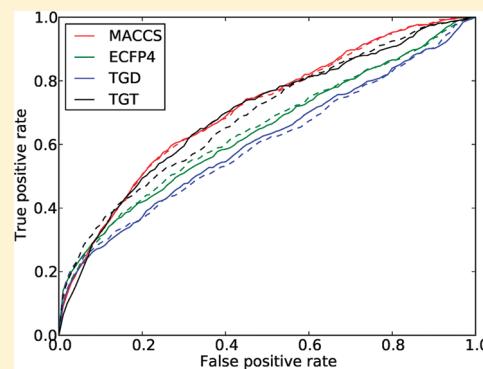


Introduction of the Conditional Correlated Bernoulli Model of Similarity Value Distributions and its Application to the Prospective Prediction of Fingerprint Search Performance

Martin Vogt and Jürgen Bajorath*

Department of Life Science Informatics, B-IT, LIMES Program Unit Chemical Biology and Medicinal Chemistry, Rheinische Friedrich-Wilhelms-Universität Bonn, Dahlmannstrasse 2, D-53113 Bonn, Germany

ABSTRACT: A statistical approach named the *conditional correlated Bernoulli model* is introduced for modeling of similarity scores and predicting the potential of fingerprint search calculations to identify active compounds. Fingerprint features are rationalized as dependent Bernoulli variables and conditional distributions of Tanimoto similarity values of database compounds given a reference molecule are assessed. The conditional correlated Bernoulli model is utilized in the context of virtual screening to estimate the position of a compound obtaining a certain similarity value in a database ranking. Through the generation of receiver operating characteristic curves from cumulative distribution functions of conditional similarity values for known active and random database compounds, one can predict how successful a fingerprint search might be. The comparison of curves for different fingerprints makes it possible to identify fingerprints that are most likely to identify new active molecules in a database search given a set of known reference molecules.



1. INTRODUCTION

In virtual screening, methods are typically evaluated retrospectively in benchmark calculations using sets of known active compounds and background databases assumed to contain only inactive compounds.¹ However, it is much more challenging to estimate (ligand-based) virtual screening performance in advance of practical applications. For example, for a given search method, active reference compounds, and a screening database, it would be interesting and helpful to estimate how likely it might be to identify new active compounds in a virtual screen. For similarity-based methods, such predictions are complicated by the fact that there is no well-defined relationship between molecular similarity (however calculated) and similarity with respect to biological activity of test compounds.² In other words, calculated similarity values cannot be reliably utilized as an indicator of biological activity. Thus, methods that attempt to assign activity probabilities to calculated similarity values or predict similarity search performance must essentially derive statistical relationships from molecular calibration sets and fit predictive models to activity data. However, such methods are currently rare. For the prediction of virtual screening performance, an approach has been introduced that combines Kullback–Leibler divergence analysis from information theory³ with Bayesian modeling to estimate the ability of different descriptors to discriminate between active and inactive compounds.⁴ The Kullback–Leibler divergence between activity classes and a screening database was shown to scale with the recall rates of active compounds in Bayesian screening using property descriptors⁴ or in similarity searching using fingerprints.⁵

On the basis of these findings, linear models were calibrated to predict the fingerprint search performance for new compound classes.⁵ In order to assign probabilities of activity to calculated similarity values, another approach has been introduced that generates probability assignment curves based on the pairwise similarity of active and database compounds calculated using different fingerprints and similarity measures.⁶ Similarity values are then converted to probabilities of activity utilizing the concept of belief theory, hence subjecting different compound rankings to data fusion yielding a quantitative estimate for the probability of activity.⁶ This data fusion approach has also been extended to structure-based virtual screening where scoring functions replace similarity measures.⁷ Although the methodologies that are based on Kullback–Leibler divergence and belief theory differ conceptually, they both require different sets of known active and database compounds for model generation and fitting.

Herein, we present a methodological concept for statistical modeling of similarity value distributions on the basis of fingerprint features. Several approaches to derive similarity value distributions from fingerprint feature parameters have been introduced,⁸ and we utilize a similar route to model similarity values (vide infra). Going beyond the estimation of similarity value distributions, the method introduced herein can be readily applied to predict fingerprint search performance in a prospective manner. Different from the two previously introduced

Received: July 25, 2011

Published: September 05, 2011

approaches discussed above, the newly developed methodology does not require different data sets of active compounds for model calibration and is hence more general and less dependent on prior knowledge of activity data and the specifics of calibration sets.

2. METHODOLOGY

2.1. Distribution of Fingerprint Features and Similarity Values.

Fingerprints are bit string representations of molecular structure and properties.² Depending on the fingerprint design, individual bit positions account for different (e.g., structural, topological, or pharmacophoric) features. In the following, we will refer to individual fingerprint bits simply as features. For a 2D fingerprint such as MACCS⁹ that consists of substructural features the following observations concerning feature frequency can generally be made.

1. Some features are extremely frequent in compounds, for example, contains a “ring” or “nitrogen”, whereas others are rare, for example, contains an “iodine” or “four-membered ring”.
2. Other features often occur in many molecules (but not with extreme frequency), for example, “seven-membered ring”, “primary amine”, or “hydroxyl group”.
3. Many features are (in part strongly) correlated, for example, “ring” and “six-membered ring”, or “more than one oxygen” and “hydroxyl group”.

Thus, a fingerprint calculated for a given molecule typically consists of more or less frequent features that might or might not be (positively) correlated with others (more often than not, correlation will be present).

The quantitative comparison of fingerprint representations requires the application of similarity metrics. Most of the popular similarity measures¹⁰ used for fingerprint comparison rely on set-theoretic magnitudes. For instance, the popular Tanimoto coefficient (T_c)¹⁰ can be understood as the ratio of the cardinalities of the intersection and union of two fingerprints.

There are many different fingerprint designs² with different feature characteristics that influence the distribution of similarity scores. The main numerical parameters influencing fingerprint search behavior are as follows:

1. The *feature frequency*. Depending on the fingerprint, features account for many different structural or molecular properties to which test compounds respond in different ways. The frequencies with which features are found in database compounds determine the average cardinalities of the intersection and union of two fingerprints and thus are central parameters of the score distribution.
2. The *feature correlation*. Depending on the fingerprint, features might not occur independently of each other. Frequently, individual features are correlated. This is especially the case for combinatorial fingerprint representations such as pharmacophore pattern fingerprints that enumerate all possible constellations of two, three, or four pharmacophore features. In this case, systematic correlation between features is inherent in the design. Correlations between features play a major role in the variance of the score distributions, that is, they determine how much similarity scores tend to vary from the average.
3. The *fingerprint size*, that is, the total number of possible fingerprint features. The size can be constant ranging from

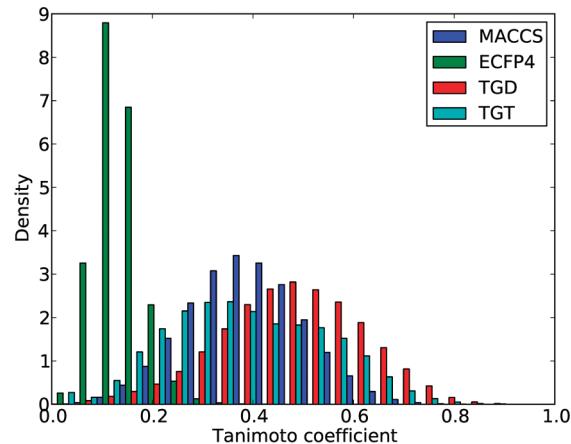


Figure 1. Similarity value distributions for different fingerprints. Shown is the distribution of T_c values for four different fingerprints in the ZINC database using a reference compound randomly selected from ZINC. Fingerprint descriptions are provided in section 3.2.

less than 100 to tens of thousands or even millions of features. Alternatively, the number of features can be variable and sometimes virtually unlimited, for example, when features are represented by arbitrary integers, a theoretical total of $2^{32} \approx 4.3 \times 10^9$ possibilities would exist.

4. The *fingerprint density*, that is, the number of features present in a molecule over the total number of possible features.

These factors often lead to very different similarity value (score) distributions, as illustrated in Figure 1, that is, different fingerprint designs yield different distributions of T_c values in database searching.

In a fingerprint search, database compounds are ranked according to their similarity to one or more given reference compounds and similarity score distributions typically vary depending on the fingerprint settings (i.e., feature distribution and density) of the reference compounds. Figure 2 illustrates the dependency of T_c value distributions on a chosen reference compound. Depending on the nature of this distribution, individual T_c values (or value ranges) have different significance for compound ranking. As a general rule, the higher the densities of fingerprints are, the more the database T_c distributions are shifted toward higher values.

2.2. Modeling the Distribution of Similarity Values. We begin with the mathematical formulation. For a molecule A we define the following:

1. $BV(A) = \mathbf{a}$ is the *bit vector* of the fingerprint, that is, $\mathbf{a} = (a_i)_{i=1,\dots,n}, a_i \in \{0,1\}$ for a fingerprint with n features.
2. $FS(A) = \mathcal{A}$ is the *feature set* of the fingerprint, that is, $\mathcal{A} = \{i | a_i = 1, 1 \leq i \leq n\}$
3. $|A| = |FS(A)| = \sum_{i=1}^n a_i$ is the cardinality of the feature set for fingerprint A (cardinality refers to the size of a set).

In addition, we define the cardinalities of the intersection and the union of two fingerprints for molecules A and B :

1. $I(A,B) = |FS(A) \cap FS(B)| = \sum_{i=1}^n a_i b_i$ is the cardinality of the intersection of two fingerprints.
2. $U(A,B) = |FS(A) \cup FS(B)| = |A| + |B| - I(A,B) = \sum_{i=1}^n a_i + \sum_{i=1}^m b_i - \sum_{i=1}^n a_i b_i$ is the cardinality of the union of two fingerprints.
3. $Tc(A,B) = I(A,B)/U(A,B)$ is the Tanimoto coefficient for the two molecules with respect to the fingerprint.

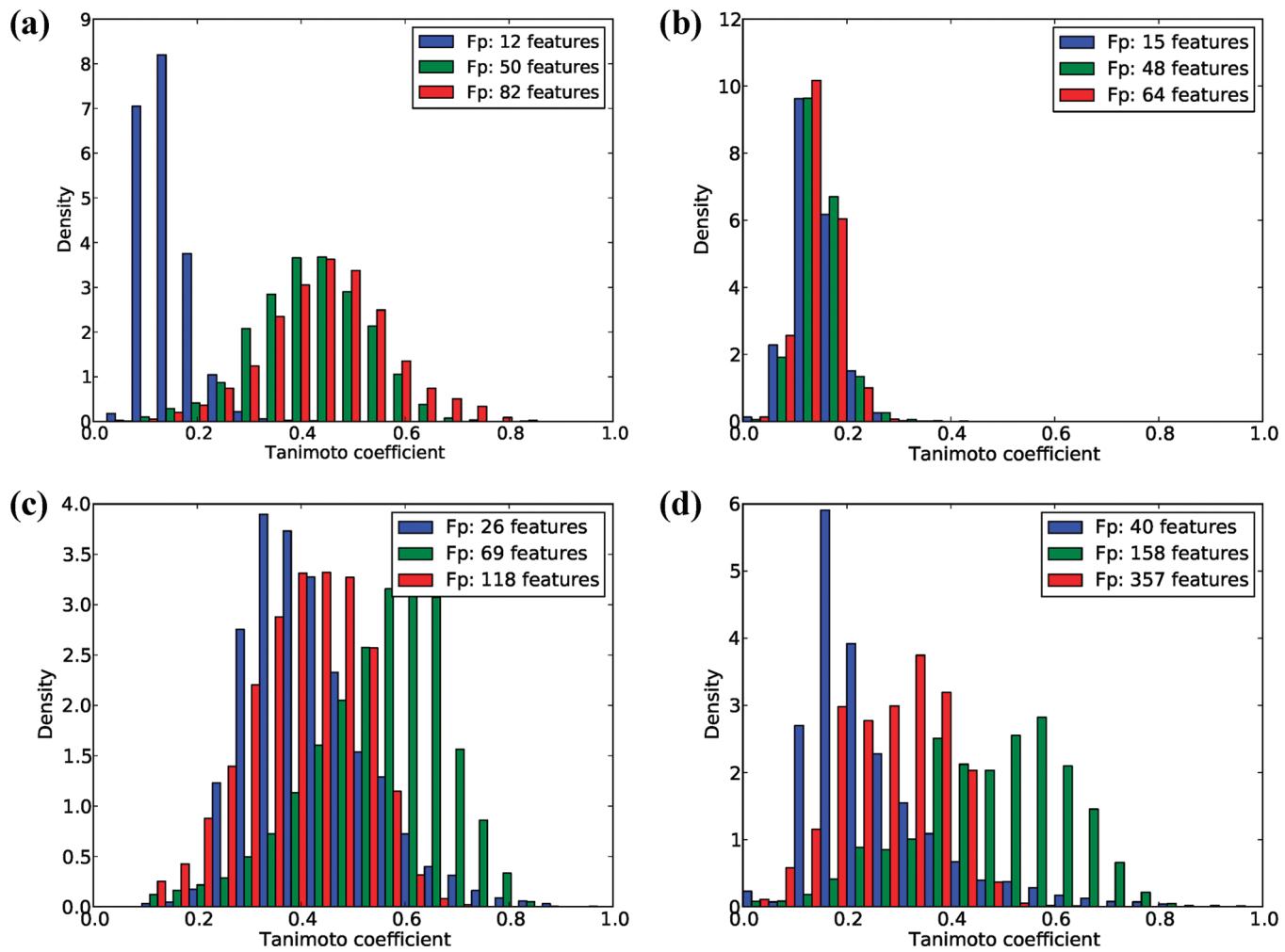


Figure 2. Similarity value distributions for different reference compounds. Shown is the distribution of Tc values for (a) MACCS, (b) ECFP4, (c) TGD, and (d) TGT fingerprints in the ZINC database using three different reference compounds randomly selected from ZINC. For each reference compound, the number of fingerprint features it produces is reported. Fingerprint descriptions are provided in section 3.2.

Next, we focus on model derivation. Here, the primary objective is to model the distribution of Tc values on the basis of probability distributions. In mathematical terms, the underlying question is how $Tc(X, Y)$ is distributed if X and Y are two “random” molecules or, more specifically and interestingly, how $Tc(A, X)$ is distributed for a chosen reference molecule A and any (random) database molecule X. The first question considers the overall distribution of Tc values and the second the *conditional distribution* of these values given a reference compound. As illustrated above, these distributions might significantly differ. Importantly, the conditional distribution reflects how Tc values are distributed in a typical similarity search where a known active reference compound is utilized. On the basis of a conditional distribution model, the significance of Tc values with respect to the database distribution, that is, the significance for compound ranking, can be quantitatively assessed, which makes it also possible to compare Tc values across different fingerprints and reference compounds.

To obtain a meaningful model, we must first specify what “random” means in this context. It is straightforward to model fingerprints as a vector of Bernoulli variables $(X_i)_{i=1,\dots,n}$ where each feature X_i has a probability of $p_i \in [0,1]$, $i = 1, \dots, n$, that is,

$\Pr(X_i = 1) = p_i$. For the Bernoulli variables, independence is not assumed. Instead, these variables might have a nonzero covariance, taking the correlation between fingerprint features into account. We denote the covariance of X_i and X_j by $c_{ij} \in [-0.25, 0.25]$, $i, j = 1, \dots, n$.

For the derivation of this model, the feature probabilities and covariances should best be estimated from a large and representative compound database. The nature of the model as a vector of Bernoulli variables makes it possible to analytically determine expectation values and variances for the cardinalities of intersections and unions of fingerprints. These cardinalities can be expressed as the sum of Bernoulli variables. As such, the distributions of the cardinalities of intersections and unions are approximated by normal distributions. One should note here that the quality of this approximation cannot be guaranteed by the central limit theorem because the Bernoulli variables are not independent. However, the presence of dependence does not necessarily invalidate this approximation. In fact, evaluation of the model reveals that the resulting distribution accurately represents Tc values for different fingerprint types (vide infra).

By the equation given above, the Tanimoto coefficient is determined as the ratio of the cardinalities of the intersection

and the union of two fingerprints (vide supra). As such, its distribution is modeled as the ratio of the two normal distributions. Here, it should be noted that the intersection and union will generally not be independent of each other and hence the correlation between these two cardinalities also needs to be taken into account.

Previously, Baldi and Nasr⁸ have also modeled Tc values from ratios of normal distributions. In contrast to the conditional correlated Bernoulli model presented herein, the main result of Baldi and Nasr focused on assuming a conditional normal uniform model for the feature distribution. This assumption implies that the cardinality of fingerprints can be modeled by normal distributions and furthermore that each fingerprint of a given cardinality is equally probable. Therefore, this model is not applicable to fingerprints where individual features have very different frequencies. In such cases, the conditional normal uniform model tends to significantly underestimate the cardinality of the intersection (and thus overestimate the cardinality of the union), thereby yielding much lower theoretical Tc values than observed for such fingerprints in practical applications. However, Baldi and Nasr have not attempted to compare the performance of different fingerprint representations based on their model.

2.3. Conditional Correlated Bernoulli Model. Let us assume a reference compound A with $FS(A) = \mathcal{A}$ is given. First, we determine the model for the intersection $I(A,X)$ given reference A. To model the cardinality of the intersection as a normal distribution, the mean and the variance of the distribution need to be determined, that is, $E(I(A,X)|A)$ and $\text{var}(I(A,X)|A)$. For a random molecule X, the bit vector $BV(X)$ is a vector of random Bernoulli variables $(X_i)_{i=1,\dots,n}$ with feature probabilities $(p_i)_{i=1,\dots,n}$. Furthermore, let \mathcal{X} denote the feature set of X. For each feature $i \in \mathcal{A}$, the feature is part of the intersection $\mathcal{A} \cap \mathcal{X}$ if and only if it is present in the feature set \mathcal{X} . Thus

$$\mu_I^A = E(I(A,X)|A) = E\left(\sum_{i \in \mathcal{A}} X_i\right) = \sum_{i \in \mathcal{A}} E(X_i) = \sum_{i \in \mathcal{A}} p_i$$

For the variance, we obtain

$$\begin{aligned} (\sigma_I^A)^2 &= \text{var}(I(A,X)|A) = \text{var}\left(\sum_{i \in \mathcal{A}} X_i\right) = \sum_{i,j \in \mathcal{A}} \text{cov}(X_i, X_j) \\ &= \sum_{i \in \mathcal{A}} \text{var}(X_i) + \sum_{i,j \in \mathcal{A}, i \neq j} \text{cov}(X_i, X_j) \\ &= \sum_{i \in \mathcal{A}} p_i(1-p_i) + \sum_{i,j \in \mathcal{A}, i \neq j} c_{ij} \end{aligned}$$

The union $U(A,X)$ given reference A can be derived analogously. Again, the mean and variance must be determined. The cardinality of the union $\mathcal{A} \cup \mathcal{X}$ consists of the $a = |\mathcal{A}|$ features of the reference A plus the features in \mathcal{X} that do not occur in \mathcal{A} . The expected number of features of \mathcal{X} that are not in \mathcal{A} is $E(\sum_{i \notin \mathcal{A}} X_i) = \sum_{i \notin \mathcal{A}} p_i$ so that

$$\mu_U^A = E(U(A,X)|A) = E(a + \sum_{i \notin \mathcal{A}} X_i) = a + \sum_{i \notin \mathcal{A}} p_i$$

In analogy to the intersection the variance of the cardinality of the union is

$$\begin{aligned} (\sigma_U^A)^2 &= \text{var}(U(A,X)|A) = \text{var}\left(\sum_{i \notin \mathcal{A}} X_i\right) = \sum_{i,j \notin \mathcal{A}} \text{cov}(X_i, X_j) \\ &= \sum_{i \notin \mathcal{A}} p_i(1-p_i) + \sum_{i,j \notin \mathcal{A}, i \neq j} c_{ij} \end{aligned}$$

Finally, the covariance between $I(A,X)$ and $U(A,X)$ given reference A is determined as

$$\begin{aligned} \text{cov}_{IU}^A &= \text{cov}(I(A,X), U(A,X)|A) \\ &= \text{cov}\left(\sum_{i \in \mathcal{A}} X_i, \sum_{i \notin \mathcal{A}} X_i\right) = \sum_{i \in \mathcal{A}, j \notin \mathcal{A}} \text{cov}(X_i, X_j) \\ &= \sum_{i \in \mathcal{A}, j \notin \mathcal{A}} c_{ij} \end{aligned}$$

A key feature of the conditional correlated Bernoulli model is that it makes it possible to derive Tc distributions for the database molecules that are conditional on the fingerprint feature settings of a given reference compounds. The ensuing distribution then provides the basis to estimate the significance of a Tc value obtained for any other test compound (relative to the reference) for compound ranking (vide infra).

The conditional model applies when a reference compound is known. This usually is the situation in virtual screening. However, one can also ask how scores are distributed when no reference molecule is given, i.e. how scores are distributed for two random molecules from a database. In the following, we derive the relevant parameters for the *unconditional* correlated Bernoulli model that yields a distribution of Tc values for two random molecules. Because of the dependence of Tc value distributions on the reference compound, the unconditional model is unsuitable for the prediction of ranks in a similarity search. Nevertheless, the resulting unconditional model provides a meaningful reference for the evaluation of the correlated Bernoulli model.

For the generation of the unconditional model, we begin with the intersection. For two random molecules X and Y with bit vectors $(X_i)_{i=1,\dots,n}$ and $(Y_i)_{i=1,\dots,n}$ and corresponding feature sets \mathcal{X} and \mathcal{Y} , the size of the intersection is given by $\sum_{i=1}^n X_i Y_i$ where X_i and Y_i are (pairwise independent) Bernoulli variables with feature probability p_i .

$$\begin{aligned} \mu_I &= E(I(X,Y)) = E\left(\sum_{i=1}^n X_i Y_i\right) = \sum_{i=1}^n E(X_i) E(Y_i) \\ &= \sum_{i=1}^n p_i^2 \end{aligned}$$

The variance then is

$$\begin{aligned} \sigma_I^2 &= \text{var}(I(X,Y)) = \text{var}\left(\sum_{i=1}^n X_i Y_i\right) = \sum_{i=1}^n \sum_{j=1}^n \text{cov}(X_i Y_i, X_j Y_j) \\ &= \sum_{i=1}^n \sum_{j=1}^n (E(X_i Y_i X_j Y_j) - E(X_i Y_i) E(X_j Y_j)) \\ &= \sum_{i=1}^n \sum_{j=1}^n (E(X_i X_j) E(Y_i Y_j) - E(X_i Y_i) E(X_j Y_j)) \end{aligned}$$

Given

$$\begin{aligned} E(X_i Y_i) &= p_i^2 \text{ and } E(X_i X_j) = \text{cov}(X_i, X_j) + E(X_i) E(X_j) \\ &= c_{ij} + p_i p_j \end{aligned}$$

we obtain

$$\begin{aligned} \sigma_I^2 &= \text{var}(I(X,Y)) = \sum_{i=1}^n \sum_{j=1}^n ((c_{ij} + p_i p_j)^2 - p_i^2 p_j^2) \\ &= \sum_{i=1}^n \sum_{j=1}^n (c_{ij}^2 + 2c_{ij} p_i p_j) \end{aligned}$$

Next we derive the mean and variance for the union. The size of the union is given by $U(X, Y) = |X| + |Y| - I(X, Y) = \sum_{i=1}^n (X_i + Y_i - X_i Y_i)$

$$\begin{aligned}\mu_U &= E(U(X, Y)) = E\left(\sum_{i=1}^n (X_i + Y_i - X_i Y_i)\right) \\ &= \sum_{i=1}^n (2p_i - p_i^2)\end{aligned}$$

The variance then is

$$\begin{aligned}\sigma_U^2 &= \text{var}(U(X, Y)) = \text{var}\left(\sum_{i=1}^n (X_i + Y_i - X_i Y_i)\right) \\ &= 2 \sum_{i=1}^n \sum_{j=1}^n \text{cov}(X_i, X_j) - 4 \sum_{i=1}^n \sum_{j=1}^n \text{cov}(X_i, X_j Y_j) \\ &\quad + \sum_{i=1}^n \sum_{j=1}^n \text{cov}(X_i Y_i, X_j Y_j)\end{aligned}$$

Using $E(X_i X_j) = c_{ij} + p_i p_j$, we obtain

$$\begin{aligned}\text{cov}(X_i, X_j Y_j) &= E(X_i X_j Y_j) - E(X_i)E(X_j Y_j) \\ &= E(X_i X_j)E(Y_j) - E(X_i)E(X_j Y_j) \\ &= (c_{ij} + p_i p_j)p_j - p_i p_j^2 = c_{ij}p_j\end{aligned}$$

So for σ_U^2 , we get

$$\begin{aligned}\sigma_U^2 &= \text{var}(U(X, Y)) = \text{var}\left(\sum_{i=1}^n (X_i + Y_i - X_i Y_i)\right) \\ &= 2 \sum_{i=1}^n \sum_{j=1}^n c_{ij} - 4 \sum_{i=1}^n \sum_{j=1}^n c_{ij}p_j + \sigma_I^2\end{aligned}$$

The covariance then is

$$\begin{aligned}\text{cov}_{IU} &= \text{cov}(I(X, Y), U(X, Y)) \\ &= \text{cov}\left(\sum_{i=1}^n X_i Y_i, \sum_{i=1}^n (X_i + Y_i - X_i Y_i)\right) \\ &= 2 \sum_{i=1}^n \sum_{j=1}^n \text{cov}(X_i, X_j Y_j) - \sum_{i=1}^n \sum_{j=1}^n \text{cov}(X_i Y_i, X_j Y_j) \\ &= 2 \sum_{i=1}^n \sum_{j=1}^n c_{ij}p_j - \sigma_I^2\end{aligned}$$

2.4. Similarity Value Distribution and Rank Probability. In $T_c = I/U$, I is the cardinality of the intersection and U the cardinality of the union of two fingerprints. The distributions of I and U are modeled by normal distributions $I \approx N(\mu_I, \sigma_I^2)$ and $U \approx N(\mu_U, \sigma_U^2)$ with correlation $\rho = \text{cov}(I, U) / (\sigma_I \sigma_U)$. The probability density function of the T_c can then be given in analytical form as described previously.^{8,11,12} To assess the significance of individual T_c values for the prediction of rank positions in a fingerprint search, the cumulative distribution function (CDF) is utilized, which can be approximated as follows:¹¹

$$F(t) \approx \Phi\left(\frac{\mu_U t - \mu_I}{\sigma_I \sigma_U a(t)}\right)$$

where

$$a(t) = \left(\frac{t^2}{\sigma_I^2} - \frac{2\rho t}{\sigma_I \sigma_U} + \frac{1}{\sigma_U^2}\right)^{1/2}$$

$\Phi(u) = \int_{-\infty}^u (2\pi)^{-1/2} \exp(-x^2/2) dx$ is the CDF of the standard normal distribution.

From the CDF $F(t) = \Pr(T_c \leq t)$, a p -value can immediately be calculated as

$$p = 1 - F(t) = \Pr(T_c > t)$$

giving the probability of achieving a T_c value larger than t . The p -value enables us to estimate how highly a database compound obtaining a certain T_c value will be ranked in a similarity search: rank $\sim pN$ for a database of size N .

2.5. Prediction of Fingerprint Search Performance. On the basis of the CDF and the ensuing p -value, statistical measures like the receiver operating characteristic (ROC) curve,¹³ area under the curve (AUC), or compound recovery rates can be estimated from the distribution models. This approach has previously been described by Baldi and Nasr.⁸ A ROC curve can be generated from the distribution of scores for active and database compounds. If $F_{FP}^A(t)$ is the CDF of T_c values for a fingerprint FP and a reference compound A for screening database compounds and $G_{FP}^A(t)$ is the corresponding CDF of T_c values for other known active compounds, the ROC curve is given by the parametric equations

$$x = 1 - F_{FP}^A(t) \text{ and } y = 1 - G_{FP}^A(t)$$

In principle, the CDFs can either be obtained through empirical determination of the T_c distribution or by estimating the distributions using the conditional Bernoulli model. The empirical approach requires the determination of the distribution by screening the database for each reference compound A. By contrast, the conditional model can easily be applied to estimate the distribution for different reference compounds.

The prediction of the performance of a fingerprint search for a specific class of active compounds then involves the following steps:

1. For a given fingerprint and the screening database, the conditional correlated Bernoulli model is derived. Therefore, the feature probabilities and their covariances need to be determined on the basis of the screening database. For a database of N molecules $(A_k)_{k=1,\dots,N}$ with fingerprints $\mathbf{a}_k = (a_{ki})_{i=1,\dots,m}$, $a_{ki} \in \{0,1\}$ the feature probabilities are estimated by

$$p_i = \frac{1}{N} \sum_{k=1}^N a_{ki}$$

and the covariance between feature i and feature j can be estimated by

$$c_{ij} = \frac{1}{N} \sum_{k=1}^N a_{ki} a_{kj} - p_i p_j$$

Here, probabilities and covariances are calculated from the large ZINC¹⁴ database, which represents a comprehensive collection of currently available pharmaceutically relevant small molecules. Note that the parameters of the conditional model only need to be derived once for a given fingerprint type and screening database.

It should be noted that the estimation of all pairwise covariances is practically feasible for fingerprints of up to a few thousand features (bit positions). For much larger fingerprints, it is infeasible to calculate all pairwise covariances. In these cases pairwise covariances are calculated

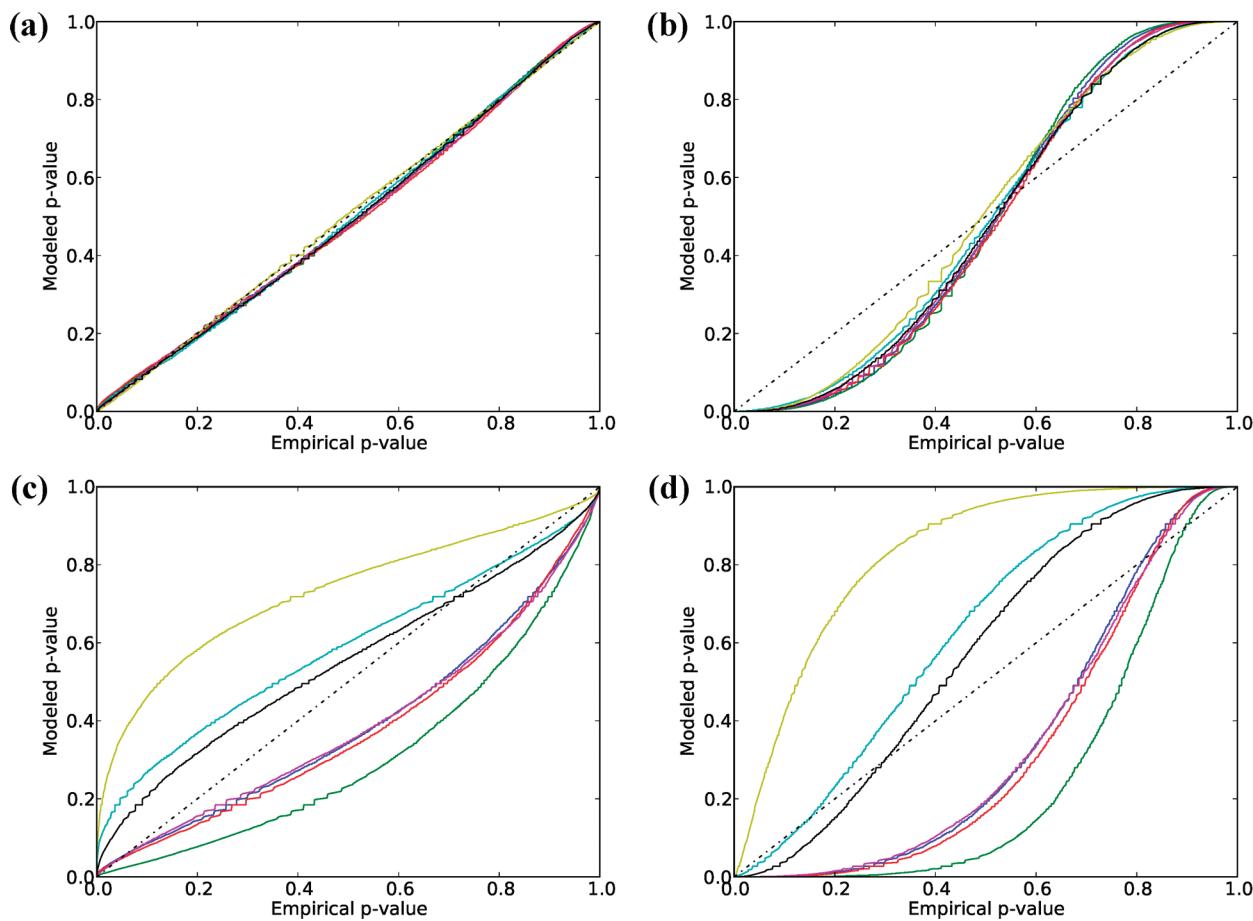


Figure 3. Modeled and empirical *p*-values for different models. Graphs are shown where the predicted *p*-value (y-axis) is plotted against the empirical *p*-value (x-axis) for different score distribution models for the MACCS fingerprint. Each curve was obtained by selecting a randomly chosen ZINC compound as a reference and by calculating the Tc values for a subset of 10 000 random ZINC compounds. The compounds were ranked in decreasing order of their Tc values yielding the empirical *p*-value, that is, the rank corresponding to a specific Tc-value. On the y-axis the corresponding *p*-value resulting from a predictive model is shown. A perfect model would yield a straight diagonal from the bottom left to the top right (indicated as a dashed line). For each model, seven differently colored curves are drawn that correspond to seven randomly selected compounds. Graphs are shown for the following models: (a) correlated conditional Bernoulli model, (b) independent conditional Bernoulli model (i.e., assuming the absence of correlation between features), (c) unconditional correlated Bernoulli model (i.e., assuming independence of the score distribution from the reference compound), and (d) unconditional independent Bernoulli model (i.e., assuming both the absence of correlation between features and independence of the score distribution from the reference compound).

only for the most frequently occurring features. For instance, most of the features of large fingerprints like ECFP4¹⁵ are extremely rare and their covariances are very small. In this case, the pairwise covariances between the 2000 most frequently occurring features were estimated according to the formula above. Furthermore, for each of the 2000 most frequent features, the average covariance of a frequent feature and all infrequent features was estimated. Finally, the average covariance between two infrequent features was also estimated. Overall, covariances involving infrequent features only contributed marginally to the parameter estimates. Thus, the accuracy of the models would not be expected to degenerate notably even if these covariances were completely ignored.

- For a set of known active compounds, each compound is taken once as the reference compound to derive the conditional Tc value distribution of the screening database using the conditional correlated Bernoulli model. For the remaining active compounds, the Tc values relative to

the reference compound are calculated and their database ranks are estimated (*vide supra*).

- On the basis of the union of all rank positions, a ROC curve is calculated for this class of active compounds. From the ROC curve or AUC, it can be predicted how likely it is to identify other active compounds that might be available in the screening database (the underlying assumption being that these potential hits are similar in fingerprint space to the known active compounds).
- The comparison of ROC curves for different fingerprints and a set of known active compounds makes it possible to prioritize fingerprints that are expected to yield best search performance.

3. MODEL EVALUATION

3.1. Different Variants of the Bernoulli Model.

To assess the quality of Tc value distributions derived via the correlated conditional Bernoulli model, CDFs of Tc value distributions

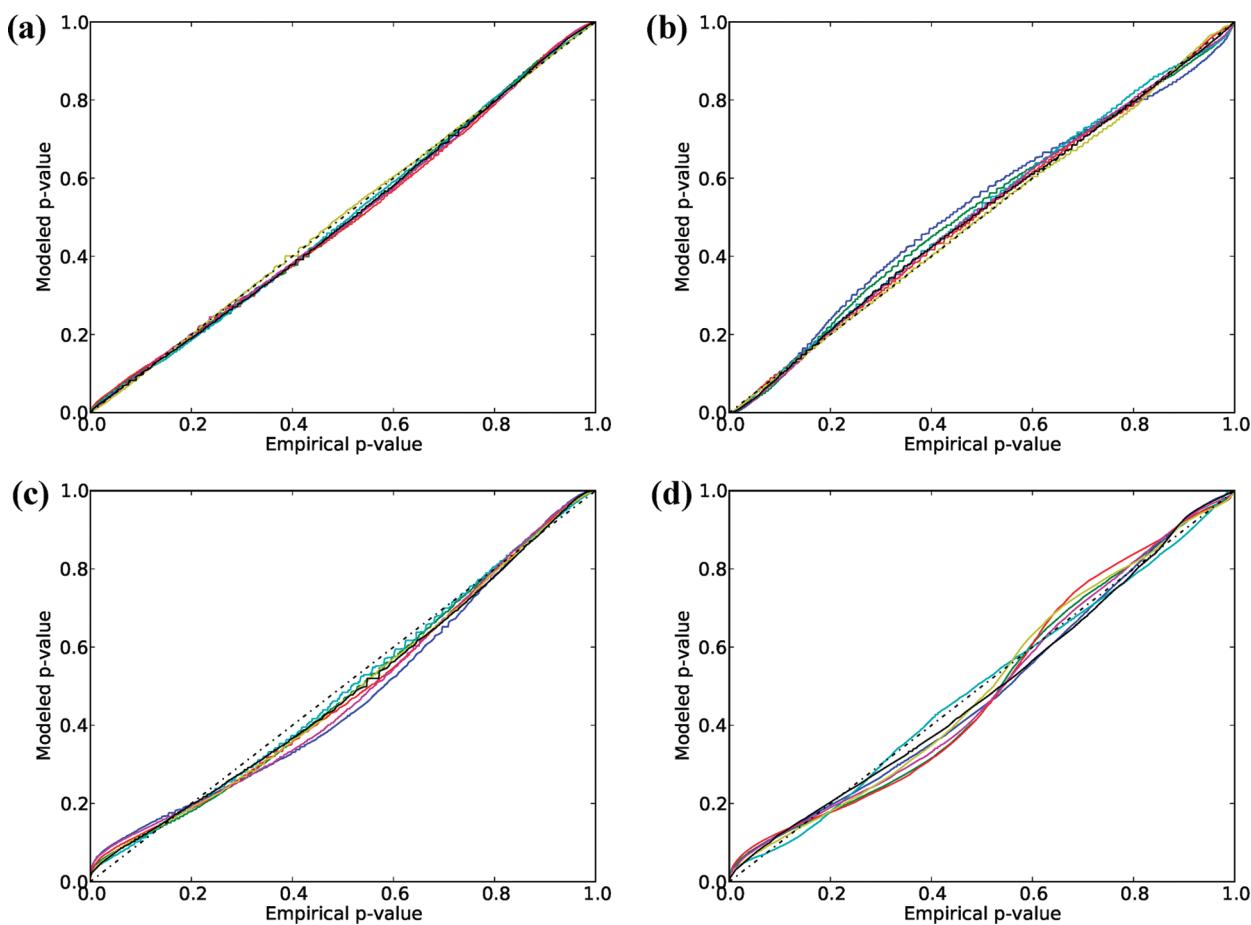


Figure 4. Modeled and empirical *p*-values for different fingerprints. Graphs are shown according to Figure 3 for the conditional correlated Bernoulli model where the predicted *p*-value (*y*-axis) is plotted against the empirical *p*-value (*x*-axis) for different fingerprints: (a) MACCS (same as in Figure 3a), (b) ECFP4, (c) TGD, and (d) TGT.

were empirically determined and compared to modeled CDFs. Simplified distribution models without conditional distribution modeling or without feature correlation were also evaluated. For these calculations, seven randomly chosen ZINC compounds were independently utilized as reference compounds and 10 000 other random ZINC molecules as the database. For these comparisons, MACCS was consistently utilized as a fingerprint. The empirical score distribution for a given reference compound was calculated for the 10 000 database compounds. By ranking the compounds in the order of decreasing T_c values, empirical *p*-values for similarity scores were determined. The empirical *p*-values were then compared to the *p*-values predicted by the model. The results for different models are reported in Figure 3. The quality of a modeled distribution can be judged by its closeness to the diagonal in the graph. The correlated conditional Bernoulli model yielded excellent agreement between empirical and predicted *p*-values, as shown in Figure 3a.

In Figure 3b, the results are shown for a simplified conditional model where the Bernoulli variables are assumed to be independent. The use of this model led to a significant underestimation of the variances for the cardinalities of the intersection and union, which resulted in a typical sigmoidal shape of the curves. In Figure 3c, the unconditional correlated Bernoulli model is applied, which does not consider the distribution dependence on the reference compound. The graphs clearly demonstrate that

conditioning of the model was essential for its quality. Individual unconditional models under- and/or overestimated the *p*-values, depending on the reference compound. Figure 3d reports the results obtained for the uncorrelated plus unconditional Bernoulli models where deviation from empirical values even further increased, as one would expect. Taken together, the results clearly demonstrate the predictive value of the correlated conditional Bernoulli model and the need to take into account both correlation and conditional effects.

3.2. Conditional Correlated Bernoulli Model for Different Fingerprints. We then compared different fingerprint types using the conditional correlated Bernoulli model including MACCS⁹ (166 bits), ECFP4¹⁵ (variable length), TGD¹⁶ (420 bits), and TGT¹⁶ (1704 bits). MACCS represents structural keys describing substructures or patterns consisting of 1–10 non-hydrogen atoms. ECFP4 is an extended connectivity fingerprint with a maximum bond diameter of four that monitors layered atom environments in molecules. It generates molecule-specific variable feature numbers and is implemented in Pipeline Pilot.¹⁷ The typed-graph distances (TGD) fingerprint generates atom pairs where each atom is assigned one of seven atom types (acid, base, hydrogen bond donor, hydrogen bond acceptor, both hydrogen bond donor and acceptor, hydrophobic, other) and interatomic distances are divided into 15 different bond distances. The typed-graph triangles (TGT) fingerprint is a 2D

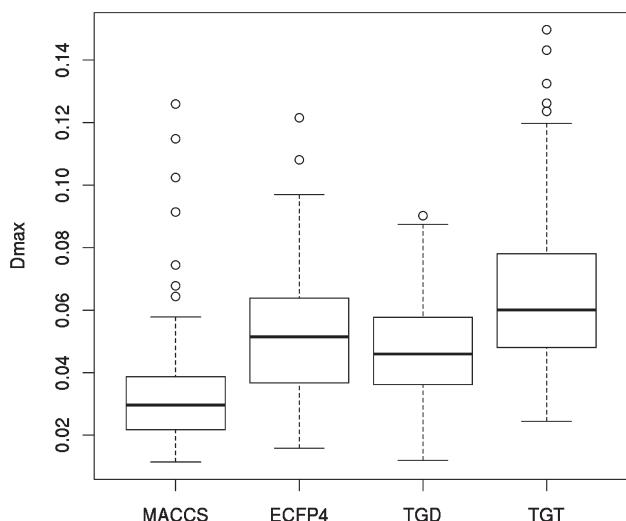


Figure 5. Conditional distributions for different fingerprints. The quality of 100 modeled conditional distributions based on randomly selected database compounds is evaluated for different fingerprints. The goodness-of-fit of a modeled distribution is evaluated by calculating the Kolmogorov–Smirnov statistic (D_{\max}). Boxplots are shown for MACCS, ECFP4, TGD, and TGT.

three-point pharmacophore fingerprint that generates atom triplets. Each atom is assigned one of four atom types (hydrogen bond donor or base, hydrogen bond acceptor or acid, both hydrogen bond acceptor and donor, or hydrophobic) and interatomic distances are divided into six different bond distances.

For the different fingerprints, empirical and predicted p -values were compared under the assumption of the conditional correlated Bernoulli model for reference and database compounds randomly chosen from ZINC (in analogy to the calculations described above for the comparison of model variants). The results are shown in Figure 4. Although these four fingerprints are of different design and size and exhibit different feature correlation effects, the curves closely follow the diagonal and also display only little variation for different reference compounds. The closest match to the diagonal was observed for MACCS (Figure 4a) and largest deviations were observed for TGT (Figure 4d), which by design introduces systematic feature correlation through the calculation of predefined (atom triplet) pharmacophore patterns. Hence, this fingerprint type was least suitable for distribution modeling due to the presence of systematic feature correlation effects, as anticipated. However, the conditional correlated Bernoulli model predicted the T_c distributions and p -values overall very well for different types of fingerprints.

To quantitatively assess the quality of the modeled distribution, the Kolmogorov–Smirnov (K–S) statistic¹⁸ was applied. The K–S statistic D_{\max} represents the largest deviation of the empirical distribution from the modeled distribution. Formally, D_{\max} is defined as follows:

$$D_{\max} = \max_x |F_{\text{BM}}(x) - F_{\text{emp}}(x)|$$

where F_{BM} is the modeled CDF and F_{emp} is the empirical CDF.

Visually, D_{\max} can be interpreted as the largest deviation of the curves from the diagonal in Figures 3 and 4. For each of the four fingerprints MACCS, ECFP4, TGD, and TGT, D_{\max} values were calculated for 100 modeled conditional distributions. As above,

each of the modeled distributions was derived from randomly sampled reference compounds and the empirical distributions were determined for 10 000 randomly sampled database compounds. Figure 5 reports the results. The box plots show that MACCS yielded overall the best models, consistent with visual analysis of the curves in Figure 4. The median K–S statistic was ~ 0.03 , which means that for 50% of the modeled distributions, the least accurate prediction of a p -value was within 3% of the empirical p -value. ECFP4 and TGD displayed similar quality in their models with $\sim 75\%$ of the models yielding a D_{\max} of below 0.06. TGT performed slightly worse (again consistent with visual inspection of the curves in Figure 4), with a median of ~ 0.08 and a number of models with D_{\max} larger than 0.10. This can be attributed to the fact that the cardinalities of the intersections and the union were not as accurately approximated by normal distributions as for the other fingerprints, due to the presence of systematic correlation effects in TGT (vide supra).

4. PREDICTION OF FINGERPRINT SEARCH PERFORMANCE

We then predicted the search performance of four fingerprints (MACCS, ECFP4, TGD, and TGT) on nine previously reported activity classes of varying intraclass structural diversity.¹⁹ ROC curves comparing the predictions according to section 2.5. with empirical fingerprint search results are presented in Figure 6. The activity classes contained 22–27 compounds (and are designated in the legend of Figure 6). For the benchmark calculations, each active compound was once used as a reference for fingerprint searching and the remaining active molecules were added to a background database of 1 000 000 randomly selected ZINC compounds. For calculation of the empirical ROC curve, all observed compound ranks were combined, in analogy to the prediction. As expected and shown in Figure 6, fingerprint search performance in benchmark calculations generally decreased with increasing structural diversity of the activity classes. In addition, in some cases, the search performance of individual fingerprints considerably varied. Thus, these search calculations covered a wide performance range and yielded compound class-specific fingerprint differences. However, the comparison in Figure 6 revealed very good to excellent agreement between our predictions and the empirical search results in all cases, without an exception.

It should be noted that prediction of ROC curves is not essential in typical benchmark investigations, where ROC curves are usually determined on the basis of reference and test sets. However, in these cases, the prediction using the conditional correlated model is computationally much more efficient. For example, for a nonparallelized implementation of these methods in the Scala programming language²⁰ on a standard desktop computer with a 2.67 GHz Intel Xeon quad-core processor (running Linux), the prediction of ROC curves using the Bernoulli model is approximately 50–650 times faster than the standard determination through data collection, depending on the particular fingerprint. However, computational efficiency is not the key aspect in this context. Rather, we emphasize that fingerprint performance can actually be predicted using the conditional correlated model in prospective situations when only a compound reference set is considered once a model for the score distribution has been established. For a given set of known active compounds, a particular fingerprint, and screening database, it can then be calculated following the protocol described above at which positions active compounds would be placed in

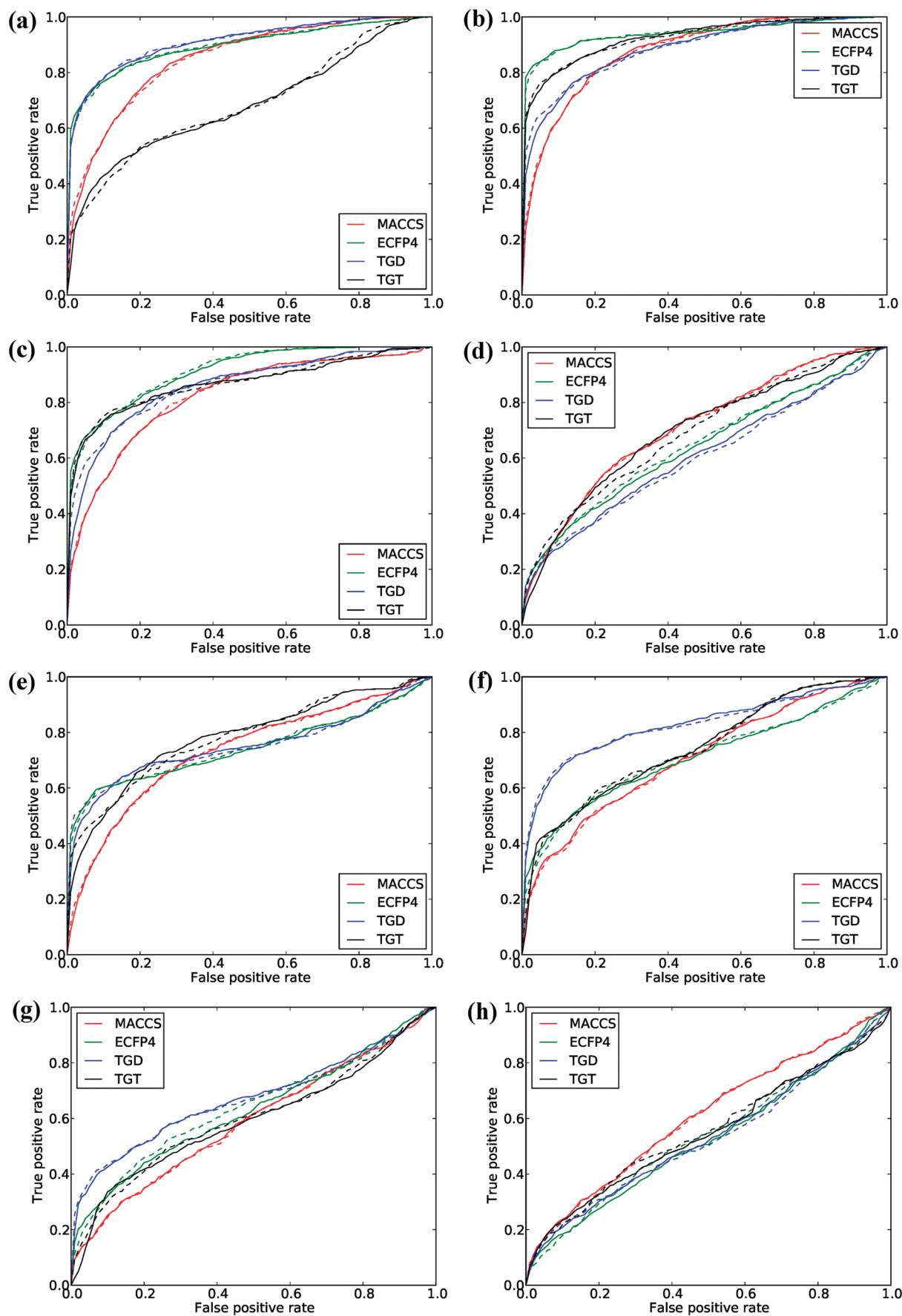


Figure 6. Continued

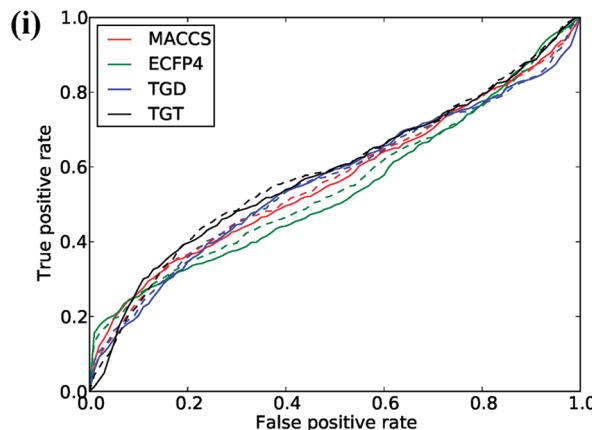


Figure 6. Predicted and empirical ROC curves. In each graph, the ROC curves for four fingerprints are shown as predicted by the conditional correlated Bernoulli model (solid lines). The dashed lines show the empirically determined ROC curves for a benchmark screen. Results for nine different activity classes are shown. From (a) to (i), these classes are arranged in the order of increasing intraclass structural diversity.¹⁹ (a) Angiotensin-II antagonists (ANG), (b) Renin inhibitors (REN), (c) HIV protease inhibitors (HIV), (d) thrombin inhibitors (THR), (e) IL-1 β -converting enzyme inhibitors (IL1), (f) endothelin antagonists (ETA), (g) aldose reductase inhibitors (ARI), (h) leukotriene synthesis inhibitors (LSI), and (i) cyclooxygenase-2 inhibitors (COX).

database rankings. If these rank positions are high, the fingerprint search would have a high likelihood to detect active compounds, provided the screening database contains other active compounds that are chemically similar to the reference set, which is unknown in practical applications, different from benchmark settings. Thus, by carrying out these calculations it can be estimated at which similarity values and rank positions active compounds would be detected, which makes it possible to compare different fingerprints and select those for prospective applications that place active compounds highest in the database ranking.

5. CONCLUDING REMARKS

With the conditional correlated Bernoulli model, we have introduced a method for modeling similarity value distributions of fingerprint search calculations that takes feature correlation and conditionality of value distributions on reference compounds into account. The methodology has been presented in detail and its theoretical foundations have been rationalized. Evaluation of the model has revealed that both reference conditionality and feature correlation must be considered for high-quality modeling of similarity value distributions. Furthermore, as we have discussed, the development of approaches to prospectively estimate the potential success of similarity search calculations (or, in general, virtual screening) is still a wide open field and only very few concepts have been introduced to estimate virtual screening performance in practical applications. However, the conditional correlated Bernoulli model is also readily applicable to predict fingerprint search performance. We have compared rank predictions made by our model with empirical search results of different fingerprints on compound activity classes with varying structural diversity. On the basis of this comparison, the predictions have been consistently accurate. We conclude that the conditional correlated Bernoulli model introduced herein represents a promising approach for modeling of conditional similarity value distributions and for predicting similarity search performance in the context of practical virtual screening applications.

AUTHOR INFORMATION

Corresponding Author

*Tel: +49-228-2699-306. Fax: +49-228-2699-341. E-mail: bajorath@bit.uni-bonn.de.

REFERENCES

- (1) Geppert, H.; Vogt, M.; Bajorath, J. Current Trends in Ligand-Based Virtual Screening: Molecular Representations, Data Mining Methods, New Application Areas, and Performance Evaluation. *J. Chem. Inf. Model.* **2010**, *50*, 205–216.
- (2) Stumpfe, D.; Bajorath, J. Similarity Searching. *Wiley Interdiscip. Rev.: Comput. Mol. Sci.* **2011**, *1*, 260–282.
- (3) Cover, T. M.; Thomas, J. A. *Elements of Information Theory*; Wiley: New York, NY, 1991.
- (4) Vogt, M.; Bajorath, J. Introduction of an Information-Theoretic Method to Predict Recovery Rates of Active Compounds for Bayesian in Silico Screening: Theory and Screening Trials. *J. Chem. Inf. Model.* **2007**, *47*, 337–341.
- (5) Vogt, M.; Bajorath, J. Introduction of a Generally Applicable Method to Estimate Retrieval of Active Molecules for Similarity Searching using Fingerprints. *ChemMedChem* **2007**, *2*, 1311–1320.
- (6) Muchmore, S. W.; Debe, D. A.; Metz, J. T.; Brown, S. P.; Martin, Y. C.; Hajduk, P. J. Application of Belief Theory to Similarity Data Fusion for Use in Analog Searching and Lead Hopping. *J. Chem. Inf. Model.* **2008**, *48*, 941–948.
- (7) Swann, S. L.; Brown, S. P.; Muchmore, S. W.; Patel, H.; Martin, Y. C.; Hajduk, P. A Unified Probabilistic Framework for Ligand- and Structure-Based Virtual Screening. *J. Med. Chem.* **2011**, *54*, 1223–1233.
- (8) Baldi, P.; Nasr, R. When Is Chemical Similarity Significant? The Statistical Distribution of Chemical Similarity Scores and Its Extreme Values. *J. Chem. Inf. Model.* **2010**, *50*, 1205–1222.
- (9) MACCS Structural Keys; Accelrys: San Diego, CA, 2011.
- (10) Willett, P.; Barnard, J. M.; Downs, G. M. Chemical Similarity Searching. *J. Chem. Inf. Comput. Sci.* **1998**, *38*, 983–996.
- (11) Hinkley, D. V. On the Ratio of Two Correlated Normal Random Variables. *Biometrika* **1969**, *56*, 635–639.
- (12) George Marsaglia, G. Ratios of Normal Variables and Ratios of Sums of Uniform Variables. *J. Am. Stat. Assoc.* **1965**, *60*, 193–204.
- (13) Bradley, A. P. The Use of the Area under the ROC Curve for the Evaluation of Machine Learning Algorithms. *Pattern Recognit.* **1997**, *30*, 1145–1159.

- (14) Irwin, J. J.; Shoichet, B. K. ZINC—A Free Database of Commercially Available Compounds for Virtual Screening. *J. Chem. Inf. Model.* **2005**, *45*, 177–182.
- (15) Rogers, D.; Hahn, M. Extended-Connectivity Fingerprints. *J. Chem. Inf. Model.* **2010**, *50*, 742–754.
- (16) TGD, TGT. In *MOE (Molecular Operating Environment)*; Chemical Computing Group Inc.: Montreal, Quebec, Canada, 2009.
- (17) *Pipeline Pilot*, student ed.; Accelrys, Inc.: San Diego, CA, 2009.
- (18) Birnbaum, Z. W.; Tingey, F. H. One-Sided Confidence Contours for Probability Distribution Functions. *Ann. Math. Stat.* **1951**, *22*, 592–596.
- (19) Tovar, A.; Eckert, H.; Bajorath, J. Comparison of 2D Fingerprint Methods for Multiple-Template Similarity Searching on Compound Activity Classes of Increasing Structural Diversity. *ChemMedChem* **2007**, *2*, 208–217.
- (20) The Scala Programming Language. <http://www.scala-lang.org/> (accessed Aug 1, 2011).