

# Prediction of Drug-target interactions from heterogeneous information network based on LINE embedding model

**Bo-Ya Ji**

Xinjiang Technical Institutes of Physics and Chemistry, Chinese Academy of Sciences, Urumqi 830011, China. University of Chinese Academy of Sciences, Beijing 100049, China

**Zhu-Hong You** (✉ [zhuhongyou@ms.xjb.ac.cn](mailto:zhuhongyou@ms.xjb.ac.cn))

Xinjiang Technical Institutes of Physics and Chemistry, Chinese Academy of Sciences, Urumqi 830011, China. University of Chinese Academy of Sciences, Beijing 100049, China <https://orcid.org/0000-0001-8458-1751>

**Han-Jing Jiang**

Xinjiang Technical Institutes of Physics and Chemistry, Chinese Academy of Sciences, Urumqi 830011, China. University of Chinese Academy of Sciences, Beijing 100049, China

**Zhen-Hao Guo**

Xinjiang Technical Institutes of Physics and Chemistry, Chinese Academy of Sciences, Urumqi 830011, China. University of Chinese Academy of Sciences, Beijing 100049, China

**Kai Zheng**

School of Computer Science and Engineering, Central South University, Changsha, 410083, China

---

## Research

**Keywords:** drug-target interactions, heterogeneous information network, LINE, Random forest

**Posted Date:** May 13th, 2020

**DOI:** <https://doi.org/10.21203/rs.3.rs-16492/v2>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

**Version of Record:** A version of this preprint was published on September 7th, 2020. See the published version at <https://doi.org/10.1186/s12967-020-02490-x>.

# Prediction of Drug-target interactions from heterogeneous information network based on LINE embedding model

Bo-Ya Ji<sup>1,2</sup>, Zhu-Hong You<sup>1,2,\*</sup>, Han-Jing Jiang<sup>1,2</sup>, Zhen-Hao Guo<sup>1,2</sup>, Kai Zheng<sup>3</sup>

<sup>1</sup>Xinjiang Technical Institutes of Physics and Chemistry, Chinese Academy of Sciences, Urumqi 830011, China

<sup>2</sup>University of Chinese Academy of Sciences, Beijing 100049, China

<sup>3</sup>School of Computer Science and Engineering, Central South University, Changsha, 410083, China.

\* **Correspondence:** [zhuhongyou@ms.xjb.ac.cn](mailto:zhuhongyou@ms.xjb.ac.cn)

## Abstract

**Background:** The prediction of potential drug-protein target interactions (DTIs) not only provides a better comprehension of biological processes but also is critical for identifying new drugs. However, due to the disadvantages of expensive and high time-consuming traditional experiments, only a small section of interactions between drugs and targets in the database were verified experimentally. Therefore, it is meaningful and important to develop new computational methods with good performance for DTIs prediction. At present, many existing computational methods only utilize the single type of interactions between drugs and proteins without paying attention to the associations and influences with other types of molecules.

**Methods:** In this work, we developed a novel network embedding-based heterogeneous information integration model to predict potential drug-target interactions. Firstly, a heterogeneous information network is built by combining the known associations among protein, drug, lncRNA, disease, and miRNA. Secondly, the Large-scale Information Network Embedding (LINE) model is used to learn

27 behavior information (associations with other nodes) of drugs and proteins in the  
28 network. Hence, the known drug-protein interaction pairs can be represented as a  
29 combination of attribute information (e.g. protein sequences information and drug  
30 molecular fingerprints) and behavior information of themselves. Thirdly, the Random  
31 Forest classifier is used for training and prediction.

32 **Results:** In the results, under the 5-fold cross validation, our method obtained  
33 85.83% prediction accuracy with 80.47% sensitivity at the AUC of 92.33%. Moreover,  
34 in the case studies of three common drugs, the top 10 candidate targets have 8  
35 (Caffeine), 7 (Clozapine) and 6 (Pioglitazone) are respectively verified to be  
36 associated with corresponding drugs.

37 **Conclusions:** In short, these results indicate that our method can be a powerful  
38 tool for predicting potential drug-protein interactions and finding unknown targets for  
39 certain drugs or unknown drugs for certain targets.

40 **Keywords:** drug-target interactions; heterogeneous information network; LINE;  
41 Random forest

## 42 1 INTRODUCTION

43 Predicting potential drug-target interactions (DTIs) plays an important part in  
44 drug research and discovery. It not only helps researchers better understand biological  
45 processes but also reduces the failure rates and costs in the development of new drugs  
46 [1, 2]. However, there are still many difficulties in the prediction of drug-target  
47 interactions. For example, drugs have many positive and negative effects that are  
48 difficult to detect and clarify. In addition, different people respond differently to drugs,

49 even if the gene products are slightly different [3-6]. Moreover, the biological  
50 interactions in the human body are extremely complex, making it difficult to trace the  
51 effect of drugs. In the past few years, humans have made great efforts in predicting  
52 drug-target interactions to overcome these difficulties. With the completion of the  
53 Human Genome Project and the development of molecular medicine, more and more  
54 unknown drug-target interactions have been discovered. However, due to the high  
55 time-consuming, high cost and small research scope of the previous traditional  
56 experimental methods, the number of experimentally validated drug-target pairs is  
57 still very small. Therefore, this has spurred researchers to develop new computational  
58 methods to overcome these limitations to predict potential drug-target interactions  
59 [7-9].

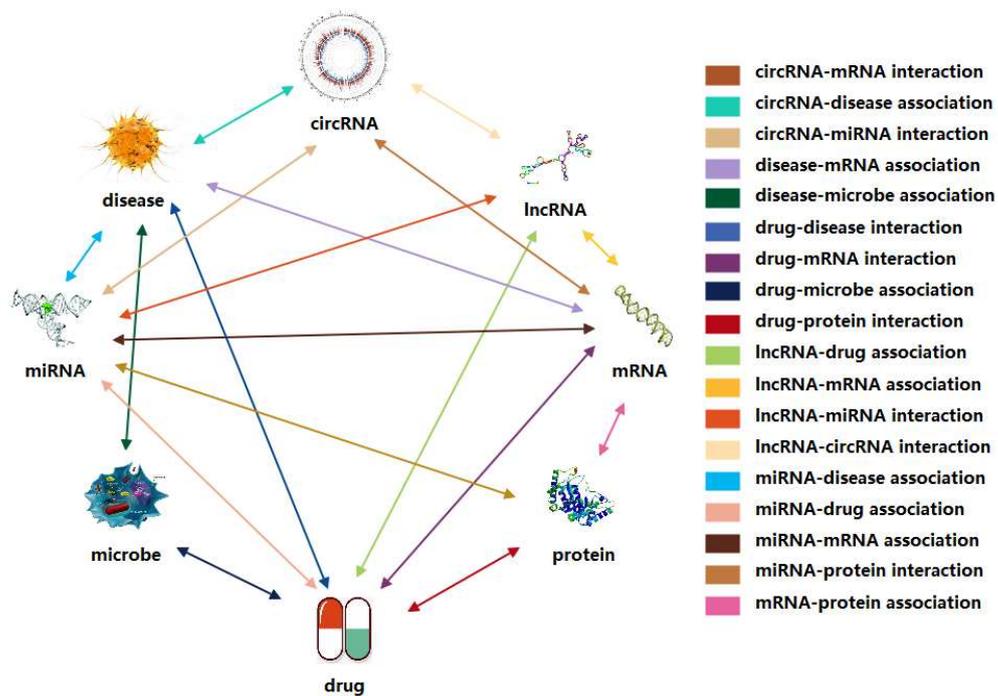
60 At present, a number of public online drug-target interaction databases, such as  
61 DrugBank [10], STITCH [11], KEGG [12] and ChEMBL [13], all store the major  
62 information about drugs and their interacting targets. These databases greatly facilitate  
63 the study of new methods involving drug-target interactions, and many existing  
64 calculation models are based on the known drug-target interactions in these databases  
65 to predict potential drug-target interactions. More specifically, these methods can be  
66 roughly divided into two categories: docking simulation and machine learning.  
67 However, the docking simulation method usually requires a three-dimensional (3D)  
68 structure of the target (traditional docking) or a larger set of drugs (reverse docking).  
69 Because of these limitations of the less known 3D structure of the target or the small  
70 size of the existing drug data sets or the high time-consuming, this method is often

71 difficult to conduct. Therefore, machine learning methods are more commonly used in  
72 the prediction of drug-target interactions. For example, Wang *et al.* [14] encoded the  
73 protein sequence as a position-specific scoring matrix (PSSM) descriptor to represent  
74 biological evolution information of proteins and encoded the drug molecules as a  
75 fingerprint feature vector to indicate the presence of a specific functional group or  
76 fragment. After that, the Rotation Forest classifier was adapted for the prediction of  
77 potential drug-target interactions. Wang *et al.* [15] used the stacked auto-encoder  
78 model in deep learning to fully extract drug molecular structure and protein sequence  
79 information. In this way, they generated highly representative features through  
80 multiple layers of iteration. Finally, the Rotation Forest classifier was used for the  
81 prediction of potential drug-target interactions and achieved good results. Meng *et al.*  
82 [16] developed a novel prediction model for the potential drug-target interactions  
83 based on the protein sequence. This method combined position-specific scoring  
84 matrix (PSSM), principal component analysis (PCA) with relevance vector machine  
85 (RVM) and bi-gram probabilities (BIGP), and had good effectiveness and robustness.  
86 Li *et al.* [17] proposed a computational model for the prediction of drug-target  
87 interactions, which used the position-specific scoring matrix (PSSM) of the target  
88 protein sequence information, the discriminant vector machine (DVM) classifier, the  
89 local binary pattern (LBP) histogram descriptor and the high-identification  
90 information of the drug-target interactions. The experimental results show that this  
91 method can effectively predict the potential drug-protein interactions. Huang *et al.* [18]  
92 exploited the pseudo substitution matrix representation (Pseudo-SMR) descriptors to

93 represent the protein sequence and used a new fingerprint feature vector to represent  
94 the drug signatures. After that, the two vector spaces are connected to represent the  
95 drug-protein interaction pairs. The final experimental results indicated that this  
96 method has a good performance for the prediction of the potential drug-protein  
97 interactions. Wen *et al.* [19] developed an algorithm framework based on deep  
98 learning to predict the potential drug-protein interactions. This approach solves the  
99 shortcomings of many traditional methods, which relied heavily on descriptors  
100 describing proteins and drugs, and can accurately predict the potential interactions  
101 between drugs and targets.

102 However, many existing computational methods only utilize the single-type of known  
103 drug-target association information without paying more attention to the associations  
104 between drugs and proteins and other biomolecules. In this work, we propose a novel  
105 computational model for predicting potential drug-target interactions. Firstly, we  
106 comprehensively analyzed and constructed a heterogeneous information network by  
107 combining known associations among disease, protein, drug, lncRNA, and miRNA  
108 from multiple databases as shown in Figure 1. In the network, the nodes and  
109 undirected edges among these nodes respectively represent lncRNAs, miRNAs,  
110 diseases, drugs and proteins, and interactions among them. In this way, the  
111 heterogeneous information network can help people more clearly understand the  
112 various life activities of living things [20, 21]. Secondly, the LINE [22] method is  
113 conducted to extract the association information between drugs and proteins and other  
114 nodes in the network, which we call the behavior information of drugs and proteins.

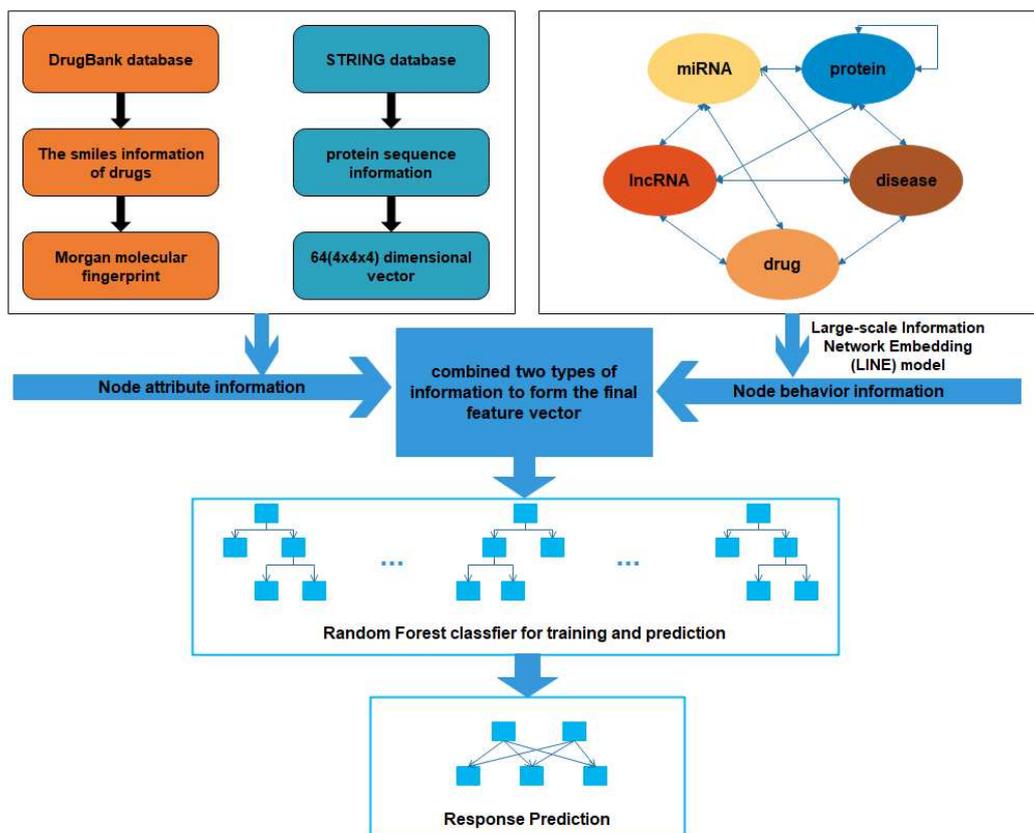
115 The LINE method can map tightly connected nodes in large networks to similar  
116 low-dimensional vector space locations. Thirdly, we integrate the attribute  
117 information (sequences of proteins and drugs' molecular fingerprints) and behavior  
118 information (associations with other molecules) to represent known drug-protein  
119 interaction pairs. Finally, the Random Forest classifier is applied for the training and  
120 prediction of the drug-target interactions. For the training samples in our model,  
121 11107 known drug-protein interaction pairs obtained from DrugBank 3.0 [10]  
122 databases are selected as positive sample sets, and the negative sample sets consist of  
123 the same number of randomly selected pairs of unrelated drugs and proteins. Figure 2  
124 shows the computation framework of our proposed model. In the results, our method  
125 was estimated under the five-fold cross-validation and achieved average the areas  
126 under the ROC curve (AUC) and the areas under the PR curve (AUPR) of 0.9233 and  
127 0.9301, respectively. In addition, we also compared the performance of different  
128 classifiers and different feature combinations of our method. Besides, in order to  
129 further estimate the performance of our model, we also conduct case studies of three  
130 major drugs. All these results fully demonstrate that our method has a good  
131 performance for drug-target interaction prediction in practical applications.



132

133

**Figure 1.** The heterogeneous association information network



134

135

**Figure 2.** Computation framework of our model

## 2 Materials and Methods

### 2.1 Combine eight kinds of associations to construct the heterogeneous information network

The heterogeneous association network is composed of known relationships among protein, drug, disease, miRNA, and lncRNA. We download these known associations from multiple databases and unify identifiers, remove redundant items, simplify and delete unrelated items. The final detailed data is shown in Table 1. In addition, we further counted the number of each node in the network. The final statistical results are shown in Table 2.

**Table 1.** The association information in the network

Association	Database	Amount
miRNA-lncRNA	lncRNASNP2[23]	8374
miRNA-disease	HMDD v3.0[24]	16427
miRNA-protein	miRTarBase:update 2018[25]	4944
lncRNA-disease	LncRNADisease[26], lncRNASNP2[23]	1264
drug-disease	CTD: update 2019[27]	18416
lncRNA-protein	LncRNA2Target v2.0[28]	690
protein-protein	STRING: in 2017[29]	19237
protein-disease	DisGeNET[30]	25087
Total	N/A	94439

**Table 2.** The node information in the network

Node	Amount
Drug	134
MiRNA	1023
Disease	2062
Protein	613
LncRNA	769
Total	4601

### 2.2. Drug Molecular Fingerprint

171 The Simplified Molecular Input Line Entry Specification (SMILES) of drugs mainly  
172 utilizes letters and symbols to indicate the structure of the compound for computer  
173 input. It is very different from traditional chemical formulas and has special writing  
174 rules. We download the drug's smiles from the DrugBank 3.0 [10] database and then  
175 convert the drug's smile to the relevant Morgan Molecular Fingerprint through using  
176 the RDKit python package.

### 177 **2.3 Protein Sequence Information**

178 The protein sequence information is derived from the STRING [29] database and used  
179 to represent the attribute information of the protein. After that, we choose the method  
180 in the article by Shen et al [31] to encode them. In this paper, according to the polarity  
181 of the side chain, 20 amino acids are divided into four categories including (Arg, Lys,  
182 and His); (Gly, Cys, Ser, Gln, Thr, Asn, and Tyr); (Ala, Ile, Trp, Val, Leu, Phe, Pro and  
183 Met); (Glu and Asp). In this way, each protein sequence can be represented as a  
184 64-dimensional vector, and each dimension denotes the occurrence frequency of a  
185 3-mer (e.g. UCC, AGU).

### 186 **2.4 Large-scale Information Network Embedding (LINE)**

187 As a novel network embedding method, LINE [32] mainly solves the problem of  
188 embedding large information networks into low-dimensional vector spaces. It can  
189 map closely connected nodes in a large network to similar low-dimensional vector  
190 space positions and is fully used for visualization, node classification, and link  
191 prediction. The LINE method is suitable for any type of information network and  
192 optimizes a well-designed objective function to retain both local and global network  
193 structure information. It not only considers the first-order proximity of nodes, that is,

194 two points are directly connected with an edge of higher power value, they are  
 195 considered to be more similar, but also considers the second-order proximity of nodes,  
 196 that is, two points may not be directly connected but is considered similar if they have  
 197 more public first-order proximity friends. Based on these two perspectives, the LINE  
 198 model can be divided into the following two categories:

199 Model 1: LINE with First-order Proximity

200 It should be noted that this model is only applicable to undirected graphs. For an  
 201 undirected edge  $(i, j)$ , the joint probabilities of the two vertex  $v_i$  and  $v_j$  defining this  
 202 edge is as follows:

$$203 \quad p_1(v_i, v_j) = \frac{1}{1 + \exp(-\vec{u}_i^T \vec{u}_j)} \quad (1)$$

204 where  $\vec{u}_i$  and  $\vec{u}_j$  are the low-dimensional vector representation of vertex  $v_i$  and  $v_j$ .

205 It is equivalent to describe the intimacy between vertices from the perspective of  
 206 embedding. Formula (1) defines the distribution  $p(*, *)$  on the space  $V \times V$ , and its  
 207 empirical probability can be defined as:

$$208 \quad \hat{p}_1(i, j) = \frac{w_{ij}}{W} \quad (2)$$

209 where  $w_{ij}$  represents the weight of the edge between vertex  $v_i$  and  $v_j$ , and  $W$   
 210 represents the sum of all weights of edges in the network. Our optimization goal is to  
 211 make the difference between  $p_1$  and  $\hat{p}_1$  as small as possible, so the objective  
 212 function can be defined as follows:

$$213 \quad O_1 = d(p_1(*, *), \hat{p}_1(*, *)) \quad (3)$$

214 where  $d()$  function is used to measure the difference between the two distributions.

215 Generally, the **Kullback-Leibler (KL)** divergence can be selected to replace the  $d(*,*)$ .

216 In this way, the KL divergence is brought into the above formula, and the constants

217 can be omitted (e.g.  $W$ ), the final optimized form can be obtained:

$$218 \quad O_1 = -\sum_{(i,j) \in E} w_{ij} \log p_1(v_i, v_j) \quad (4)$$

219 Therefore, we can represent each vertex in the  $d$ -dimensional space by finding the

220  $\{\vec{u}_i\}_{i=1 \dots |V|}$  which minimizes the objective in Eq.(4).

221 Model 2: LINE with Second-order Proximity

222 This model considers the effects of second-order relationships between nodes and is

223 suitable for both directed and undirected graphs. For a directed edge  $(i, j)$  (from  $i$  to  $j$ ),

224 the probability that vertex  $v_j$  is a neighbor of  $v_i$  can be represented as follows:

$$225 \quad p_2(v_j | v_i) = \frac{\exp(\vec{u}_j'^T \cdot \vec{u}_i)}{\sum_{k=1}^{|V|} \exp(\vec{u}_k'^T \cdot \vec{u}_i)} \quad (5)$$

226 where  $|V|$  represents the number of vertices. Next, in order to make the conditional

227 distribution of context  $p_2(\cdot | v_i)$  specified by the low-dimensional representation be

228 closed to the empirical distribution  $\hat{p}_2(\cdot | v_i)$ , which is defined as follows:

$$229 \quad \hat{p}_2(v_j | v_i) = \frac{w_{ij}}{d_i} \quad (6)$$

230 where  $d_i$  represents the out-degree of vertex  $i$  and  $w_{ij}$  represents the weight of the

231 edge, it is necessary to minimize the following formula:

$$232 \quad O_2 = \sum_{i \in V} \alpha_i d(\hat{p}_2(*, *), p_2(*, *)) \quad (7)$$

233 where  $\alpha_i$  represents the prestige of vertex  $i$  and can be measured by the degree or

234 estimated through an algorithm such as PageRank [33]. In this article, for convenience,

235 we set  $\alpha_i$  as the degree of vertex  $i$  and replace  $d(*, *)$  with KL-divergence. The Eq.(7)

236 can be finally optimized as follows:

$$237 \quad O_2 = -\sum_{(i,j) \in E} w_{ij} \log p_2(v_j | v_i) \quad (8)$$

238 Therefore, we can represent each vertex  $v_i$  with a  $d$ -dimensional vector  $\vec{u}_i$  via  
239 learning  $\{\vec{u}_i\}_{i=1 \dots |V|}$  and  $\{\vec{u}'_i\}_{i=1 \dots |V|}$  which minimizes this objective.

## 240 **2.5 The Receiver Operating Characteristic (ROC) and Precision-Recall (PR)** 241 **curve**

242 The Receiver Operating Characteristic (ROC) curve is a very important and common  
243 statistical analysis method. It sorts and predicts samples according to the prediction  
244 results of the classifier. In addition, it calculates the values of two important quantities  
245 each time: True Positive Rate (TPR) and False Positive Rate (FPR), which are  
246 respectively plotted on the horizontal and vertical coordinates. The AUC value is  
247 defined as the areas under the ROC curve and can be used as a numerical value to  
248 intuitively evaluate the quality of the classifier. Generally, the larger the AUC value,  
249 the more accurate the prediction result and the better the classification effect of the  
250 model. The Precision-Recall (PR) curve is also a method to test the capability of a  
251 classifier. Compared with the ROC curve, the PR curve can better reflect the  
252 performance of the classification when the proportion of positive and negative  
253 samples is large.

254

## 255 **2.6 Node Representation**

256 **Drugs and proteins are respectively represented by attribute information and behavior**  
257 **information (association information with other molecules) in the network we**  
258 **constructed.** Their attribute information is respectively sequences of proteins and  
259 molecular fingerprints of drugs. Besides, in this article, we choose a network

260 embedding model LINE to get the behavior information of them. In this way, the final  
261 128-dimensional feature vector contains 64-dimensional attribute information (protein  
262 sequences information and drug molecular fingerprints) and 64-dimensional behavior  
263 information (associations with other molecules) of drugs and targets. These two types  
264 of information are functionally similar and collaboratively provide information for the  
265 classifier to predict the potential associations between drugs and targets.

## 266 **3 Result and Discussion**

### 267 **3.1. Evaluation of our model under five-fold cross validation**

268 Cross-validation is a statistical analysis method for verifying the performance of a  
269 classifier to obtain a reliable and stable model. In this work, 5-fold cross-validation is  
270 conducted to estimate the performance of our model. 11107 known drug-target  
271 interaction pairs obtained from DrugBank 3.0 [10] database are used as training  
272 samples. In this way, we take 4/5 samples (training set) to build the model and leave  
273 1/5 sample (test set) to predict the newly built model. We repeat this experiment 5  
274 times so that the model can effectively avoid over- or under-learning, and the results  
275 obtained are more persuasive. In this article, we choose the following six common  
276 parameters as the evaluation indicators of our model: Accuracy (Acc.), Specificity  
277 (Spec.), Sensitivity (Sen.), Precision (Prec.), Matthews Correlation Coefficient (MCC),  
278 Areas under the ROC Curve (AUC). The detailed results of our method are shown in  
279 Table 3, and the last row of Table 3 shows the average value and their standard  
280 deviation of the results across 5 runs of the classifier.

281  
282

**Table 3.** Evaluation of our model under five-fold cross-validation

Fold	ACC.(%)	Spec.(%)	Prec.(%)	MCC(%)	Sen.(%)	AUC(%)
0	86.45	91.49	90.54	73.28	81.41	92.90
1	85.87	90.86	89.85	72.10	80.87	92.31
2	85.08	90.82	89.63	70.63	79.34	92.05
3	85.13	91.27	90.05	70.79	78.98	91.71
4	86.64	91.53	90.61	73.63	81.75	92.66
<b>Average</b>	<b>85.83±0.72</b>	<b>91.19±0.34</b>	<b>90.14±0.43</b>	<b>72.09±1.38</b>	<b>80.47±1.24</b>	<b>92.33±0.47</b>

283

284 Figure 3 and Figure 4 respectively show the ROC curves and AUC values, PR curves

285 and AUPR values of our model under five-fold cross validation. It can be seen from

286 the figure that the mean AUC and AUPR of our model are 0.9233 and 0.9301,

287 respectively. The results fully demonstrate that our proposed model has a good

288 performance for potential drug-target interactions prediction. Besides, the variance of

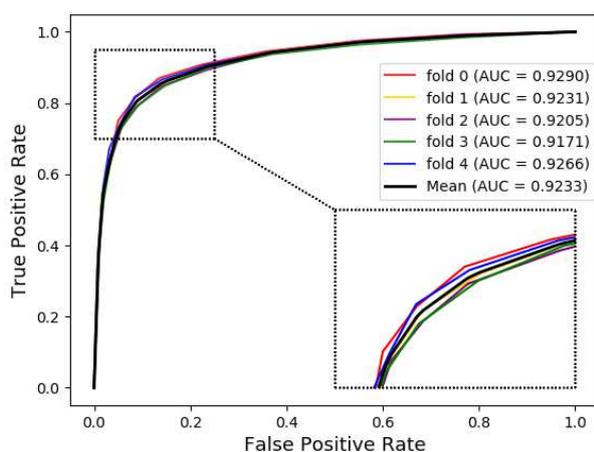
289 a model can describe the generalization ability of it. Generally, the larger the variance,

290 the easier the model is disturbed. On the contrary, the smaller the variance, the more

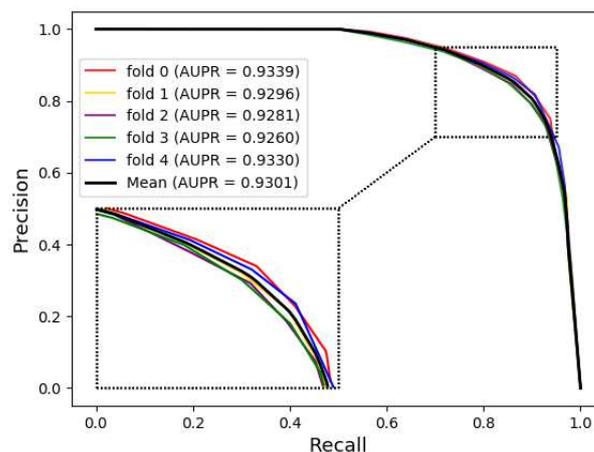
291 stable the model. In this work, the variance of the AUC for 5 runs of our model is

292 0.002%. The small variance can also prove that our method is stable for the prediction

293 of potential drug-target interactions.



**Figure 3.** The ROC curves of our model under five-fold cross-validation



**Figure 4.** The PR curves of our model under five-fold cross-validation

294

295 **3.2. Comparison of Different Feature Combinations**

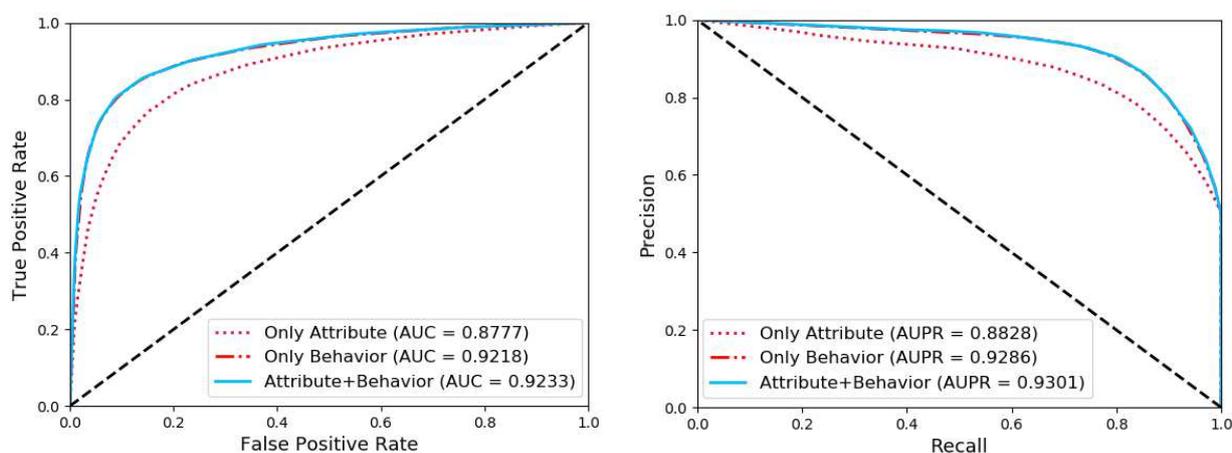
296 As we mentioned before, our approach utilizes a combination of attribute and  
 297 behavior information to represent known drug-protein interaction pairs. Hence, to test  
 298 the performance of different feature combinations on the results, we further conducted  
 299 experiments with three different feature combinations. **More specifically, we use only**  
 300 **attribute information, only behavior information, and the combination of attribute and**  
 301 **behavior information to respectively represent the drug and protein nodes.** After that,  
 302 the five-fold cross-validation experiment was conducted respectively. The  
 303 experimental environment and parameters of the three modes are consistent. Table 4  
 304 and Figure 5 show the detailed results of three models, and the classification results  
 305 are better when we utilize both the attribute and behavior information.

306  
307

**Table 4.** Comparison of different feature combinations

Feature	Acc.(%)	Spec.(%)	Prec.(%)	MCC(%)	Sen.(%)	AUC(%)
Attribute	80.73±0.79	84.36±1.05	83.14±1.04	61.63±1.61	77.11±0.60	87.77±0.83
Behavior	85.75±0.59	91.12±0.90	90.06±0.92	71.92±1.21	80.37±0.68	92.18±0.51
<b>Both</b>	<b>85.83±0.72</b>	<b>91.19±0.34</b>	<b>90.14±0.43</b>	<b>72.09±1.38</b>	<b>80.47±1.24</b>	<b>92.33±0.47</b>

308



309 **Figure 5.** Comparison of different feature combinations under five-fold cross validation

310

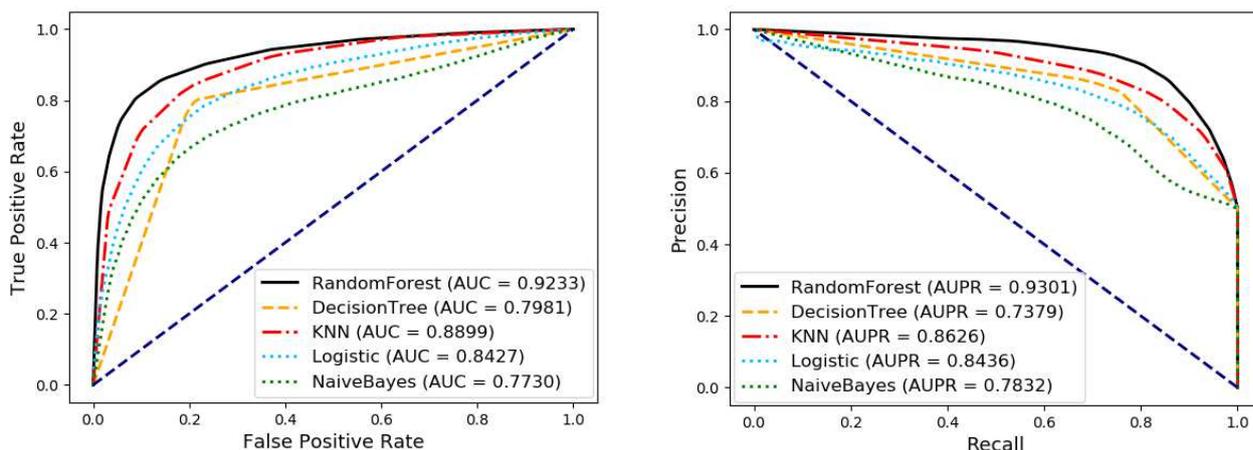
### 311 **3.3 Comparison of Different Machine Learning Classifiers**

312 To estimate the impact of different classifiers on the final results, we further  
313 respectively use Logistic, KNN, Naive Bayes, Decision Tree, and Random Forest  
314 classifier to perform five-fold cross-validation on our proposed model. **In particular,**  
315 **all the variables in the experiment are the same for the five classifiers, and all the**  
316 **classifiers use default parameters to make the comparative results more fair and**  
317 **reliable.** The detailed results can be founded in Table 5 and Figure 6. As can be seen  
318 from the results, the Random Forest classifier is not as good as KNN in sensitivity, but  
319 it has better performance for AUC and accuracy, which can better reflect the  
320 performance of our model. In conclusion, the Random Forest has a better performance  
321 than other classifiers and is more suitable for our method.

322 **Table 5.** Comparison of different machine learning classifiers

<b>Classifier</b>	<b>ACC.(%)</b>	<b>Spec.(%)</b>	<b>Prec.(%)</b>	<b>MCC(%)</b>	<b>Sen.(%)</b>	<b>AUC(%)</b>
Logistic	77.63±1.03	81.19±0.75	79.74±0.91	55.40±2.04	74.06±1.41	84.27±1.30
KNN	82.04±1.19	79.83±2.26	80.72±1.74	64.15±2.32	84.24±0.78	88.99±0.81
Naive Bayes	72.57±1.16	73.74±1.09	73.11±1.04	45.15±2.31	71.39±1.91	77.30±1.57
DecisionTree	79.81±0.66	79.73±1.29	79.78±1.01	59.63±1.32	79.89±0.60	79.81±0.66
<b>RandomForest</b>	<b>85.83±0.72</b>	<b>91.19±0.34</b>	<b>90.14±0.43</b>	<b>72.09±1.38</b>	<b>80.47±1.24</b>	<b>92.33±0.47</b>

323



**Figure 6.** Comparison of different machine learning classifiers under five-fold cross-validation

324

325

326

### 327 3.4. Case studies

328 To further estimate the performance of our model in practical applications, we select

329 three common drugs (Caffeine, Clozapine, and Pioglitazone) for case studies. These

330 three drugs are all closely related to human health and are often chosen by many

331 computational methods for case studies.

332 The chemical composition of Caffeine is 1,3,7-trimethylamine, which can be founded

333 in tea, coffee, cocoa, guarana and kola [34]. Recently, many researches have been

334 reported that caffeine may have an anti-cancer effect [35-37] and orally applied

335 caffeine can protect the skin from skin cancer caused by ultraviolet (UV) rays [38, 39].

336 Besides, transdermally applied caffeine can be used to treat skin cancer locally and

337 systemically.

338 Clozapine is a second-generation psychiatric drug. In addition, it has been proved that

339 clozapine is effective for psychotic positive and negative symptoms. Contrary to

340 concerns that typical antipsychotics may aggravate drug abuse, recent reports indicate

341 that clozapine has a reduced effect on nicotine, alcohol or other drug abuse in patients

342 with schizophrenia [40, 41]. Clozapine can also alleviate the emotional symptoms  
343 associated with schizophrenia (depression, guilt, anxiety), as well as the excitement  
344 and illusion of treatment for mania or other psychotic disorders.

345 Pioglitazone is a hypoglycemic drug that can be used alone or in combination with  
346 other hypoglycemic agents for the treatment of type 2 diabetes. The main function of  
347 this medicine is to reduce the insulin resistance in the body and enhance the  
348 sensitivity of the cells to insulin so that the body can make full use of the existing  
349 insulin to achieve the purpose of lowering blood sugar. At the same time, pioglitazone  
350 can improve the blood fat and pressure of the patient and reduce the blood vessels of  
351 the heart [42]. The drug has been well-tolerated by adult patients of all ages in clinical  
352 studies [43].

353 **Therefore, the identification of these three drugs' targets is of great importance.** More  
354 specifically, we utilize **the known drug-protein interactions in the DrugBank 3.0**  
355 **database of Knox *et al.* [10] as the training data set in the case studies.** One important  
356 fact that must be noted is that the known associations with the corresponding drug  
357 have been removed from the training data set to illustrate the applicability of our  
358 method to new drugs (drugs with no known related proteins). For the test data set, it  
359 contains proteins and corresponding drug interaction pairs in the heterogeneous  
360 association information network. After the prediction is complete, we rank all the  
361 proteins based on the predicted association scores and select the top 10 predicted  
362 targets to validate them using two databases on the relationship between drug and  
363 target, SuperTarget [44] and DrugBank 5.0 [45].

364 Table 6 shows the prediction result of the top 10 targets associated with caffeine, and  
 365 8 of which were successfully confirmed by the database. For example, the interaction  
 366 between cytochrome P450 1A2 (CYP1A2) and caffeine has been confirmed by  
 367 previous experiments [46]. The experiment proves that there is an interaction between  
 368 caffeine and CYP1A2 by studying the expression of CYP1A2 in mouse striatum.

369 Table 7 shows the prediction result of our method of the top 10 targets associated with  
 370 clozapine, 7 of which were successfully confirmed by the database. For example, the  
 371 interaction between cytochrome P450 1A2 and clozapine has been confirmed by  
 372 previous experiments [47].

373 Table 8 shows the prediction result of the top 10 targets associated with pioglitazone  
 374 using our method, 6 of which were successfully confirmed by the database. For  
 375 example, the interaction between cytochrome P450 3A4 and pioglitazone has been  
 376 confirmed by previous experiments [48]. This study evaluated the effect of  
 377 pioglitazone on the activity of cytochrome P450 3A4 (CYP3A4), demonstrating that  
 378 pioglitazone has a concentration-dependent inhibitory effect on CYP3A4 enzyme  
 379 activity.

380 **Table 6.** Prediction of the top 10 targets associated with Caffeine.

UniProt ID	Target	Evidence
9606.ensp00000342007	Cytochrome P450 1A2	SuperTarget
9606.ensp00000360372	Cytochrome P450 2C19	Unconfirmed
9606.ensp00000337915	Cytochrome P450 3A4	SuperTarget
9606.ensp00000478255	ATP-dependent translocase ABCB1	DrugBank

9606.ensp00000360317	Cytochrome P450 2C8	SuperTarget
9606.ensp00000260682	Cytochrome P450 2C9	SuperTarget
9606.ensp00000324648	Cytochrome P450 2B6	Unconfirmed
9606.ensp00000440689	Cytochrome P450 2E1	SuperTarget
9606.ensp00000353820	Cytochrome P450 2D6	SuperTarget
9606.ensp00000222982	Cytochrome P450 3A5	SuperTarget

381

382

383

**Table 7.** Prediction of the top 10 targets associated with Clozapine.

UniProt ID	Target	Evidence
9606.ensp00000478255	ATP-dependent translocase ABCB1	DrugBank
9606.ensp00000342007	Cytochrome P450 1A2	SuperTarget
9606.ensp00000360372	Cytochrome P450 2C19	SuperTarget
9606.ensp00000260682	Cytochrome P450 2C9	SuperTarget
9606.ensp00000337915	Cytochrome P450 3A4	SuperTarget
9606.ensp00000324648	Cytochrome P450 2B6	Unconfirmed
9606.ensp00000353820	Cytochrome P450 2D6	SuperTarget
9606.ensp00000222982	Cytochrome P450 3A5	SuperTarget
9606.ensp00000295897	Serum albumin	Unconfirmed
9606.ensp00000480571	Cytochrome P450 3A7	Unconfirmed

384

385

**Table 8.** Prediction of the top 10 targets associated with Pioglitazone.

UniProt ID	Target	Evidence
9606.ensp00000337915	Cytochrome P450 3A4	SuperTarget

9606.ensp00000478255	ATP-dependent translocase ABCB1	Unconfirmed
9606.ensp00000353820	Cytochrome P450 2D6	SuperTarget
9606.ensp00000367102	Solute carrier family 22 member 6	Unconfirmed
9606.ensp00000222982	Cytochrome P450 3A5	Unconfirmed
9606.ensp00000260682	Cytochrome P450 2C9	SuperTarget
9606.ensp00000360372	Cytochrome P450 2C19	DrugBank
9606.ensp00000369050	Cytochrome P450 1A1	Unconfirmed
9606.ensp00000360317	Cytochrome P450 2C8	SuperTarget
9606.ensp00000256958	Solute carrier organic anion transporter family member 1B1	DrugBank

386

387

## 388 **4 Conclusion**

389 The prediction of drug-target (protein) interactions is an important part of  
390 understanding the biological process and detecting new drugs. In this work, we put  
391 forward a novel network embedding-based heterogeneous information integration  
392 model for drug-protein interaction prediction. More specifically, we utilize the  
393 network embedding method LINE to obtain the behavior information (associations  
394 with other nodes) of drug and protein node in the network and then combine it with  
395 the intrinsic attribute information of them to represent the known drug-protein  
396 interaction pairs. Finally, the Random Forest classifier is selected to train and predict  
397 the transformed feature vectors. As a result, our proposed method has good  
398 performance for the potential drug-target interactions prediction under the five-fold  
399 cross-validation, and the prediction results are better than the model of using only

400 behavior information or attribute information. Besides, to further estimate the  
401 performance of our model, we also conduct case studies of three common drugs  
402 (Caffeine, Clozapine, and Pioglitazone). The results of case studies further indicate  
403 that our model performs well in predicting the potential drug-target interactions and  
404 targets associated with a given drug. Generally speaking, our proposed model can be  
405 an efficient tool for the prediction of potential drug-target interactions in the future.

## 406 **Declarations**

### 407 **Ethics approval and consent to participate**

408 Not applicable  
409

### 410 **Consent for publication**

411 Not applicable  
412

### 413 **Availability of data and material**

414 The datasets analyzed during the current study are available from the corresponding  
415 author on reasonable request.  
416

### 417 **Competing interests**

418 The authors declare that they have no competing interests.  
419

## 420 **Funding**

421 This work is supported by the NSFC Excellent Young Scholars Program, under  
422 Grants 61722212, in part by the National Science Foundation of China under Grants  
423 61873212, 61861146002, 61732012, in part by the West Light Foundation of the  
424 Chinese Academy of Sciences, Grants 2017-XBZG-BR-001.

425

## 426 **Authors' contributions**

427 B.Y.J. designed and carried out the experiment, prepared the data set and wrote the  
428 manuscript. Z.H.Y., H.J.J., Z.H.G. and K.Z. processed the data set and analyzed the  
429 experiment. All the authors contributed to the text of the manuscript.

430

## 431 **Acknowledgements**

432 ZHY was supported by the NSFC Excellent Young Scholars Program, under Grants  
433 61722212, in part by the National Science Foundation of China under Grants  
434 61873212, 61861146002, 61732012, in part by the West Light Foundation of the  
435 Chinese Academy of Sciences, Grants 2017-XBZG-BR-001. The authors would like  
436 to thank the editors and anonymous reviewers for their reviews.

437

438

## 439 **References**

- 440 1. Wang Y-C, Yang Z-X, Wang Y, Deng N-Y: **Computationally probing drug-protein interactions**  
441 **via support vector machine.** *Letters in Drug Design & Discovery* 2010, **7**:370-378.
- 442 2. Xia Z, Wu L-Y, Zhou X, Wong STC: **Semi-supervised drug-protein interaction prediction from**  
443 **heterogeneous biological spaces.** *BMC Systems Biology* 2010, **4**:S6.
- 444 3. Wang J-F, Wei D-Q, Li L, Zheng S-Y, Li Y-X, Chou K-C: **3D structure modeling of cytochrome**  
445 **P450 2C19 and its implication for personalized drug design.** *Biochemical and Biophysical*  
446 *Research Communications* 2007, **355**:513-519.
- 447 4. Wei D-Q, Wang J-F, Chen C, Li Y, Chou K-C: **Molecular modeling of two CYP2C19 SNPs and its**  
448 **implications for personalized drug design.** *Protein and peptide letters* 2008, **15**:27-32.
- 449 5. Wang J-F, Wei D-Q, Chou K-C: **Pharmacogenomics and personalized use of drugs.** *Current*  
450 *topics in medicinal chemistry* 2008, **8**:1573-1579.
- 451 6. Wang J-F, Zhang C-C, Chou K-C, Wei D-Q: **Structure of cytochrome p450s and personalized**  
452 **drug.** *Current medicinal chemistry* 2009, **16**:232-244.
- 453 7. Li Q, Lai L: **Prediction of potential drug targets based on simple sequence properties.** *Bmc*

- 454 *Bioinformatics* 2007, **8**:353.
- 455 8. Overington JP, Al-Lazikani B, Hopkins AL: **How many drug targets are there?** *Nature reviews*  
456 *Drug discovery* 2006, **5**:993.
- 457 9. Landry Y, Gies JP: **Drugs and their molecular targets: an updated overview.** *Fundamental &*  
458 *clinical pharmacology* 2008, **22**:1-18.
- 459 10. Knox C, Law V, Jewison T, Liu P, Ly S, Frolkis A, Pon A, Banco K, Mak C, Neveu V: **DrugBank 3.0:**  
460 **a comprehensive resource for 'omics' research on drugs.** *Nucleic acids research* 2010,  
461 **39**:D1035-D1041.
- 462 11. Kuhn M, von Mering C, Campillos M, Jensen LJ, Bork P: **STITCH: interaction networks of**  
463 **chemicals and proteins.** *Nucleic Acids Research* 2007, **36**:D684-D688.
- 464 12. Kanehisa M, Goto S, Sato Y, Furumichi M, Tanabe M: **KEGG for integration and interpretation**  
465 **of large-scale molecular data sets.** *Nucleic acids research* 2011, **40**:D109-D114.
- 466 13. Gaulton A, Bellis LJ, Bento AP, Chambers J, Davies M, Hersey A, Light Y, McGlinchey S,  
467 Michalovich D, Al-Lazikani B: **ChEMBL: a large-scale bioactivity database for drug discovery.**  
468 *Nucleic acids research* 2011, **40**:D1100-D1107.
- 469 14. Wang L, You Z-H, Chen X, Yan X, Liu G, Zhang W: **Rfdt: A rotation forest-based predictor for**  
470 **predicting drug-target interactions using drug structure and protein sequence information.**  
471 *Current Protein and Peptide Science* 2018, **19**:445-454.
- 472 15. Wang L, You Z-H, Chen X, Xia S-X, Liu F, Yan X, Zhou Y, Song K-J: **A computational-based**  
473 **method for predicting drug–target interactions by using stacked autoencoder deep neural**  
474 **network.** *Journal of Computational Biology* 2018, **25**:361-373.
- 475 16. Meng F-R, You Z-H, Chen X, Zhou Y, An J-Y: **Prediction of drug–target interaction networks**  
476 **from the integration of protein sequences and drug chemical structures.** *Molecules* 2017,  
477 **22**:1119.
- 478 17. Li Z, Han P, You Z-H, Li X, Zhang Y, Yu H, Nie R, Chen X: **In silico prediction of drug-target**  
479 **interaction networks based on drug chemical structure and protein sequences.** *Scientific*  
480 *reports* 2017, **7**:11174.
- 481 18. Huang Y-A, You Z-H, Chen X: **A systematic prediction of drug-target interactions using**  
482 **molecular fingerprints and protein sequences.** *Current Protein and Peptide Science* 2018,  
483 **19**:468-478.
- 484 19. Wen M, Zhang Z, Niu S, Sha H, Yang R, Yun Y, Lu H: **Deep-learning-based drug–target**  
485 **interaction prediction.** *Journal of proteome research* 2017, **16**:1401-1409.
- 486 20. Hrdlickova B, de Almeida RC, Borek Z, Withoff S: **Genetic variation in the non-coding genome:**  
487 **Involvement of micro-RNAs and long non-coding RNAs in disease.** *Biochimica et Biophysica*  
488 *Acta (BBA)-Molecular Basis of Disease* 2014, **1842**:1910-1922.
- 489 21. Barabasi A-L, Oltvai ZN: **Network biology: understanding the cell's functional organization.**  
490 *Nature reviews genetics* 2004, **5**:101.
- 491 22. Tang J, Qu M, Wang M, Zhang M, Yan J, Mei Q: **Line: Large-scale information network**  
492 **embedding.** In *Proceedings of the 24th international conference on world wide web.*  
493 International World Wide Web Conferences Steering Committee; 2015: 1067-1077.
- 494 23. Miao Y-R, Liu W, Zhang Q, Guo A-Y: **lncRNASNP2: an updated database of functional SNPs**  
495 **and mutations in human and mouse lncRNAs.** *Nucleic acids research* 2017, **46**:D276-D280.
- 496 24. Huang Z, Shi J, Gao Y, Cui C, Zhang S, Li J, Zhou Y, Cui Q: **HMDD v3. 0: a database for**  
497 **experimentally supported human microRNA–disease associations.** *Nucleic acids research*

- 498 2018, **47**:D1013-D1017.
- 499 25. Chou C-H, Shrestha S, Yang C-D, Chang N-W, Lin Y-L, Liao K-W, Huang W-C, Sun T-H, Tu S-J, Lee  
500 W-H: **miRTarBase update 2018: a resource for experimentally validated microRNA-target**  
501 **interactions**. *Nucleic acids research* 2017, **46**:D296-D302.
- 502 26. Chen G, Wang Z, Wang D, Qiu C, Liu M, Chen X, Zhang Q, Yan G, Cui Q: **LncRNADisease: a**  
503 **database for long-non-coding RNA-associated diseases**. *Nucleic acids research* 2012,  
504 **41**:D983-D986.
- 505 27. Davis AP, Grondin CJ, Johnson RJ, Sciaky D, McMorran R, Wieggers J, Wieggers TC, Mattingly CJ:  
506 **The comparative toxicogenomics database: update 2019**. *Nucleic acids research* 2018,  
507 **47**:D948-D954.
- 508 28. Cheng L, Wang P, Tian R, Wang S, Guo Q, Luo M, Zhou W, Liu G, Jiang H, Jiang Q:  
509 **LncRNA2Target v2. 0: a comprehensive database for target genes of lncRNAs in human and**  
510 **mouse**. *Nucleic acids research* 2018, **47**:D140-D144.
- 511 29. Szklarczyk D, Morris JH, Cook H, Kuhn M, Wyder S, Simonovic M, Santos A, Doncheva NT, Roth  
512 A, Bork P: **The STRING database in 2017: quality-controlled protein-protein association**  
513 **networks, made broadly accessible**. *Nucleic acids research* 2016:gkw937.
- 514 30. Piñero J, Bravo À, Queralt-Rosinach N, Gutiérrez-Sacristán A, Deu-Pons J, Centeno E,  
515 García-García J, Sanz F, Furlong LI: **DisGeNET: a comprehensive platform integrating**  
516 **information on human disease-associated genes and variants**. *Nucleic acids research*  
517 2016:gkw943.
- 518 31. Shen J, Zhang J, Luo X, Zhu W, Yu K, Chen K, Li Y, Jiang H: **Predicting protein-protein**  
519 **interactions based only on sequences information**. *Proceedings of the National Academy of*  
520 *Sciences* 2007, **104**:4337-4341.
- 521 32. Tang J, Qu M, Wang M, Zhang M, Yan J, Mei Q: **LINE: Large-scale Information Network**  
522 **Embedding**. In *Proceedings of the 24th International Conference on World Wide Web*. pp.  
523 1067–1077. Florence, Italy: International World Wide Web Conferences Steering Committee;  
524 2015:1067–1077.
- 525 33. Page L, Brin S, Motwani R, Winograd T: **The PageRank citation ranking: Bringing order to the**  
526 **web**. Stanford InfoLab; 1999.
- 527 34. Murray SD, Hansen PJ: **The extraction of caffeine from tea: An old undergraduate**  
528 **experiment revisited**. *Journal of chemical education* 1995, **72**:851.
- 529 35. Sarkaria JN, Busby EC, Tibbetts RS, Roos P, Taya Y, Karnitz LM, Abraham RT: **Inhibition of ATM**  
530 **and ATR kinase activities by the radiosensitizing agent, caffeine**. *Cancer research* 1999,  
531 **59**:4375-4382.
- 532 36. Sabisz M, Skladanowski A: **Modulation of cellular response to anticancer treatment by**  
533 **caffeine: inhibition of cell cycle checkpoints, DNA repair and more**. *Current pharmaceutical*  
534 *biotechnology* 2008, **9**:325-336.
- 535 37. Tsuchiya H, Wan S, Sakayama K, Yamamoto N, Nishida H, Tomita K: **Reconstruction using an**  
536 **autograft containing tumour treated by liquid nitrogen**. *The Journal of bone and joint*  
537 *surgery British volume* 2005, **87**:218-225.
- 538 38. Lu Y-P, Lou Y-R, Lin Y, Shih WJ, Huang M-T, Yang CS, Conney AH: **Inhibitory effects of orally**  
539 **administered green tea, black tea, and caffeine on skin carcinogenesis in mice previously**  
540 **treated with ultraviolet B light (high-risk mice): relationship to decreased tissue fat**. *Cancer*  
541 *research* 2001, **61**:5002-5009.

- 542 39. Lu Y-P, Lou Y-R, Peng Q-Y, Xie J-G, Nghiem P, Conney AH: **Effect of caffeine on the ATR/Chk1**  
543 **pathway in the epidermis of UVB-irradiated mice.** *Cancer research* 2008, **68**:2523-2529.
- 544 40. Marcus P, Snyder R: **Reduction of comorbid substance abuse with clozapine.** *The American*  
545 *journal of psychiatry* 1995.
- 546 41. McEvoy JP, Freudenreich O, Levin ED, Rose JE: **Haloperidol increases smoking in patients**  
547 **with schizophrenia.** *Psychopharmacology* 1995, **119**:124-126.
- 548 42. Sanyal AJ, Chalasani N, Kowdley KV, McCullough A, Diehl AM, Bass NM, Neuschwander-Tetri  
549 BA, Lavine JE, Tonascia J, Unalp A: **Pioglitazone, vitamin E, or placebo for nonalcoholic**  
550 **steatohepatitis.** *New England Journal of Medicine* 2010, **362**:1675-1685.
- 551 43. Gillies PS, Dunn CJ: **Pioglitazone.** *Drugs* 2000, **60**:333-343.
- 552 44. Günther S, Kuhn M, Dunkel M, Campillos M, Senger C, Petsalaki E, Ahmed J, Urdiales EG,  
553 Gewiess A, Jensen LJ: **SuperTarget and Matador: resources for exploring drug-target**  
554 **relationships.** *Nucleic acids research* 2007, **36**:D919-D922.
- 555 45. Wishart DS, Feunang YD, Guo AC, Lo EJ, Marcu A, Grant JR, Sajed T, Johnson D, Li C, Sayeeda Z:  
556 **DrugBank 5.0: a major update to the DrugBank database for 2018.** *Nucleic acids research*  
557 2017, **46**:D1074-D1082.
- 558 46. Singh S, Singh K, Gupta SP, Patel DK, Singh VK, Singh RK, Singh MP: **Effect of caffeine on the**  
559 **expression of cytochrome P450 1A2, adenosine A2A receptor and dopamine transporter in**  
560 **control and 1-methyl 4-phenyl 1, 2, 3, 6-tetrahydropyridine treated mouse striatum.** *Brain*  
561 *research* 2009, **1283**:115-126.
- 562 47. Olesen OV, Linnet K: **Contributions of five human cytochrome P450 isoforms to the N -**  
563 **demethylation of clozapine in vitro at low and high concentrations.** *The Journal of Clinical*  
564 *Pharmacology* 2001, **41**:823-832.
- 565 48. Choi J-S, Choi I, Choi D-H: **Effects of pioglitazone on the pharmacokinetics of nifedipine and**  
566 **its main metabolite, dehydronifedipine, in rats.** *European journal of drug metabolism and*  
567 *pharmacokinetics* 2016, **41**:231-238.
- 568

# Figures

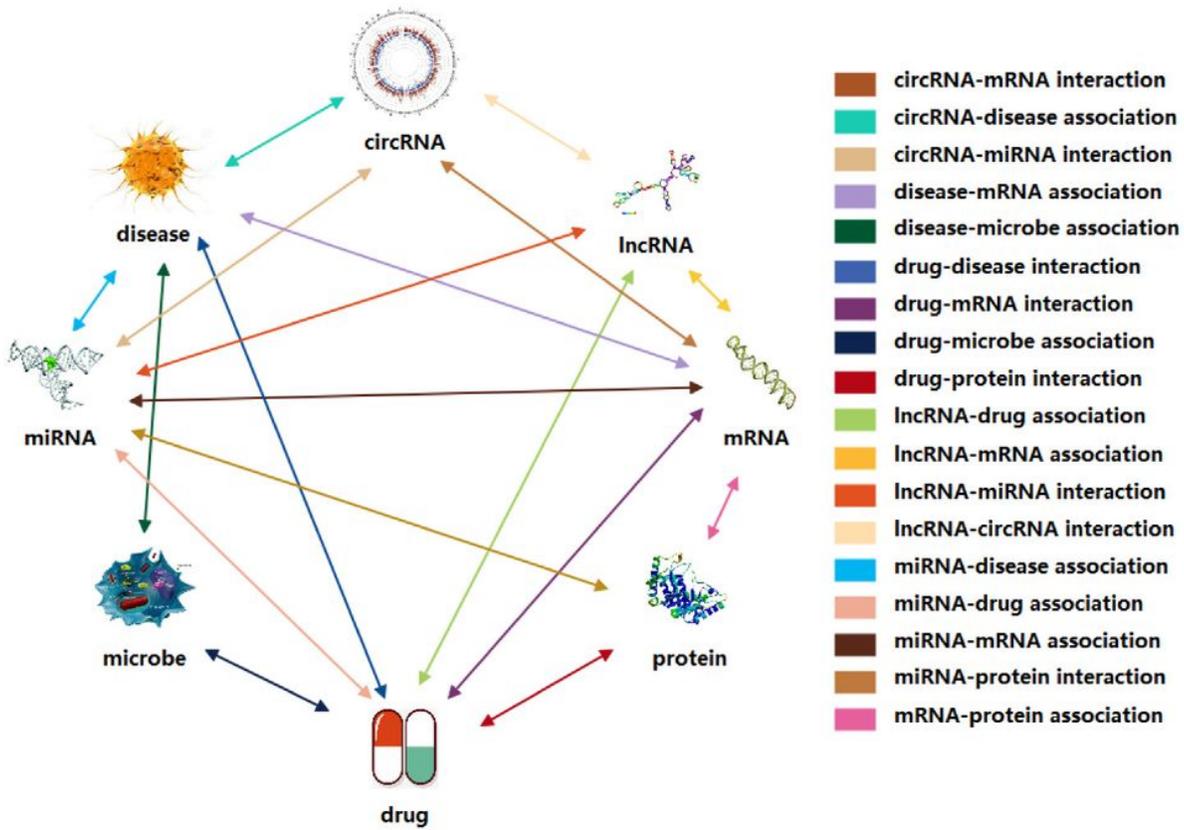


Figure 1

The heterogeneous association information network

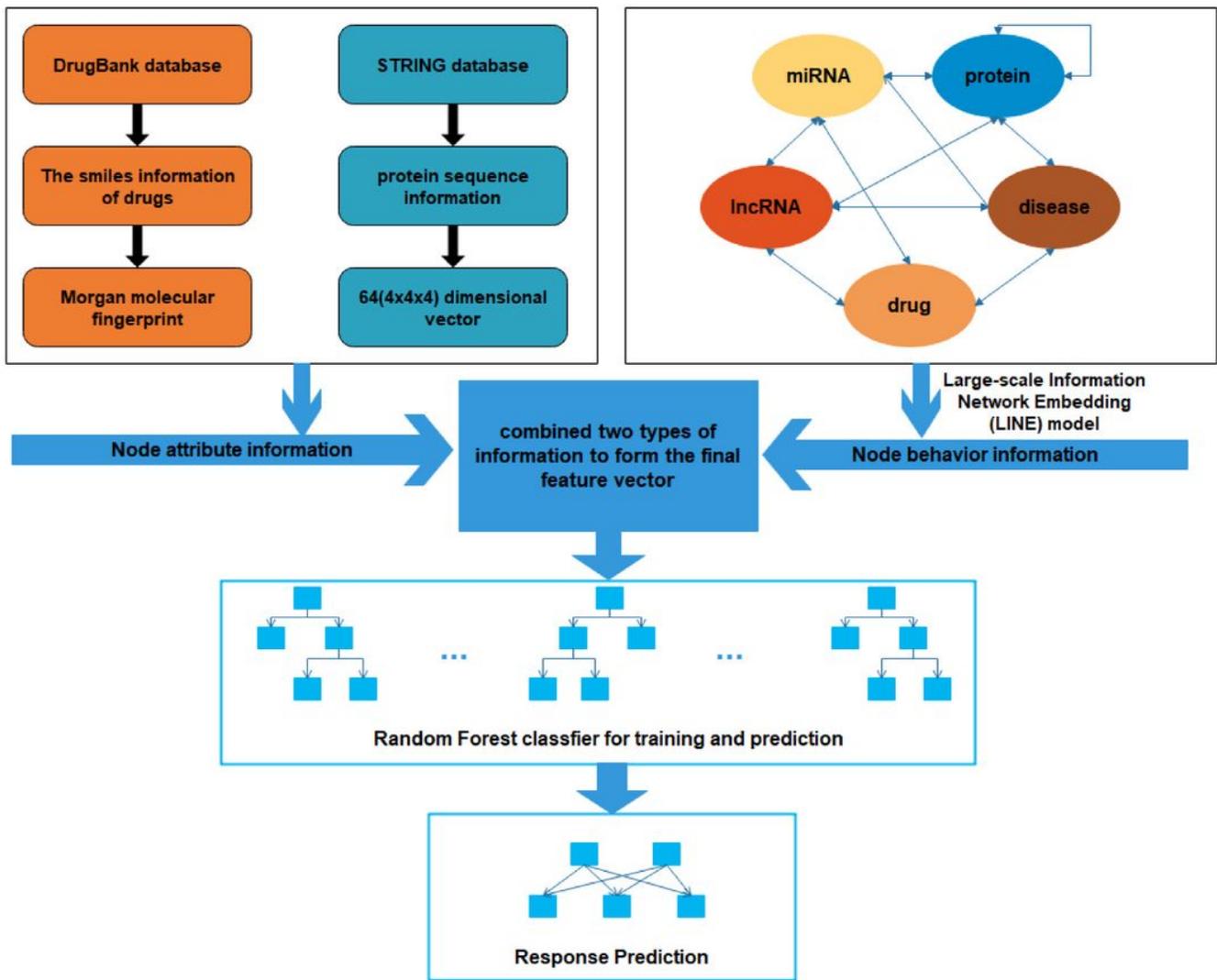


Figure 2

Computation framework of our model

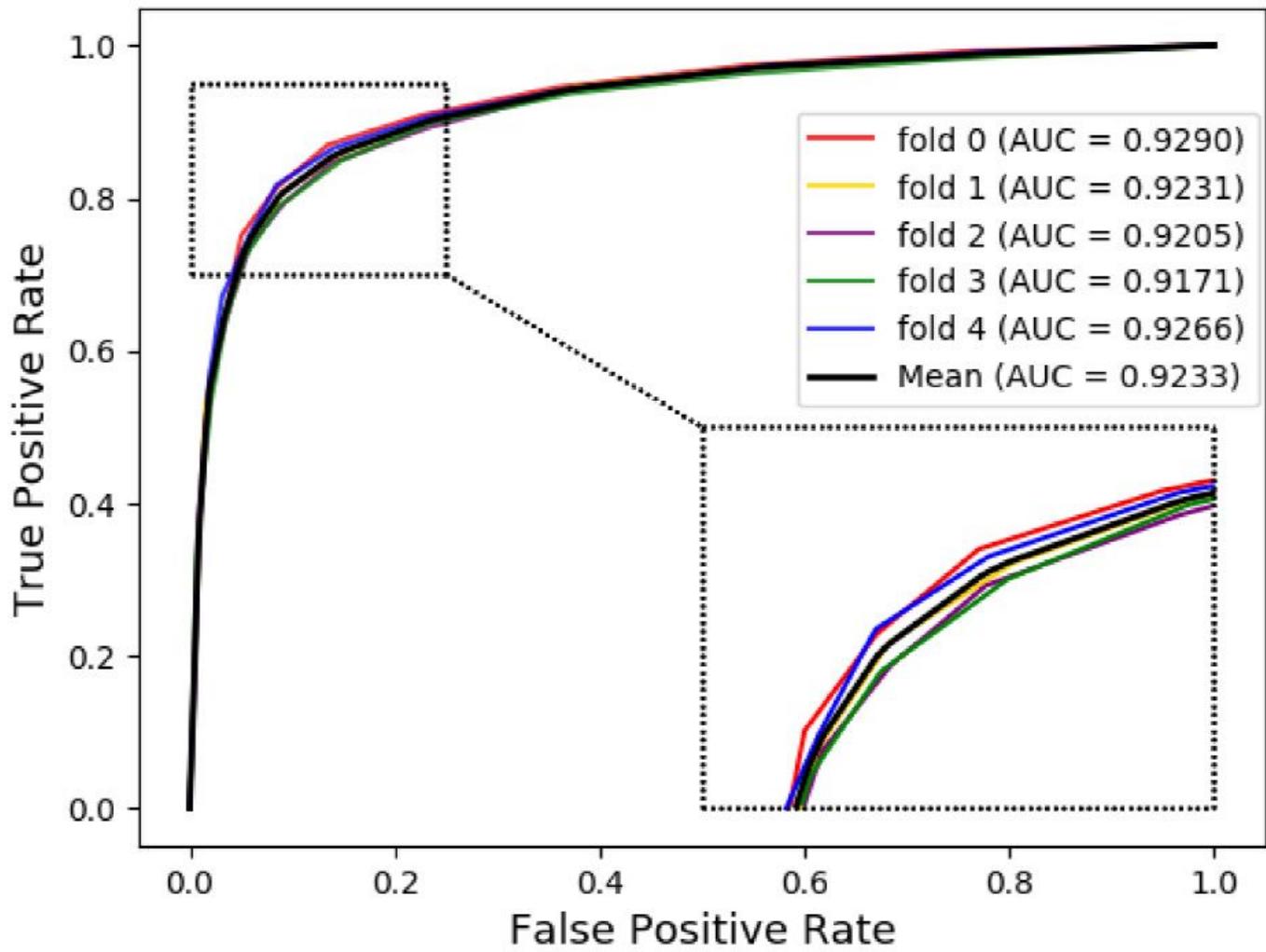


Figure 3

The ROC curves of our model under five-fold cross-validation

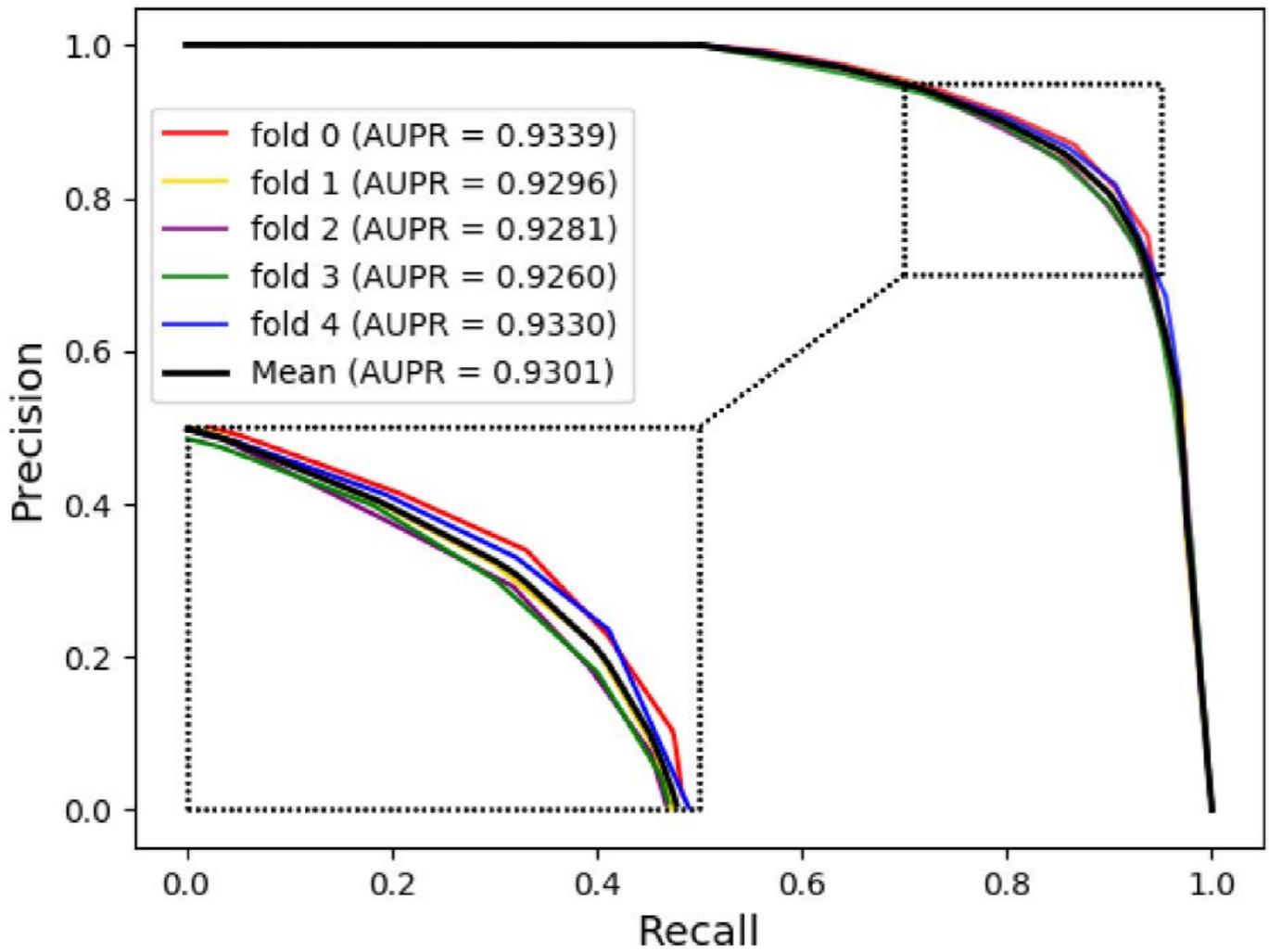
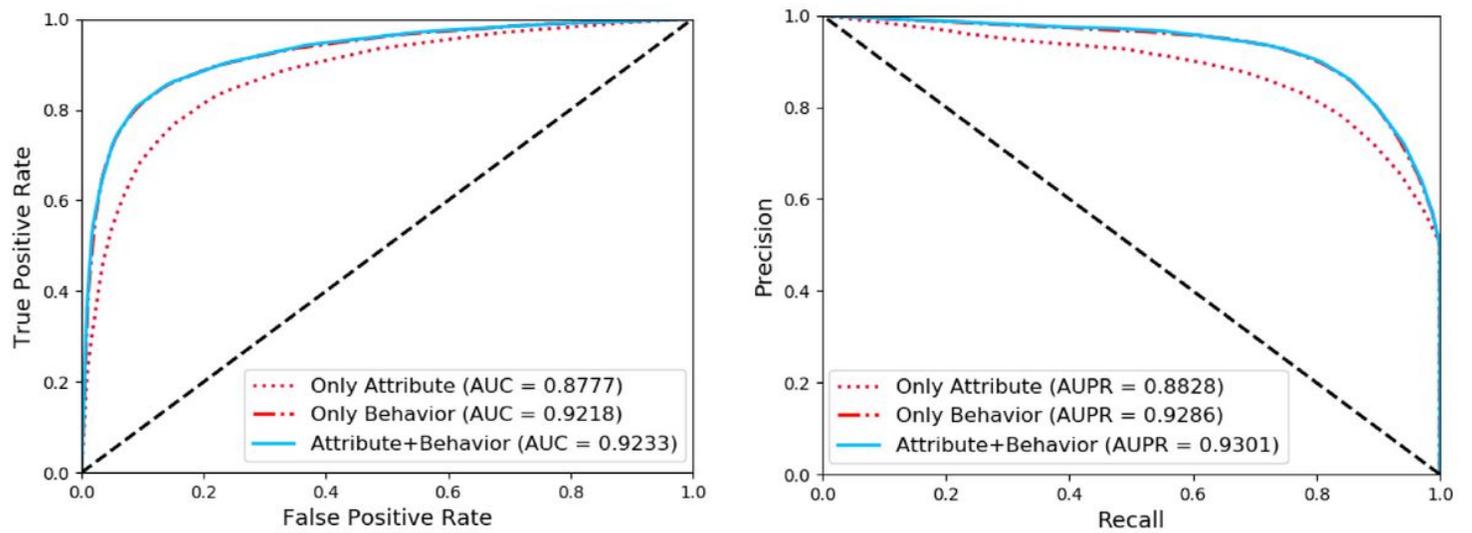


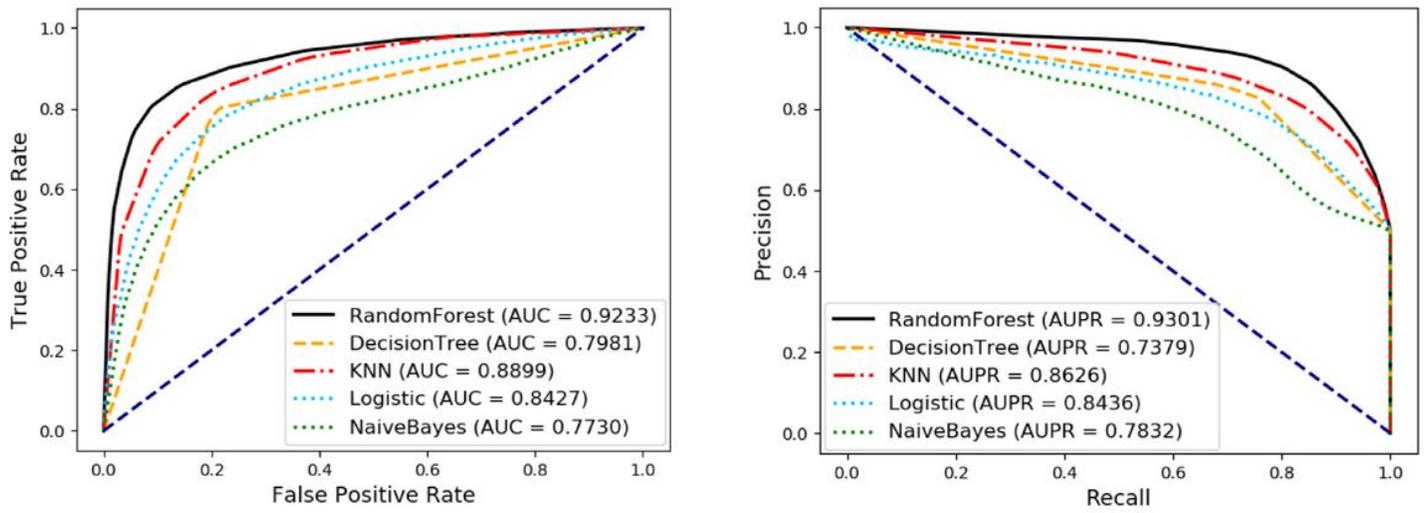
Figure 4

The PR curves of our model under five-fold cross-validation



**Figure 5**

Comparison of different feature combinations under five-fold cross validation



**Figure 6**

Comparison of different machine learning classifiers under five-fold cross-validation