

Kinome-wide Activity Modeling from Diverse Public High-Quality Data Sets

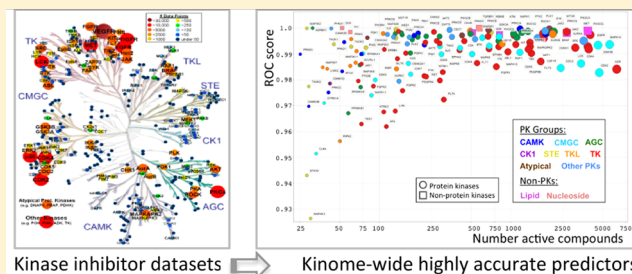
Stephan C. Schürer^{*,†} and Steven M. Muskal[‡]

[†]Department of Molecular and Cellular Pharmacology, Miller School of Medicine and Center for Computational Science, University of Miami, Miami, Florida 33136, United States

[‡]Eidogen-Sertanty, Inc., 3460 Marron Road No. 103-475, Oceanside, California 92056, United States

S Supporting Information

ABSTRACT: Large corpora of kinase small molecule inhibitor data are accessible to public sector research from thousands of journal article and patent publications. These data have been generated employing a wide variety of assay methodologies and experimental procedures by numerous laboratories. Here we ask the question how applicable these heterogeneous data sets are to predict kinase activities and which characteristics of the data sets contribute to their utility. We accessed almost 500 000 molecules from the Kinase Knowledge Base (KKB) and after rigorous aggregation and standardization generated over 180 distinct data sets covering all major groups of the human kinome. To assess the value of the data sets, we generated hundreds of classification and regression models. Their rigorous cross-validation and characterization demonstrated highly predictive classification and quantitative models for the majority of kinase targets if a minimum required number of active compounds or structure–activity data points were available. We then applied the best classifiers to compounds most recently profiled in the NIH Library of Integrated Network-based Cellular Signatures (LINCS) program and found good agreement of profiling results with predicted activities. Our results indicate that, although heterogeneous in nature, the publically accessible data sets are exceedingly valuable and well suited to develop highly accurate predictors for practical Kinome-wide virtual screening applications and to complement experimental kinase profiling.



INTRODUCTION

Over 500 human protein kinases¹ (<http://www.kinase.com/>) as well as many nonprotein kinases are involved in virtually every signal-transduction process, which are controlled by phosphorylation cascades. Because of their ubiquity, kinases are one of the most intensely pursued classes of drug targets. Numerous distinct kinase targets and over 300 kinase inhibitors are in clinical development,^{2,3} and 15 protein kinase inhibitor drugs have been approved so far (compiled from different sources). While currently most of the clinical kinases drug targets are being investigated for treatment of cancer, a growing number of other disorders including immunological, neurological, metabolic, and infectious diseases have been associated with dysregulation of protein phosphorylation.⁴ This suggests that the number of kinases as potential drug targets is substantial and has generated huge interest in the development of small-molecule inhibitors, most of them targeting the ATP binding site of the kinase catalytic domain.² The ATP binding region is highly conserved among protein kinases, which has important consequences for the drug discovery process. Achieving selectivity of a small molecule inhibitor against kinase off-targets to avoid adverse reactions has generally been considered challenging. On the other hand, the clinical efficacy of many of the current kinase inhibitor oncology drugs is related to their polypharmacology—their ability to inhibit

multiple kinases at the same time.⁵ It is now well-understood that the development of novel efficacious and safe compounds requires a finely tuned balance of polypharmacology and selectivity.⁶ Despite the high degree of conservation in the ATP binding site, reasonably selective inhibitors with favorable pharmacological properties can be developed.⁷ This is helped by the increasing understanding of structural and functional relationships across the kinome.^{8,9} Today, it is quite common to profile inhibitors against an extensive set of kinase targets¹⁰ at an early stage of development. Kinase profiling technologies have generated valuable data sets and provided insights into the determinants of selectivity and promiscuity of clinical inhibitors.^{11–13} One such public effort currently takes place under the NIH Libraries of Network-based Cellular Signatures (LINCS) program.¹⁴ The LINCS program develops a library of molecular signatures based on gene expression and other cellular changes in response to perturbing agents across a variety of cell types using various high-throughput screening approaches. LINCS profiled compounds include many kinase inhibitors, because of their translational potential. LINCS kinase profiling results are currently available from the HMS

Received: August 24, 2012

Published: December 21, 2012

LINCS DataBase¹⁵ and also via the LINCS Information Framework (LIFE).¹⁶

High-throughput technologies and conventional screening over two decades of small molecule kinase inhibitor research have generated a huge amount of data available in peer-reviewed publications and patents. In addition to numerous individual virtual screening approaches, the availability of large data sets of small molecule kinase inhibitors—in particular in the pharmaceutical industry—has spurred efforts to utilize the available data sets to understand and model kinase polypharmacology for example by analysis of inhibitor cross reactivity^{17–19} and analysis of kinase chemotype selectivity patterns.²⁰

Here, we focus on the data that are available in the public domain. The reliability of heterogeneous data sets generated under different screening conditions and using different assay methods and technologies are sometimes questioned. We want to understand how useful these types of data sets really are. To do this, we evaluate results from different machine learning techniques applied to the data. Earlier examples among numerous machine learning approaches to classify kinase inhibitors include neural networks trained using BCUTS descriptors,²¹ Naïve Bayesian modeling using extended 2D topological fingerprints and basic molecular descriptors,²² and a survey of different machine methods using fragment-based Ghose-Crippen descriptors.²³ In an effort toward virtual polypharmacology, fragments and fragment counts have successfully been used in prospectively exploring kinase inhibitor activity space.^{24,25} Here, we generated and characterized hundreds of kinase classification and regression models based on a large number of data sets extracted from the Kinase Knowledge Base (KKB),²⁶ which span the entire human kinome. The data sets were curated from several thousand peer-reviewed journal and patent publications from numerous laboratories comprising various assay technologies, assay designs, and procedures. We investigated the applicability of these data sets to generate predictive models, which modeling technique(s) were best suited and applicable for virtual screening, and which characteristics of data sets led to good predictors. We applied the best predictors to compounds profiled in the LINCS program and found that these profiling results were in good agreement with predicted actives.

METHODS

Kinase Data Sets. The Q4 2009 release of the KKB incorporated >430 000 bioactivity data points from over 20 000 assay experiments curated from more than 1800 journal articles and more than 4400 patent publications. It covers over 500 000 unique molecules. The KKB includes various metadata annotations including standardized target and assay format. From the available assay experiments only data from biochemical assays (enzymatic assays, performed with purified protein) of human species were chosen. Any mutant kinase targets were excluded. Only high-quality concentration–response end points (e.g., IC_{50} , K_i) were kept. Because in the KKB, chemical structures are stored as published, we standardized all structures using an in-house Pipeline Pilot protocol.²⁷ Salts/addends and duplicate fragments were removed (using an in-house salt library) so that each structure consisted of only one fragment. Stereochemistry and charges were standardized, the structures were then ionized at pH = 7.4 and tautomers were canonicalized. For the purpose of this modeling study in which we use extended connectivity

fingerprint of length four (ECFP4) descriptors (see below), stereochemistry and E/Z geometric configurations were also removed. All data points were first transformed into p-values (i.e., $pIC_{50} = -\log_{10}[IC_{50}]$ in molar concentration) and then aggregated, first within each experiment and then across experiments by unique structures and kinase targets. Kinase protein targets were identified by unique standardized Entrez Gene symbols, which are annotated in the KKB. The kinase gene symbols were mapped to Uniprot accessions. Although detailed experimental descriptions are available in the KKB, we did not filter any assay technologies or experimental conditions. Identical target-structure data points were aggregated using the median for exact data points (to minimize the effect of outliers introduced by different assay methods and experiments) and in case of qualified data (greater than, less than, range) the most conservative (inclusive) ranges were kept. This aggressive data aggregation procedure resulted in 233 667 unique (structure–target) data points for 126 114 unique structures covering 411 unique kinase targets. We also standardized the structures of all KKB molecules resulting in a total of 489 373 unique chemical compounds. Access to these data sets along with the entire KKB database for academic research groups can be obtained through the portal www.kinasedb.com. For the naïve Bayesian classification models, active compounds were defined as p-transformed activity concentration of greater or equal to 6 (i.e., $IC_{50} \leq 1 \mu M$). Using this definition of “active”, 189 kinase targets have at least 10 active compounds (see Supporting Information Table S1; not including TBK1, which includes 169 actives and no inactives). We built naïve Bayesian classification models treating the data sets in two different ways: in one case, the classifiers are trained employing only compounds that are explicitly defined as active or inactive (KA-KI; known active-known inactive). In the other case, we use all unique KKB compounds as decoys and presume as inactive all compounds that are not annotated as actives for any specific kinase (KA-PI; known active-presumed inactive). For the regression models (see below) only data points with exact activity values (no <, >, or range data) were kept. They included data sets for 168 kinases with at least 20 structure–activity data points (Supporting Information Table S2, which also includes data set statistics). To evaluate the activity range of kinase inhibitors by structural series and the coverage of structural series across different kinase targets, we clustered all compounds into 336 clusters (twice the number of kinase data sets). We used the partitioning algorithm implemented in the Pipeline Pilot 8.0 (Accelrys) modeling collection with the Tanimoto distance function on ECFP4 fingerprints (see below); cluster centers were selected by maximum dissimilarity. After generating the clusters, we visualized the pIC_{50} activity ranges for each kinase data set for each cluster (Supporting Information Table S3).

KINOMEScan kinase profiling results from the LINCS project were downloaded from the HMS LINCS DataBase.¹⁵ In total 25 064 data points were obtained with 60 unique compounds (by HMSL_ID), 43 of them with defined/known chemical structure, and 486 different targets (including mutations for several kinases); not all compounds were tested against all targets. Kinase activity was screened at 10 μM compound concentration. Reported activities were transformed into percent inhibition. The kinase targets were manually annotated with their corresponding Uniprot accessions based on their symbol/description. The targets were mapped to the KKB data sets (and models) based on their Uniprot accessions. A total of 4796 kinase compounds pairs were mapped for

compounds with known structures after removing mutant kinase targets from the KINOMEScan data sets.

Laplacien-Corrected Naïve Bayesian Classifiers. Naïve Bayes is a statistical classification method based on conditional probabilities: in this context, the probability of a compound being active given the presence of structural features computed from the frequencies of occurrence in a training set of active and inactive samples. “Naïve” refers to the assumption that features are independent and that the overall probability therefore can be computed by multiplying probabilities of the individual events. The Laplacien-corrected estimator accounts for the different sampling frequencies of different features assuming that the vast majority of features have no relation with activity. The Laplacien correction stabilizes the estimator: as the number of samples containing a feature approaches zero, the features probability contribution converges to the baseline probability. The final estimate for a particular sample is computed as the sum of the logarithm of the relative (corrected) individual feature weights.²² This classifier has several desirable criteria. It scales linearly with the number of molecules and is therefore applicable to large data sets. Bayesian classification is suitable to model in high-dimensional spaces (large number of descriptors do not cause overfitting). It is therefore appropriate for developing models from structurally dissimilar molecules and to incorporate multiple activity classes into a single model (for example different sites/binding modes or different mechanism of action). The Laplacien-corrected naïve Bayes classification method is also reasonably resistant to noise such as false positives or false negatives.²⁸ The Pipeline Pilot modeling collection includes an implementation of this classification learner, which was employed here.

Here we built Laplacien-corrected naïve Bayesian classification models based on two different methods of handling the kinase data sets as described above. Classifiers are trained from known actives and known inactives (see data sets for 188 kinases provided in Supporting Information Table S1). KA-PI classifiers were built using one data file of all unique (489 373) kinase molecules in which the activity categories are defined by an array containing the respective kinase symbol(s) for which each compound qualifies as active. The classifiers were then built by identifying for each individual kinase the active compounds and treating the remaining compounds as inactive (decoy). Both types of classification models were first built using all data sets and characterized by the Pipeline Pilot internal leave-one-out cross-validation to estimate the area under the receiver operating characteristic (ROC) curve (ROC score) and enrichment results (for 1, 5, 10, 25, 50, 75, and 90%). The ROC curve is the true positive rate (TPR) over the false positive rate (FPR). TPR is equal to sensitivity (S) of the model and defined as the number of true positives (TP) divided by the number of actives (N_{act}). Specificity (SP) is the number of true negatives (TN) divided by the number of inactives (N_{inact}). FPR is $1 - \text{SP}$, which is equal to the number of false positives (FP) divided by N_{inact} . The enrichment factor (EF) at any given percentage of compounds tested is defined as $\text{EF} = [\text{TP}/(\text{TP} + \text{FP})]/[N_{\text{act}}/N]$ where N is the overall number of samples; it is a measure of how well the predictor recovers active compounds relative to random retrieval. To validate these modeling approaches more rigorously, we randomly split the data sets into 75/25 training/test sets and built the models using the training sets and generated ROC scores and enrichment results (for 0.1, 0.5, 1, 3, 5, 10, and 20%) using the test sets. This randomized train/test evaluation procedure

was repeated 10 times, and the results were averaged over the repetitions. Supporting Information Tables S1 and S4 report the modeling results for all 188 and 189 kinase data sets corresponding to the KA-KI and KA-PI approach, respectively. Tables S1 and S4 also include the average numbers of training and test sets for each kinase data set and the maximum achievable (perfect) enrichment factors (EF_{max}), which is the maximum number of actives among the percentage of tested compounds divided by the fraction of actives in the entire data set. To evaluate enrichment, we report the ratio of $\text{EF}/\text{EF}_{\text{max}}$. This normalized enrichment factor is a useful measure, because EF values are not directly comparable across different data sets. This is because maximum possible EF by definition is limited to the ratio of total to active compounds in a data set. EF is also limited to the reciprocal of the fraction of compounds tested. For example if 0.1% of all compounds are tested, the highest possible EF is 1000 if the fraction of actives is less or equal to 0.1% of total compounds (and if all retrieved compounds are true positives). Table S1 includes some empty fields for EF at lower percentages, because enrichment factors for a given percentage of screened compounds can only be obtained if the number of tested compounds is at least one. For the KA-PI classification models using the entire set of kinase compounds as decoy, we further characterized the classifiers by the same cross-validation procedure employing training/test set ratios of 50/50 as well as 25/75. These results are reported in Supporting Information Table S5. We also built KA-PI classifiers from the same chemical structures but after randomizing the kinase activities while maintaining the number of actives for each kinase. Leave-one-out cross-validation resulted in ROC scores of close to 0.5 and enrichment factors of approximately the ratio of actives to overall number of samples, which corresponded to random classification as expected.

To test domain applicability, 53 kinase data sets that each have at least 500 active compounds were selected and each clustered into 10 series using the partitioning algorithm implemented in the Pipeline Pilot 8.0 (Accelrys) modeling collection with the Tanimoto distance function on ECFP4 fingerprints (the same as for the kinase models, see below); cluster centers were selected by maximum dissimilarity. For each kinase data set, 10 models were built, each using a different series as the (active) test set and the remaining 9 clusters combined as the (active) training set. Inactive test and training compounds were selected randomly for each model in the same ratio of active training to test compounds. This way 530 models were built. For each active test compound the Tanimoto similarity to the closest active training compound was calculated. ROC scores were computed for each model based on predictions of the test set to evaluate the performance of the models based on similarity of the test to the training set (details are provided in Supporting Information Table S6). In addition, for each active test set compound the predicted activity was recorded to evaluate true positive rate as a function of similarity to the closest training compound. In total 98 731 predictions were made across all 530 models.

To compare predicted kinase activity to the KINOMEScan profiling results, we used the EstPGood output of the KA-PI classifiers, the estimated probability that the sample is in the active category based on an assumed normal distribution within the active and inactive categories.

Regression Models. Both k nearest neighbor (kNN) and partial least square (PLS) regression was performed using the

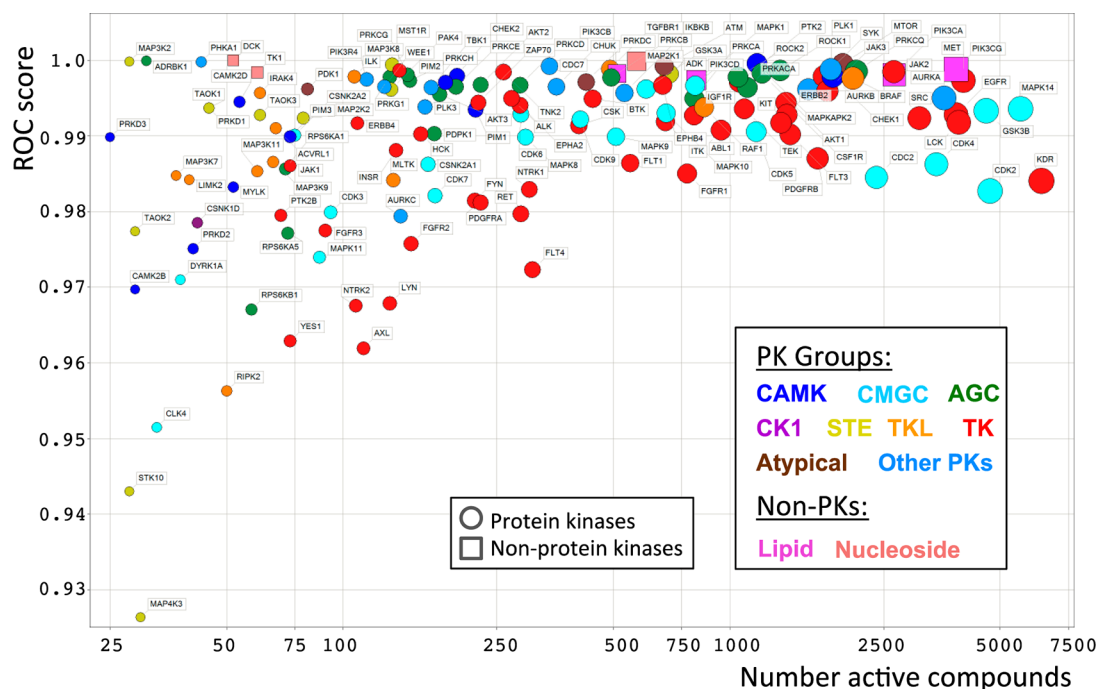


Figure 1. Characterization of 141 kinase KA-PI protein and nonprotein kinase classifiers, including all major protein kinase groups. ROC scores are shown as a function of active samples. Shape by protein vs nonprotein kinase, color-coded by kinase group, scaled by number of active data points, and annotated by HUGO kinase gene symbol.

implementations of these machine learning procedures in the Pipeline Pilot modeling collection. PLS is a statistical method to build a (linear) model of the response variable as a function of (uncorrelated) linear combinations of the input descriptors (components); it is arguably one of the most traditional, yet widely used, approaches to qualitative structure–activity relationship (QSAR) study.²⁹ In our PLS models, we restrict the number of components to 20. The kNN method is among the simplest in machine learning algorithms. In kNN regression, the response variable is computed as a weighted average of the value of the k nearest neighbors. Here we use 20 nearest neighbors and a dynamic smoothing factor (Gaussian weighting) of 0.5. For both PLS and kNN, we performed a full 10-fold cross-validation and computed R^2 for the training data, q^2 for cross-validation, and RMSE for training and cross-validation. We repeated the 10-fold full cross-validation procedure 10 times with random split of the 10 subsets. The averaged results are reported in Supporting Information Table S2.

Structural Descriptors. Because of their strong performance in previous studies,^{30–32} we employed extended connectivity fingerprints (ECFPs). Specifically we used atom type ECFPs of length four (ECFP4) implemented in Pipeline Pilot.³³ ECFPs are topological circular fingerprints characterizing each atom by its number of atomic connections, element type, charge and mass, and environment (in this case up to four neighbor atoms).

RESULTS AND DISCUSSION

KA-KI Classifiers. We first built and evaluated Laplacian-modified naïve Bayesian KA-KI classifiers. For this we selected all data sets from the preprocessed (aggregated) KKB with minimum of 10 active compounds where active is defined as p-transformed activity value of ≥ 6 (i.e., $IC_{50} < 1 \mu M$). A total of 188 kinase data sets qualify; the minimum data set size (actives

plus inactives) was 26. The classifiers were cross-validated by a leave-one-out analysis and 10 repetitions of randomized 75/25 training/test split (see Methods). For both cross-validation methods, ROC scores and enrichment factors show that these KA-KI classifiers perform well for most of the data sets (Supporting Information Table S1). Good results are obtained for data sets with more than 40 active compounds (minimum 61 compounds total), which corresponds to 130 data sets with a ROC score (leave-one-out) of 0.7 or larger; that is with the exception of one data set (MST1R) that has only two inactives and therefore cannot be meaningfully evaluated and was excluded from further analysis (leaving 129 data sets). Raising the minimum number of actives to 70 further improves the results to ROC scores of 0.84 or greater (111 data sets with at least 113 total samples). ROC scores based on leave-one-out and 75/25 (train/test) cross-validation procedures are well correlated for data sets with at least 40 actives ($R^2 = 0.73$) and increases further for data sets with 70 actives or more ($R^2 = 0.77$) as shown in Figure S1 (Supporting Information). Correlation of the ROC score of the two cross-validation procedures increases significantly for data sets with a ratio of total to active compounds of two or greater: for the data sets with at least 40 actives, R^2 is greater than 0.85, and for data sets with ≥ 70 actives, R^2 improves to ≥ 0.93 . However, this decreases the number of data sets to 57 and 47, respectively. In general the classifiers for the more balanced and larger data sets performed better.

For the 129 data sets with >40 actives, enrichment factors (EF) at 10% tested samples varied between ~ 1 and ~ 8 for both the leave-one-out and train/test cross-validation procedures. The maximum enrichment factor (EF max) depends on the fraction of actives in the data set. For all 129 data sets, the actual enrichment factors reach the maximum possible enrichment indicating well-performing classifiers. However, EF is of limited usefulness to characterize classifiers based on

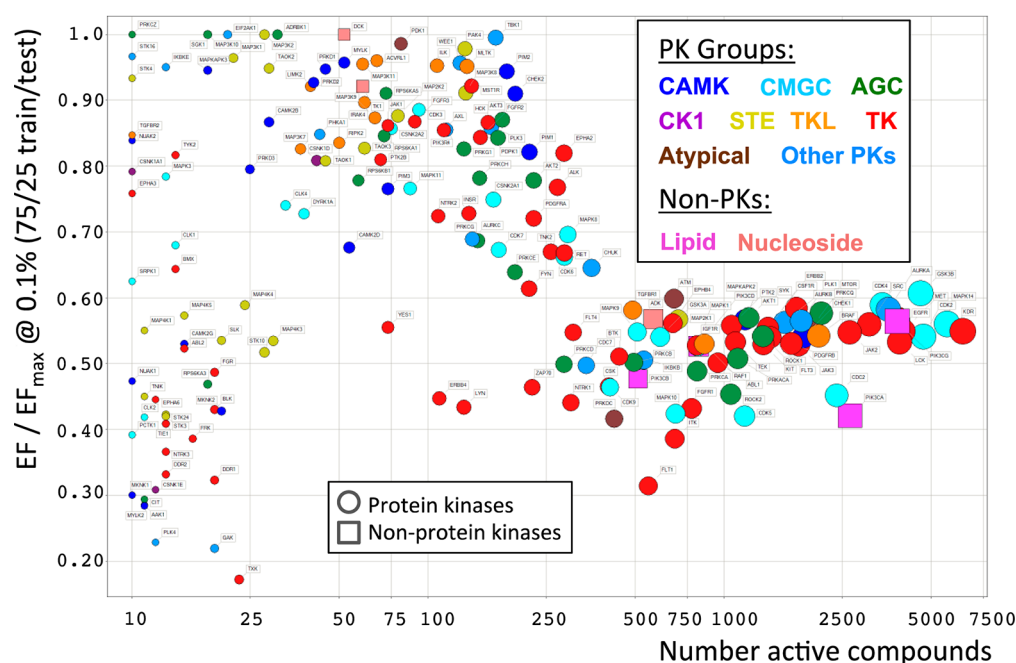


Figure 2. Normalized enrichment factors for 189 KA-PI kinase classifiers at 0.1% screened samples for 75/25 (train/test) cross-validation (10 repetitions averaged) as a function of active compounds in the data set. Shape by protein vs nonprotein kinase, color-coded by kinase group, scaled by number of active data points, and annotated by HUGO kinase gene symbol.

balanced data sets. The nature of the data sets reflects what is typically reported in the literature emphasizing active compounds and reporting only few (often structurally related) inactive compounds. In contrast, high-throughput screening (HTS) results include mostly inactive/negative results. Therefore, in practice, the KA-KI classifiers—although highly predictive—may be of limited utility. Highly miniaturized HTS allows a throughput of several hundred thousand to millions of compounds. Consequently, in order to recover a sizable fraction of hits by screening only a few hundred to thousand compounds, 100- to 1000-fold enrichment or even greater would be desirable.

KA-PI Classifiers. We investigated such a scenario by building KA-PI classifiers for the 189 kinase data sets with at least 10 actives as defined above. To build these classifiers, we employ all (489 373) KKB kinase molecules as decoys; i.e., we presume as inactive all compounds that were not specifically annotated as active for any given kinase. This resulted in 189 data sets, each consisting of 489 373 samples including 10–6388 actives. Their size and unbalanced nature corresponded quite well to real HTS data sets, for example those in PubChem.³⁴

Laplacian-modified naïve Bayesian classification is a suitable method to model our data sets of almost half a million unique kinase-related compounds (see Methods), which can bind at different sites, such as ATP competitive compounds, but also allosteric inhibitors, and with different binding modes, for example type I and type II inhibitors.³⁵ Although many of the presumed inactive compounds did not have specific activity annotations, they were all reported in patents and journal publications that focus on kinase inhibitors. They are thus closely related in terms of their biological focus and in many cases structurally.

Similar to the KA-KI classification models, the KA-PI classifiers were validated by leave-one-out and 10 repetitions of train/test cross-validation. We report ROC scores and

enrichment factors at various percentages of samples screened. In addition to the 75/25 split, we also evaluated the classifiers in a 50/50 and 25/75 training/test cross-validation (10 repetitions averaged, see Methods).

Figure 1 shows ROC score (leave-one-out) as a function of the number of active samples for all 141 kinases KA-PI classifiers with at least 25 actives. The 141 kinases cover all major groups of the human kinome and also several nonprotein kinases. All 189 models are shown in Supporting Information Figure S2. As a general trend, it can be seen that the quality of the classifiers increases with the number of active samples. In particular, data sets with more than about 50 actives resulted in much improved results compared to those with fewer actives. For classifiers with more than a few hundred actives, there appeared to be no further improvement in ROC score as the number of actives increases further. For the majority of models, the ROC scores are very high (>0.96). All details are provided in Supporting Information Table S4, which shows ROC scores and enrichment factors for the leave-one-out and 75/25 train/test cross-validation procedure for all 189 KA-PI kinase classifiers. Data set statistics are also shown. Supporting Information Figure S3 illustrates the relationship of ROC scores for leave-one-out vs 75/25 train/test cross-validation for KA-PI classifiers based on data sets with at least 10 active samples vs data sets with at least 50 actives. As with the KA-KI models, leave-one-out ROC score estimate was closely related to the ROC score obtained by 75/25 train/test validation, in particular for the data sets with a greater number of actives. Supporting Information Table S7 shows the ROC plots for the KI-PI classifiers based on a data set with >50 actives using a 75/25 train/test set. While Table S7 shows one ROC plot for each kinase/data set, it should be noted that ROC scores reported in Supporting Information Table S4 are averaged over 10 repetitions.

While ROC score is a good measure of the predictors overall performance, enrichment, in particular for low percentages of

screened compounds, is an important measure of the practical applicability of a predictor. Enrichment factors (EF) at very low percentages of screened samples for the KA-PI classifiers were very high for most kinases. Figure 2 illustrates the normalized enrichment factor, that is the ratio EF/EF_{\max} (see Methods), at 0.1% tested samples for all 189 kinase classifiers based on randomized 75/25 train/test cross-validation averaged over 10 repetitions. It can be seen that most classifiers are able to retrieve true actives very well if the number of actives in the data sets are greater than about 50. Enrichment for these classifiers is generally greater than 50% of EF_{\max} suggesting that the kinase classifiers presented here are practically applicable for virtual screening. Enrichment increased significantly for classifiers based on data sets of more than ~ 50 actives. This indicated a required minimum number of the active class of the training set to reliably retrieve true positives from the test set. This was also reflected in the ROC scores (compare Figure 1). Absolute EF values for 0.1, 0.5, and 1.0% of tested compounds are provided in Supporting Information Figure S4 and Table S4.

Figure 2 suggested the highest enrichments for data sets that have between 50 and 250 active molecules. The normalized enrichment factors were slightly lower and relatively stable around 0.5 for data sets with more than a few hundred actives. A possible reason for higher enrichment in data sets with smaller numbers of active compounds that are derived from only a few studies may lie in overrepresentation of scaffolds (analog bias; see below for domain applicability results). However, this is less likely for the larger data sets that have been extracted from a large number of articles and patents. It should also be emphasized here that the decoy (presumed inactive) compounds were all derived from the same kinase literature that are also the source of the active compounds and are therefore closely related in terms of their biological focus and in many cases also structurally. Normalized enrichment factors obtained by 75/25 train/test cross-validation at 0.1, 0.5, and 1.0% tested samples are shown in Supporting Information Figure S5. As the percentage of tested samples increases, a larger fraction of classifiers show a very high ratio EF/EF_{\max} (>0.8). Although expected, because the maximum possible enrichment decreases, the results also indicated that it is more difficult to retrieve a certain fraction of true positives in a smaller set of sampled compounds, compared to a larger, i.e. from 0.1% vs 0.5% or 1.0% tested compounds.

Following standard procedure to further validate the KA-PI classifiers with our data, we randomized the kinase activities (maintaining the number of actives for each kinase data set). Leave-one-out cross-validation resulted in ROC scores of close to 0.5 and EF of approximately the ratio of actives to overall number of samples. This corresponded to random classification as expected.

We also investigated how the ratio of training and test sets influenced the ROC and enrichment cross-validation results for the various kinase data sets. In addition to 75/25 (Supporting Information Table S4), we split the data into by 50/50 and 25/75 (Supporting Information Table S5). ROCs were slightly higher for larger compared to smaller training sets (Supporting Information Figure S6) indicating improved overall predictability, which was also consistent with the general trend of the ROC scores as a function of the number of active samples (Figure 1 and Supporting Information Figure S2). More specifically, as the ratio of training/test compounds decreased, the required number of active compounds in the data sets to

give very good ROC (>0.96) increased. This suggested a threshold of required active (training) compounds to develop very good predictors.

In contrast to increased overall predictivity (measured by ROC score) for larger (compared to smaller) train/test ratios, enrichment factors at very low percentages of tested compounds (0.1% and 0.5%) increased slightly with lower (compared to higher) ratios of train/test sets. Figure 3

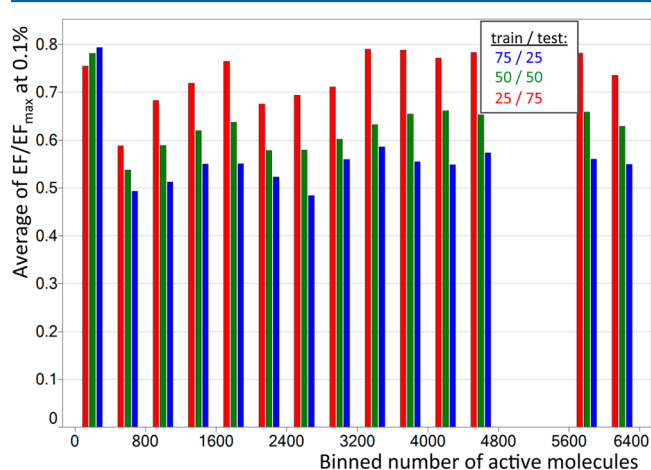


Figure 3. Average normalized enrichment (EF/EF_{\max}) at 0.1% tested samples for different train/test ratios binned by ranges of active compounds across 141 kinase data sets with at least 25 actives (EF averaged over 10 repetitions).

illustrates the average ratio of EF/EF_{\max} at 0.1% tested samples for the different train/test ratios across data sets for various ranges of numbers of active compounds. Other than for the lowest bin of active compounds (0–400), enrichment at 0.1% increases as the ratio of train/test decreases. Supporting Information Figures S7 and S8 show EF and EF_{\max} for each individual data set as a function of active compounds for the three train/test ratios at 0.1 and 0.5% tested samples, respectively. Supporting Information Figure S9 shows the normalized enrichment factors (EF/EF_{\max}) for each data set at 0.1% tested samples, illustrating the same trend. Although the ratio of active to inactive compounds is the same among the different train/test splits, these results suggest that it is “easier” for the classifier to select true actives from a test set with a larger (in absolute numbers) pool of active compounds. Because the trend holds for the data sets with the highest numbers of actives, it is likely not a trivial analog bias; although classifiers are based on structural features and therefore true positive test compounds are by definition structurally related to the active training compounds. More importantly, the results indicated again that once a certain number of active compounds are available, classifier performance does not increase significantly with additional data. This is consistent with the validation results described above based on ROC and EF.

To evaluate how similarities of test to training compounds affect model performance, we performed a simple domain applicability study (Figure 4, compare methods). Data sets with at least 500 active compounds were selected (53 kinase data sets representative of most of the human kinome), and the actives of each data set were clustered into 10 series. Ten models were generated each using one series as test set while using the remaining combined 9 clusters as training set. This way predictions are made for compounds that are structurally

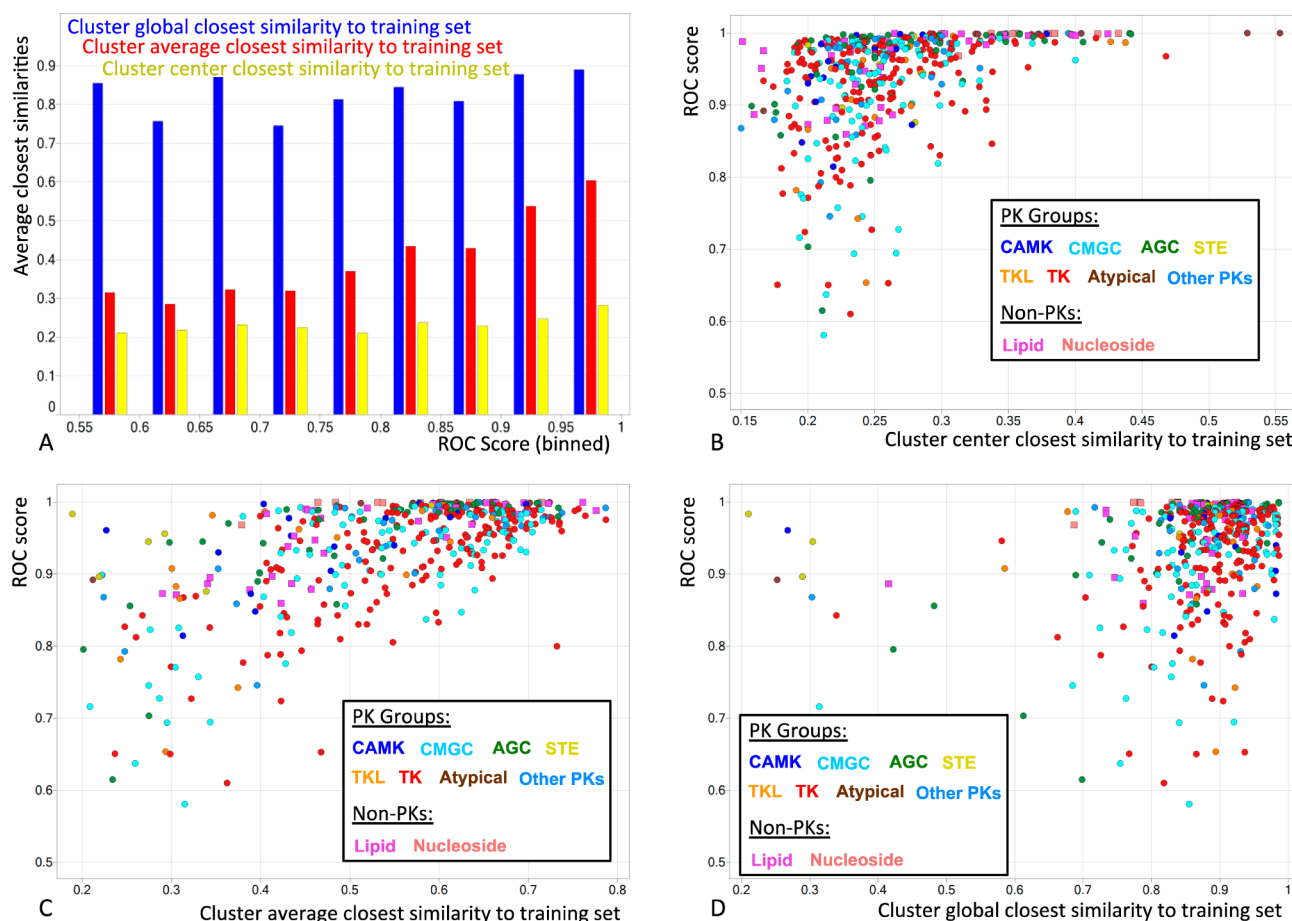


Figure 4. Model domain applicability. KA-PI model performance as measured by ROC score as a function of similarity of test to training sets. 443 models are shown for 53 kinases (representative of most of the human kinome). (A) Average closest similarities by binned ROC score. Shown are closest individual similarities (blue), cluster average closest similarities (red), and cluster center closest similarities (yellow). (B) ROC score by cluster center closest similarities. (C) ROC score by cluster average closest similarities. (D) ROC score by individual closest similarities (cluster global closest similarity).

dissimilar to training compounds. In this way, 530 models were generated and evaluated. Figure 4A shows a histogram of closest test-to-training set similarities by ROC score; results are shown for data sets with at least 10 active test compounds (443 models). Specifically we investigated ROC score as a function of the similarity of test set cluster center to the closest training compound (B), the average of the closest similarities of the test set compounds to the training set compounds (C), and the similarity of the closest training compound to any test compound (D). As can be seen, the best criteria of model applicability is the average closest similarity of the test compounds to the training compounds (panel C). The models show good predictivity ($\text{ROC} > 0.8$) for relatively low average closest similarities (>0.5).

In addition to ROC score for each model (cluster), we recorded the predictions of all active compounds over all 530 models (98 731 predictions total) and evaluated the hit rate as a function of similarity of the predicted test compound to the closest training compound (Supporting Information Figure S10). Figure S10A shows the true positive rate (TPR), and Figure S10B shows the distribution of true positives and actives as a function of closest similarity of the test compound to the training set. TPR drops off sharply as the similarity decreases below 0.5, and the TPR is greater than 0.8 for compounds with a Tanimoto similarity of >0.6 (ECFP4).

To illustrate the potential applicability of the KA-PI models for virtual screening, Figure 5 shows the enrichment results at 0.1% tested samples as the fraction of true positives obtained from the test set vs the fraction of actives in the entire data set. As expected the true positive rate increases with the ratio of active to total compounds until the latter reaches 0.1% after which the true positives identified from the test set remain relatively constant between 40 and 80%. In this plot, the enrichment factor achieved by each classifier at 0.1% is the quotient of the y and the x values. In practical terms, were the KA-PI classification models applied to prioritize 500 compounds from a library of about 500 000, one may expect to recover anywhere from 10 to 400 actives depending on the number of actual actives for a given kinase target. Thus these models appear practically applicable, assuming that our data sets reasonably well represent the kinase inhibitor chemical space.

The performance of the Laplacien-modified naïve Bayesian kinase KA-PI classifiers based on ROC scores and enrichment at very low percentages of tested compounds indicate that the data sets employed here are well-suited to build highly predictive kinase classification models. The results consistently suggest that the best classifiers are obtained if a minimum number of conservatively 50 active training compounds are available against a large decoy set. Performance of the models

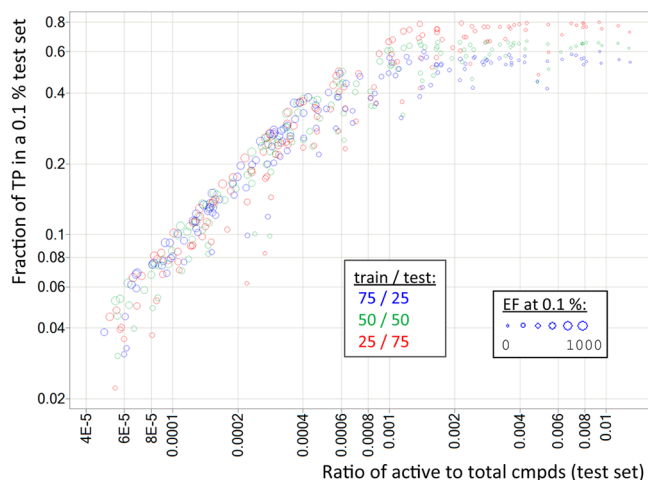


Figure 5. Fraction of true positives (TP) of the number of tested compounds in 0.1% of the test set as a function of the ratio of actives to total compounds for different ratios of training to test data. Also indicated is the enrichment factor as the size of the circles. The enrichment factor (at 0.1%) for each classifier is the quotient of y value and the x value.

can be further improved with additional active compounds of up to a few hundred but does not improve much beyond this. Given the low ratio of structure–activity data points to unique compounds in the KKB, it is likely that some of the presumed inactive (decoy) compounds are in fact active against one or several kinases. It is therefore feasible that the kinase KA-PI classifiers can be further improved as more kinase profiling data sets become available.

Regression Models. In addition to the naïve Bayes binary classifiers we also wanted to evaluate how suitable the data sets are for the development of quantitative predictors. We chose partial least square (PLS) and k nearest neighbor (kNN) regression as two fairly different learning methods to quantitatively predict a continuous property. PLS and kNN regression are considered most applicable for data sets of congeneric (structurally similar) molecules. The methods can be sensitive to outliers and require high-quality data. Here, 168 kinase data sets with at least 20 exact molecule–activity data points were extracted from the standardized and aggregated KKB data sets (see Methods). PLS and kNN QSAR models were built as described and evaluated by 10-fold cross-validation, which was further repeated 10 times randomly partitioning the data. All cross-validation results (including, q^2 , and RMSE) and data set statistics are summarized in Supporting Information Table S2. Figure 6 illustrates q^2 values for the kNN vs PLS models along with the number of structure–data points for each kinase. The kNN method in general outperformed PLS, and the performance (measured by q^2) of both methods correlated reasonably well. From Figure 6, one can also conclude a trend in which the predictive quality for both PLS and kNN models generally improved with the size of the data sets. Although this trend was not as strict as what we observed for the naïve Bayes classifiers, a cutoff for kNN $q^2 \geq 0.4$ and PLS $q^2 \geq 0.25$ leaves 91 kinase data sets—all except three having greater or equal to 50 structure–data points (also compare Figure 8). For kinases with 500 or more data points, all but one model have kNN q^2 values of >0.5 .

Another important characteristic influencing the quality of both PLS and kNN was the activity range of the data sets.

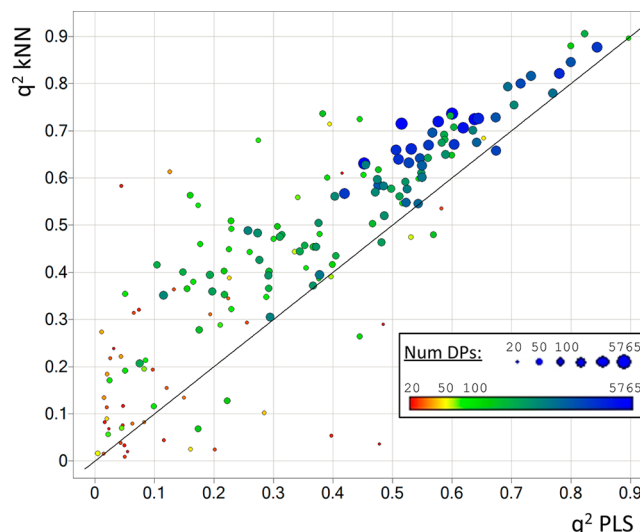


Figure 6. Quantitative regression models developed from 168 kinase data sets. q^2 values of kNN vs PLS regression models and the number of kinase activity data points indicated by the circle size and color (see text).

Figure 7 illustrates the average q^2 for both kNN and PLS as a function of the p -transformed activity range, which corresponds

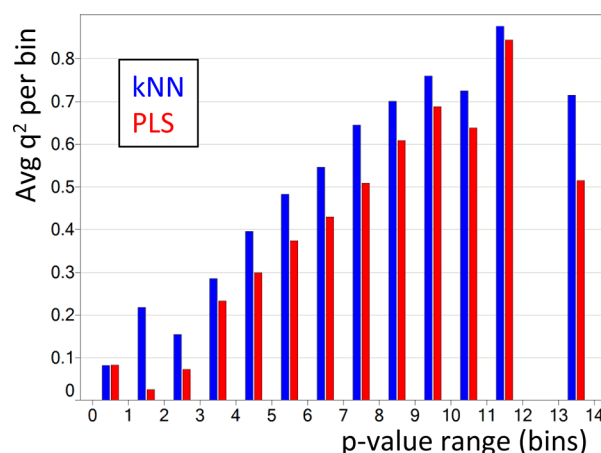


Figure 7. Average q^2 for kNN and PLS regression cross-validation results as a function of the p -value range (defined as $p\text{-value}_{\max} - p\text{-value}_{\min}$) of all 168 kinase data sets.

to the orders of magnitude between the most and least active compound. PLS and kNN q^2 values continuously increased from 0.1 to >0.8 as the activity range increased from 1 to 12 orders of magnitude. To evaluate the activity range across structural series and to see how structural series distribute across different kinase data sets, we clustered compounds into 336 clusters and calculated activity ranges for each kinase for each cluster (see Methods and Supporting Information Table S3). Supporting Information Figure S11 shows the total number of kinases and activity range of each cluster. These results show that many structural series span a wide activity range and many kinases suggesting that a scaffold bias should not be a general concern for the models built using these data sets.

Figure 8 shows kNN and PLS models (characterized by q^2 and number of structure–data points) for 91 kinases (selected from 182 total models) with kNN $q^2 \geq 0.4$ and PLS $q^2 \geq 0.25$.

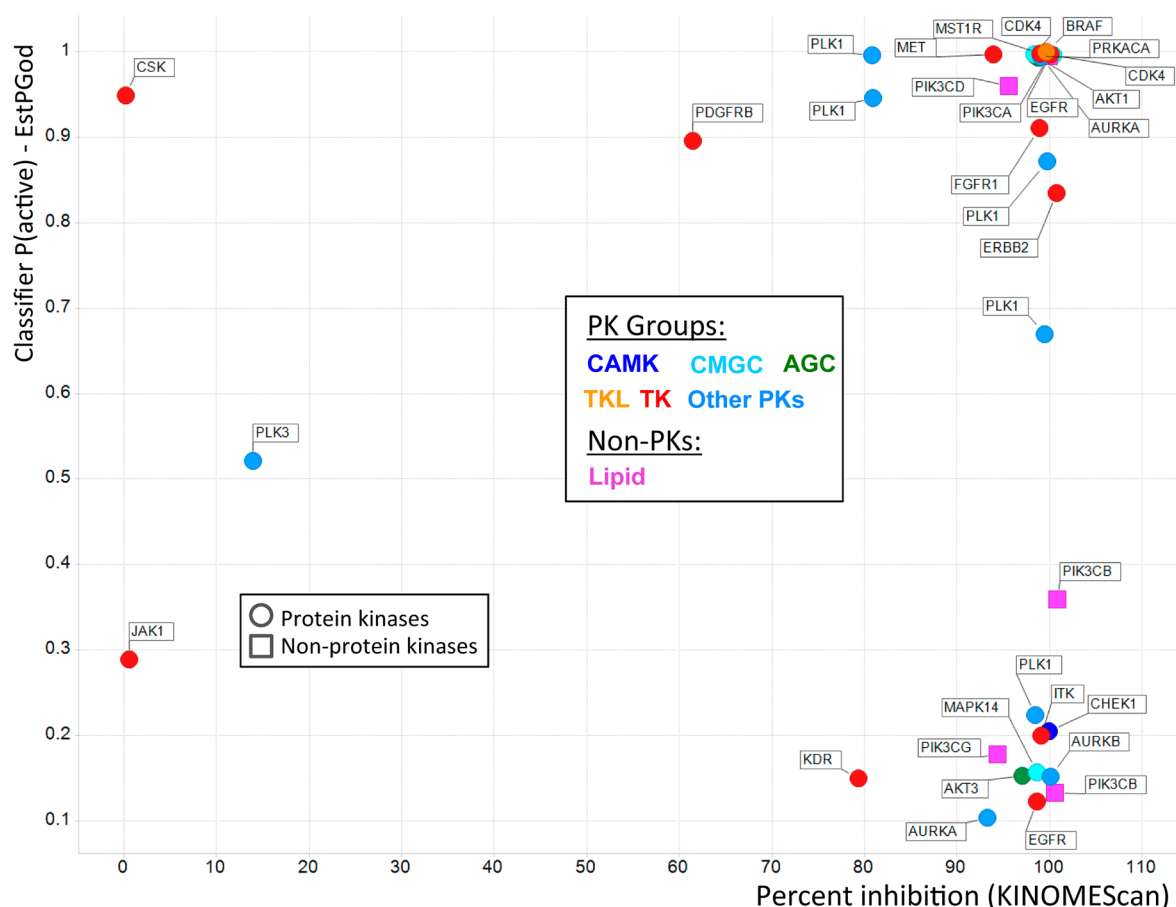


Figure 9. Probability (EstPGod) of compounds predicted active against a kinase based on KA-PI kinase classifiers and actual KINOMEScan percent inhibition values (at 10 μ M). Compare with Supporting Information Table S8. Kinases are classified by groups and protein vs nonprotein kinases.

generate and validate very high-quality kinase activity predictors for a large and diverse number of kinase targets.

From almost 500 000 kinase-related compounds, we extracted over 180 diverse structure–activity data sets covering all major groups of the human kinome. We applied different machine learning techniques including naïve Bayes binary classification and quantitative kNN and PLS regression. Rigorous cross-validation demonstrated reliable predictors for the majority of kinase targets if a minimum required number of active compounds or structure–activity data points were available. For the types of data sets investigated here, the best results for the largest number of kinases were obtained using Laplacian-corrected naïve Bayes classifiers trained (for any specific kinase) on known actives and a large background set of kinase-family focused presumed inactive compounds (KA-PI approach). This method resulted in very good ROC scores and very high enrichment rates for the majority of targets in particular for data sets with greater than 50 active compounds. The Laplacian-modified naïve Bayes classifiers generally improved when increasing the numbers of actives to a few hundred, but not much beyond that. Model domain applicability studies suggested that the classifiers are applicable to novel compounds and provide guidance to interpret virtual screening results. We applied the KA-PI classifiers to compounds recently profiled in the NIH LINCS project and found very good agreement of predicted kinase inhibition to actual screening results. Using kNN and PLS regression, we also obtained high-quality models for a large number of diverse

kinase targets; in particular for data sets with greater than 50 structure–activity records and spanning a wide activity range.

All data employed here were derived from biochemical concentration–response human (nonmutant) protein kinase assays. The data sets combined different assay methods, technologies, and various experimental procedures from numerous laboratories. After preprocessing of the data using a rigorous data standardization and aggregation procedure, our results indicate that the data sets are very well suitable to develop highly accurate predictors employing different machine learning techniques and are, by that measure, of high quality. The various kinase screening technologies and conditions appear to generate, on average, consistent results. This is supported by excellent predictivity of the various models and very good agreement of predicted kinase activity modeled on heterogeneous data compared to the KINOMEScan LINC results, which were all generated under the same conditions using the same competition binding assay. However, heterogeneity of the data may be one reason why we observe a distinct increase of model performance with a minimum number of active compounds or structure–activity data points. Depending on the number of actual actives in the specific data sets from which the models were built, the best KA-PI classifiers were able to retrieve between 40 and 80% true positives among only 0.1% of all compounds. These results suggest that they are practically applicable with great potential for virtual screening of large libraries across the kinome. Moreover, our results suggest that the KA-PI predictors are

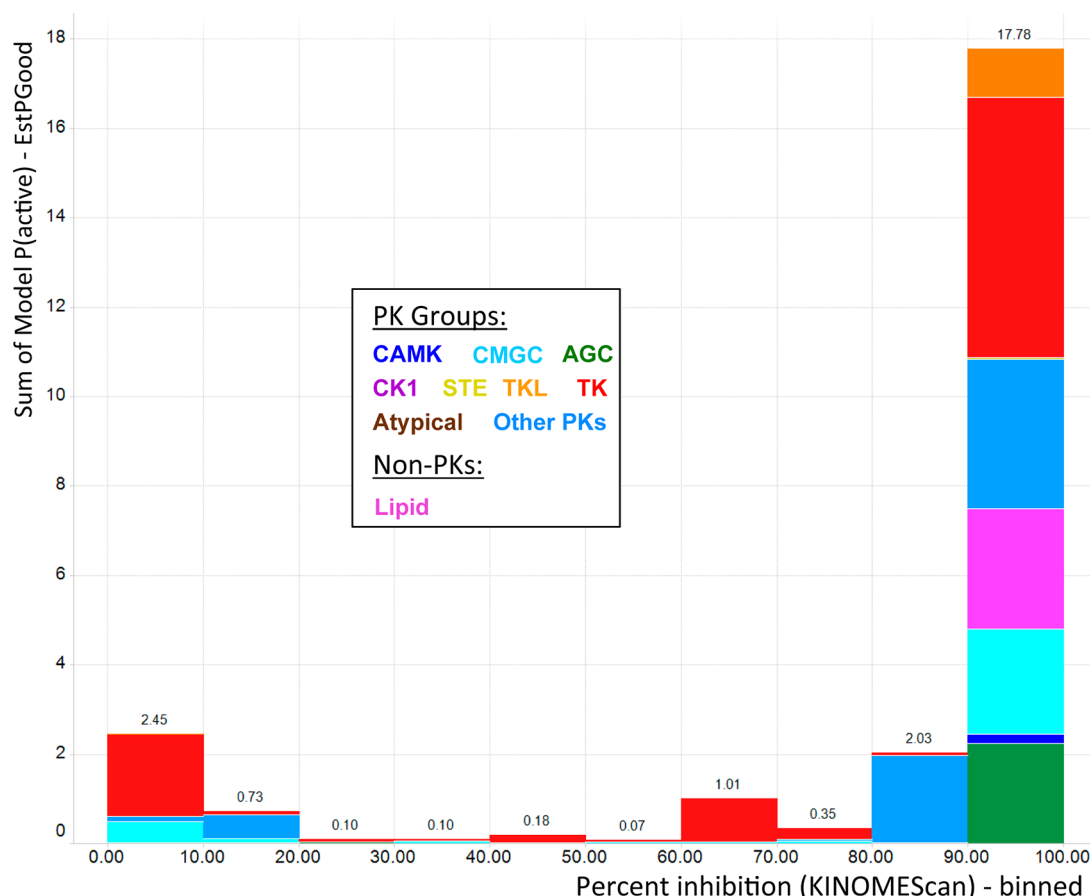


Figure 10. Aggregated predicted probabilities (EstPGood) of compounds being active against kinases (based on the KA-PI classifiers) as a function of the actual KINOMEScan percent inhibition ranges by category of kinase group and protein vs nonprotein kinase. 4796 activity data points for 43 compounds are mapped to KA-PI models (not all compounds were tested against the same number of targets).

useful to complement actual kinase profiling screening results, for example to identify likely false negatives.

■ ASSOCIATED CONTENT

Supporting Information

Additional supporting Figures S1–S12 and supporting Tables S1–S8 including all cross-validation results for all classification and quantitative models, domain applicability, ROC curves, and comparisons of predictions and kinase profiling results. This material is available free of charge via the Internet at <http://pubs.acs.org>.

■ AUTHOR INFORMATION

Corresponding Author

*E-mail: sschurer@med.miami.edu.

Notes

The authors declare no competing financial interest.

■ ACKNOWLEDGMENTS

This work was supported by the NIH grant U01 HL111561 (NIH LINCS program) and by the Center for Computational Science of the University of Miami.

■ REFERENCES

(1) Manning, G.; Whyte, D. B.; Martinez, R.; Hunter, T.; Sudarsanam, S. The protein kinase complement of the human genome. *Science* **2002**, *298*, 1912–1934.

(2) Akritopoulou-Zanze, I.; Hajduk, P. J. Kinase-targeted libraries: the design and synthesis of novel, potent, and selective kinase inhibitors. *Drug Discovery Today* **2009**, *14*, 291–297.

(3) Zhang, J.; Yang, P. L.; Gray, N. S. Targeting cancer with small molecule kinase inhibitors. *Nat. Rev. Cancer* **2009**, *9*, 28–39.

(4) Cohen, P. Protein kinases—the major drug targets of the twenty-first century? *Nat. Rev. Drug Discovery* **2002**, *1*, 309–315.

(5) Knight, Z. A.; Lin, H.; Shokat, K. M. Targeting the cancer kinome through polypharmacology. *Nat. Rev. Cancer* **2010**, *10*, 130–137.

(6) Morphy, R. Selectively Nonselective Kinase Inhibition: Striking the Right Balance. *J. Med. Chem.* **2010**, *53*, 1413–1437.

(7) Davies, S. P.; Reddy, H.; Caivano, M.; Cohen, P. Specificity and mechanism of action of some commonly used protein kinase inhibitors. *Biochem. J.* **2000**, *351*, 95–105.

(8) Zuccotto, F.; Ardini, E.; Casale, E.; Angiolini, M. Through the “gatekeeper door”: exploiting the active kinase conformation. *J. Med. Chem.* **2010**, *53*, 2681–2694.

(9) Brylinski, M.; Skolnick, J. Comprehensive Structural and Functional Characterization of the Human Kinome by Protein Structure Modeling and Ligand Virtual Screening. *J. Chem. Inf. Model* **2010**, *50*, 1839.

(10) Goldstein, D. M.; Gray, N. S.; Zarrinkar, P. P. High-throughput kinase profiling as a platform for drug discovery. *Nat. Rev. Drug Discovery* **2008**, *7*, 391–397.

(11) Karaman, M. W.; Herrgard, S.; Treiber, D. K.; Gallant, P.; Atteridge, C. E.; Campbell, B. T.; Chan, K. W.; Cicieri, P.; Davis, M. I.; Edeen, P. T.; Faraoni, R.; Floyd, M.; Hunt, J. P.; Lockhart, D. J.; Milanov, Z. V.; Morrison, M. J.; Pallares, G.; Patel, H. K.; Pritchard, S.; Wodicka, L. M.; Zarrinkar, P. P. A quantitative analysis of kinase inhibitor selectivity. *Nat. Biotechnol.* **2008**, *26*, 127–132.

- (12) Fedorov, O.; Marsden, B.; Pogacic, V.; Rellos, P.; Müller, S.; Bullock, A. N.; Schwaller, J.; Sundström, M.; Knapp, S. A systematic interaction map of validated kinase inhibitors with Ser/Thr kinases. *Proc. Natl. Acad. Sci. U.S.A.* **2007**, *104*, 20523–20528.
- (13) Fabian, M. A.; Biggs, W. H., 3rd; Treiber, D. K.; Atteridge, C. E.; Azimioara, M. D.; Benedetti, M. G.; Carter, T. A.; Ciceri, P.; Edeen, P. T.; Floyd, M.; Ford, J. M.; Galvin, M.; Gerlach, J. L.; Grotzfeld, R. M.; Herrgard, S.; Insko, D. E.; Insko, M. A.; Lai, A. G.; Lelias, J. M.; Mehta, S. A.; Milanov, Z. V.; Velasco, A. M.; Wodicka, L. M.; Patel, H. K.; Zarrinkar, P. P.; Lockhart, D. J. A small molecule-kinase interaction map for clinical kinase inhibitors. *Nat. Biotechnol.* **2005**, *23*, 329–336.
- (14) Library of Integrated Network-based Cellular Signatures (LINCS). <http://lincsproject.org/> (accessed Nov 30, 2012).
- (15) HMS LINCS DataBase. <http://lincs.hms.harvard.edu/db/> (accessed Nov 30, 2012).
- (16) LINCS Information Framework (LIFE). <http://lifekb.org/> (accessed Nov 30, 2012).
- (17) Vieth, M.; Sutherland, J. J.; Robertson, D. H.; Campbell, R. M. Kinomics: characterizing the therapeutically validated kinase space. *Drug Discovery Today* **2005**, *10*, 839–846.
- (18) Bamborough, P.; Drewry, D.; Harper, G.; Smith, G. K.; Schneider, K. Assessment of Chemical Coverage of Kinome Space and Its Implications for Kinase Drug Discovery. *J. Med. Chem.* **2008**, *51*, 7898–7914.
- (19) Vieth, M.; Higgs, R. E.; Robertson, D. H.; Shapiro, M.; Gragg, E. A.; Hemmerle, H. Kinomics-structural biology and chemogenomics of kinase inhibitors and targets. *Biochim. Biophys. Acta* **2004**, *1697*, 243–257.
- (20) Posy, S. L.; Hermsmeider, M. A.; Vaccaro, W.; Ott, K. H.; Todderud, G.; Lippy, J. S.; Trainor, G. L.; Loughney, D. A.; Johnson, S. R. Trends in Kinase Selectivity: Insights for Target Class-Focused Library Screening. *J. Med. Chem.* **2011**, *54*, 54–66.
- (21) Manallack, D. T.; Pitt, W. R.; Gancia, E.; Montana, J. G.; Livingstone, D. J.; Ford, M. G.; Whitley, D. C. Selecting screening candidates for kinase and G protein-coupled receptor targets using neural networks. *J. Chem. Inf. Comput. Sci.* **2002**, *42*, 1256–1262.
- (22) Xia, X.; Maliski, E. G.; Gallant, P.; Rogers, D. Classification of kinase inhibitors using a Bayesian model. *J. Med. Chem.* **2004**, *47*, 4463–4470.
- (23) Briem, H.; Gunther, J. Classifying “kinase inhibitor-likeness” by using machine-learning methods. *Chembiochem* **2005**, *6*, 558–566.
- (24) Sutherland, J. J.; Higgs, R. E.; Watson, I.; Vieth, M. Chemical Fragments as Foundations for Understanding Target Space and Activity Prediction. *J. Med. Chem.* **2008**, *51*, 2689–2700.
- (25) Aronov, A. M.; McClain, B.; Moody, C. S.; Murcko, M. A. Kinase-likeness and kinase-privileged fragments: toward virtual polypharmacology. *J. Med. Chem.* **2008**, *51*, 1214–1222.
- (26) Kinase Knowledge Base (Q4 2009). <http://eidogen-sertanty.com/kinasekb.php> (accessed Nov 30, 2012).
- (27) Pipeline Pilot 8.0, version 8.0; Accelrys: San Diego, CA, 2010.
- (28) Glick, M.; Jenkins, J. L.; Nettles, J. H.; Hitchings, H.; Davies, J. W. Enrichment of high-throughput screening data with increasing levels of noise using support vector machines, recursive partitioning, and Laplacian-modified naive Bayesian classifiers. *J. Chem. Inf. Model.* **2006**, *46*, 193–200.
- (29) Cramer, R. D., III Partial Least Squares (PLS): Its strengths and limitations. *Perspect. Drug Discovery Des.* **1993**, *1*, 269–278.
- (30) Hert, J.; Willett, P.; Wilton, D. J.; Acklin, P.; Azzaoui, K.; Jacoby, E.; Schuffenhauer, A. Comparison of topological descriptors for similarity-based virtual screening using multiple bioactive reference structures. *Org. Biomol. Chem.* **2004**, *2*, 3256–3266.
- (31) Rogers, D.; Brown, R. D.; Hahn, M. Using extended-connectivity fingerprints with Laplacian-modified Bayesian analysis in high-throughput screening follow-up. *J. Biomol. Screen.* **2005**, *10*, 682–686.
- (32) Sastry, M.; Lowrie, J. F.; Dixon, S. L.; Sherman, W. Large-Scale Systematic Analysis of 2D Fingerprint Methods and Parameters to Improve Virtual Screening Enrichments. *J. Chem. Inf. Model.* **2010**, *50*, 771–784.
- (33) Rogers, D.; Hahn, M. Extended-Connectivity Fingerprints. *Chem. Inf. Model.* **2010**, *50*, 742–754.
- (34) Li, Q.; Cheng, T.; Wang, Y.; Bryant, S. H. PubChem as a public resource for drug discovery. *Drug Discovery Today* **2010**, *15*, 1052–1057.
- (35) Liu, Y.; Gray, N. S. Rational design of inhibitors that bind to inactive kinase conformations. *Nat. Chem. Biol.* **2006**, *2*, 358–364.