of the predicted pairs of drugs and side-effect keywords with high scores, where we focus on relatively rare side-effect keywords, i.e., those that appear less than 20 times in our reference data set. It was predicted that baclofen, a muscle relaxant, could cause hypoventilation and sneezing. We found that these side-effects have been already reported in eHealthMe (<http://www.ehealthme.com/>), which collects 40 million latest outcomes of 45 000 drugs, vitamins, and supplements since 1977 from the FDA and community. The most structurally related drug that contributed to one of the predictions (i.e., baclofen could cause hypoventilation) was gabapentin, an antiepileptic. Similarly, lorazepam, an anxiolytic, has the most related target protein and contributed to the prediction. The both of drugs are known to cause hypoventilation, and both belong to a different class from that of baclofen.

As another example, it was predicted that desflurane (an anesthetic) could cause encephalopathy. This prediction was supported by some drugs including fluorouracil (an anti-metabolite) and acyclovir (an antiviral), which are structurally related and interact with similar proteins, respectively. We confirmed this prediction, i.e., the actual case where desflurane may have caused encephalopathy as described in LiverTox (<http://livertox.nih.gov/>), which collects clinical information on liver injury cases. Our prediction result contains many pairs of drugs and side-effects that had not been reported. For example, sibutramine, an antiobesity, was predicted to cause testicular swelling and abdominal cramps, but such side-effect reports could not be found in eHealthMe nor anywhere else. Nevertheless, we also confirmed other successful side-effect predictions for some drugs as listed in Table 3.

DISCUSSION AND CONCLUSION

In this paper we proposed a novel method to predict potential side-effect profiles of drug candidate molecules based on their

chemical structures and their target protein information. We proposed several extensions of kernel regression model for multiple responses to deal with multiple heterogeneous data sources. The originality of the proposed method lies in the integration of chemical space and biological space in a unified framework, in the applicability on a large scale, and in the simultaneous prediction of a large number of potential side-effects at a time. To our knowledge, no previous work gathers all these features in the context of drug side-effect prediction.

The proposed method is expected to be useful in various ways and at various stages of the drug development process. At early stages, among several active drug candidates, the method could help to choose the molecules that should further continue the development process and those that should be dropped. It could also help to find novel drug indications for different diseases, a process named drug repurposing. Indeed, negative side-effects of drugs used in a given pathology can be viewed as a beneficial effect in another pathology. Sildenafil is a famous example of such drug repositioning. In this study we used chemical structure profiles and target protein profiles as predictor data to predict side-effect profiles, but it should be pointed out that many other data sources can be added in the proposed framework. For example, experimental data of binding assays and high-content screening (HCS) from preclinical research can be used. It is possible to integrate all possible information about the drugs or target proteins in the proposed framework, which could improve the prediction reliability.

The proposed method depends highly on the predefinition of chemical profiles, biological profiles, and side-effect terms. Future development could evaluate the performance of using other fingerprints or commercial softwares such as Daylight and Dragon. The algorithm in our proposed methods belong to a class of kernel methods,²⁸ so the performance could be improved by using more sophisticated kernel similarity functions designed for drug structures and target proteins.

Since side-effects are correlated with each other, it is interesting to take into account the correlation of side-effects. Future possible development of a publicly available standard ontology for side-effects would be essential to systematically examine the effect of side-effect correlations on the predictive performance. In addition, not all patients suffer from side-effects; some side-effects occur often, and others rarely occur. Currently only a limited number of drugs have publicly available information on the likelihood of the occurrence of certain side-effects. We hope that such information becomes available for more drugs to conduct more detailed analysis in near future.

In practice, a limitation of the proposed method is that the target protein information is not always available and complete for all drug candidate compounds of interest. Recently, a variety of *in silico* chemogenomic approaches have been developed to predict drug-target interactions or compound-protein interactions. The underlying idea is that similar ligands are likely to interact with similar proteins, and the prediction is performed based on chemical information, genomic information, and all possible known ligand-protein interactions. The usefulness of the chemogenomic approach has been shown in many previous works.^{29–34} The combined use of the chemogenomic approach would be one solution to overcome this limitation, but it is out of scope of this paper.

From a biological viewpoint, drug side-effects are determined by drug chemical fragments, the pattern of target proteins and

the underlying biological pathways. In this sense, it would be interesting to investigate the correlation between drug substructures, target proteins and biological pathways in terms of specific side-effects of interest.

■ ASSOCIATED CONTENT

■ Supporting Information

Figures, Tables, and information as mentioned in the text. Additional information about the detailed results is available at <http://cbio.ensmp.fr/%7Eyyamanishi/integ-effect/>. This material is available free of charge via the Internet at <http://pubs.acs.org>.

■ AUTHOR INFORMATION

Corresponding Author

*E-mail: yamanishi@bioreg.kyushu-u.ac.jp.

Notes

The authors declare no competing financial interest.

■ ACKNOWLEDGMENTS

We thank Dr. Jean-Philippe Vert, Dr. Véronique Stoven, and Dr. Gemma May Kirwan for useful discussions.

■ REFERENCES

- (1) Giacomini, K. M.; Krauss, R. M.; Roden, D. M.; Eichelbaum, M.; Hayden, M. R.; Nakamura, Y. When good drugs go bad. *Nature* **2007**, *446* (7139), 975–977.
- (2) Tatonetti, N.; Liu, T.; Altman, R. Predicting drug side-effects by chemical systems biology. *Genome Biol.* **2009**, *10* (9), 238.
- (3) Campillos, M.; Kuhn, M.; Gavin, A.; Jensen, L.; Bork, P. Drug target identification using side-effect similarity. *Science* **2008**, *321* (5886), 263–266.
- (4) Whitebread, S.; Hamon, J.; Bojanic, D.; Urban, L. Keynote review: in vitro safety pharmacology profiling: an essential tool for successful drug development. *Drug Discovery Today* **2005**, *10* (21), 1421–1433.
- (5) Scheiber, J.; Chen, B.; Milik, M.; Sukuru, S.; Bender, A.; Mikhailov, D.; Whitebread, S.; Hamon, J.; Azzaoui, K.; Urban, L.; Glick, M.; Davies, J.; Jenkins, J. Gaining insight into off-target mediated effects of drug candidates with a comprehensive systems chemical biology analysis. *J. Chem. Inf. Model.* **2009**, *49* (2), 308–317.
- (6) Xie, L.; Li, J.; Xie, L.; Bourne, P. Drug discovery using chemical systems biology: identification of the protein-ligand binding network to explain the side effects of CETP inhibitors. *PLoS Comput. Biol.* **2009**, *5*, e1000387.
- (7) Mahé, P.; Ueda, N.; Akutsu, T.; Perret, J.; Vert, J. Graph kernels for molecular structure-activity relationship analysis with support vector machines. *J. Chem. Inf. Model.* **2005**, *45*, 939–951.
- (8) Varnek, A.; Kireeva, N.; Tetko, I.; Baskin, I.; Solov'ev, V. Exhaustive QSPR studies of a large diverse set of ionic liquids: how accurately can we predict melting points? *J. Chem. Inf. Model.* **2007**, *47*, 1111–1122.
- (9) Varnek, A. Fragment descriptors in structure-property modeling and virtual screening. *Methods Mol. Biol.* **2011**, *672*, 213–243.
- (10) Scheiber, J.; Jenkins, J.; Sukuru, S.; Bender, A.; Mikhailov, D.; Milik, M.; Azzaoui, K.; Whitebread, S.; Hamon, J.; Urban, L.; Glick, M.; Davies, J. Mapping adverse drug reactions in chemical space. *J. Med. Chem.* **2009**, *52* (9), 3103–3107.
- (11) Yamanishi, Y.; Kotera, M.; Kanehisa, M.; Goto, S. Drug-target interaction prediction from chemical, genomic and pharmacological data in an integrated framework. *Bioinformatics* **2010**, *26*, i246–i254.
- (12) Atias, N.; Sharan, R. An algorithmic framework for predicting side-effects of drugs. *J. Comput. Biol.* **2011**, *18*, 207–218.
- (13) Pauwels, E.; Stoven, V.; Yamanishi, Y. Predicting drug side-effect profiles: a chemical fragment-based approach. *BMC Bioinf.* **2011**, *12*, 169.
- (14) Kuhn, M.; Campillos, M.; Letunic, I.; Jensen, L.; Bork, P. A side effect resource to capture phenotypic effects of drugs. *Mol. Syst. Biol.* **2010**, *6*, 343.
- (15) Wishart, D.; Knox, C.; Guo, A.; Cheng, D.; Shrivastava, S.; Tzur, D.; Gautam, B.; Hassanali, M. DrugBank: a knowledgebase for drugs, drug actions and drug targets. *Nucleic Acids Res.* **2008**, *36*, D901–D906.
- (16) Chen, B.; Wild, D.; Guha, R. PubChem as a Source of Polypharmacology. *J. Chem. Inf. Model.* **2009**, *49* (9), 2044–2055.
- (17) Gunther, S.; Kuhn, M.; Dunkel, M.; Campillos, M.; Senger, C.; Petsalaki, E.; Ahmed, J.; Urdiales, G.; Gewiss, A.; Jensen, L.; Schneider, R.; Skoblo, R.; Russell, R.; Bourne, P.; Bork, P.; et al. SuperTarget and Matador: resources for exploring drug-target relationships. *Nucleic Acids Res.* **2008**, *36*, D919–D922.
- (18) The Uniprot Consortium, The Universal Protein Resource (UniProt) in 2010. *Nucleic Acids Res.* **2010**, *38*, D142–D148.
- (19) Yamanishi, Y.; Vert, J.-P.; Kanehisa, M. Supervised Enzyme Network Inference from the Integration of Genomic Data and Chemical Information. *Bioinformatics* **2005**, *21*, i468–i477.
- (20) Hotelling, H. Relations between two sets of variates. *Biometrika* **1936**, *28*, 321–377.
- (21) Davis, J.; Goadrich, M. The Relationship Between Precision-Recall and ROC Curves. In *Proceedings of the Twenty Third International Conference on Machine Learning*; Pittsburgh, PA, June 25–29, 2006; W.W. Cohen, A. M., Ed.; ACM Press: PA, 2006; pp 233–240.
- (22) Watson Pharma, Inc., NEXT CHOICE [package insert]. <http://dailymed.nlm.nih.gov/dailymed/lookup.cfm?setid=14052755-0ffd-419> (accessed Oct 2012).
- (23) Hexal Pharma (SA) (PTY) Ltd, NORLEVO [package insert]. <http://www.sapajournal.co.za/index.php/SAPA/article/download/151/143>, 2006 (accessed Oct 2012).
- (24) MERCK & CO., INC., CELESTONE SOLUSPAN [package insert]. <http://dailymed.nlm.nih.gov/dailymed/lookup.cfm?setid=7b5489a1-e30f-450> (accessed Oct 2012).
- (25) BERTEK PHARMACEUTICALS INC., ZAGAM [package insert]. http://www.accessdata.fda.gov/drugsatfda_docs/label/2003/020677s0061bl (accessed Oct 2012).
- (26) Aurobindo Pharma USA, Inc., ESCITALOPRAM OXALATE [package insert]. <http://dailymed.nlm.nih.gov/dailymed/lookup.cfm?setid=ba60c357-8e61-4a3> (accessed Oct 2012).
- (27) Pfizer Inc., CYTARABINE [package insert]. <http://dailymed.nlm.nih.gov/dailymed/lookup.cfm?setid=1315f89d-f530-4df> (accessed Oct 2012).
- (28) Schölkopf, B.; Tsuda, K.; Vert, J.-P. *Kernel Methods in Computational Biology*; MIT Press: Cambridge, MA, 2004.
- (29) Bock, J. R.; Gough, D. A. Virtual screen for ligands of orphan G protein-coupled receptors. *J. Chem. Inf. Model.* **2005**, *45*, 1402–1414.
- (30) Erhan, D.; L'Heureux, P.-J.; Yue, S. Y.; Bengio, Y. Collaborative filtering on a family of biological targets. *J. Chem. Inf. Model.* **2006**, *46*, 626–635.
- (31) Yamanishi, Y.; Araki, M.; Gutteridge, A.; Honda, W.; Kanehisa, M. Prediction of drug-target interaction networks from the integration of chemical and genomic spaces. *Bioinformatics* **2008**, *24*, i232–i240.
- (32) Faulon, J.; Misra, M.; Martin, S.; Sale, K.; Sapra, R. Genome scale enzyme-metabolite and drug-target interaction predictions using the signature molecular descriptor. *Bioinformatics* **2008**, *24*, 225–233.
- (33) Jacob, L.; Vert, J.-P. Protein-ligand interaction prediction: an improved chemogenomics approach. *Bioinformatics* **2008**, *24*, 2149–2156.
- (34) Bleakley, K.; Yamanishi, Y. Supervised prediction of drug-target interactions using bipartite local models. *Bioinformatics* **2009**, *25*, 2397–2403.