

Drug-Induced Cell Viability Prediction from LINCS-L1000 through WRFEN-XGBoost Algorithm

Jiaxing Lu

the College of Information Technology, Shanghai Ocean University

Ming Chen

the College of Information Technology, Shanghai Ocean University

Yufang Qin (✉ yfqin@shou.edu.cn)

Shanghai Ocean University <https://orcid.org/0000-0002-8906-8727>

Research article

Keywords: cell viability, drug sensitivity, perturbation signatures, WRFEN-XGBoost algorithm

Posted Date: December 1st, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-56410/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 January 6th, 2021. See the published version at <https://doi.org/10.1186/s12859-020-03949-w>.

RESEARCH

Drug-Induced Cell Viability Prediction from LINCS-L1000 through WRFEN-XGBoost Algorithm

Jiaxing Lu, Ming Chen* and Yufang Qin*

*Correspondence:

mchen@shou.edu.cn;

yfqin@shou.edu.cn

College of Information Technology,
Shanghai Ocean University,
Hucheng Ring Road, Shanghai,
China

Full list of author information is
available at the end of the article

Abstract

Background: Predicting the drug response of the cancer diseases through the cellular perturbation signatures under the action of specific compounds is very important in personalized medicine. In the process of testing drug responses to the cancer, traditional experimental methods have been greatly hampered by the cost and sample size. At present, the public availability of large amounts of gene expression data makes it a challenging task to use machine learning methods to predict the drug sensitivity.

Results: In this study, we introduced the WRFEN-XGBoost cell viability prediction algorithm based on LINCS-L1000 cell signatures. We integrated the LINCS-L1000, CTRP and Achilles datasets and adopted a weighted fusion algorithm based on random forest and elastic net for key gene selection. Then the FEBPSO algorithm was introduced into XGBoost learning algorithm to predict the cell viability induced by the drugs. The proposed method was compared with some new methods, and it was found that our model achieved good results with 0.83 Pearson correlation. At the same time, we completed the drug sensitivity validation on the NCI60 and CCLE datasets, which further demonstrated the effectiveness of our method.

Conclusions: The results showed that our method was conducive to the elucidation of disease mechanisms and the exploration of new therapies, which greatly promoted the progress of clinical medicine.

Keywords: cell viability; drug sensitivity; perturbation signatures; WRFEN-XGBoost algorithm

¹Background

²In recent years, the study of cell death process has always been the hot topics²
³in biology and medicine [1]. With the development of cell biology and molecular³
⁴biology, the mechanism of cell death has gradually been revealed. Programmed cell⁴
⁵death was induced by many factors, including external factors such as radiation,⁵
⁶drugs and viral infections, and internal factors such as tumors, autoimmunity and⁶
⁷degenerative diseases [2]. It has been reported that the cell viability mechanism⁷
⁸could be used to stimulate and inhibit the apoptosis of tumor cells through the⁸
⁹action of the compounds. Changes in the proportion of apoptosis and abnormal⁹
¹⁰behavior of cell proliferation are highly correlated with compound concentration and¹⁰
¹¹perturbation time, which is one of the key factors for the formation and development¹¹
¹²of tumor cells [3]. With the emergence of more canceromics data, it is still a challenge¹²
¹³to apply cell activity mechanisms to design the best intervention strategy for the¹³
¹⁴duration of the drug action, and to construct a cell signaling model to interpret¹⁴
¹⁵these data and make accurate predictions [4].¹⁵

¹⁶Cell perturbation signatures are closely related to the cell viability with the action¹⁶
¹⁷of the compounds. In the study of drug sensitivity and anticancer drug response¹⁷
¹⁸prediction, we can predict cell phenotypes from different high-coverage molecular¹⁸
¹⁹data since compounds control the expression and function of target proteins or¹⁹
²⁰enzymes in the apoptotic pathway and induce abnormal cell apoptosis. Because²⁰
²¹clinical collection of experimental data on patient and drug interactions are expen-²¹
²²sive and impractical, it was expected that the preclinical prediction models based²²
²³on large-scale pharmacogenomics of cancer cell lines could be applied. In recent²³
²⁴years, the prediction model scheme designed by machine learning method from the²⁴
²⁵perspective of cell viability research has made breakthrough progress. Based on the²⁵
²⁶genomic background of each cell lines, Michael P. Menden. et al. trained a neural²⁶
²⁷network model to predict its IC₅₀ distribution throughout the cell lines [5]. Due to²⁷
²⁸the high-dimensional and nonlinear nature of the omics data, Yongcui Wang et al.²⁸
²⁹proposed a Bayesian Neural Network (BNN) method based on the general approxi-²⁹
³⁰mation capability of feedforward neural networks to solve this problem. Compared³⁰
³¹with the deep neural network, each model might be relatively weak, but the entire³¹
³²mixed model could still perform well in data fitting and prediction [6]. They found³²
³³that the sensitivity of cancer cells to drug molecules is driven by the characteristics³³

¹of cells and drugs. Emdadi, A. and Eslahchi, C. proposed a DSPLMF method based¹
²on a recommendation system. The gene expression profile, copy number variations²
³and single nucleotide mutation information were used to calculate the similarity³
⁴of the cell lines, and the chemical structure was used to calculate the similarity of⁴
⁵the drugs. And the possibility of cell lines being sensitive to drugs was calculated⁵
⁶through the logical matrix decomposition to discover the effective characteristics⁶
⁷of the cell lines and drugs [7]. Similarly, Xie et al. used a deep learning model to⁷
⁸predict the response and efficacy of different anticancer drugs to the breast cancer,⁸
⁹and proposed an unsupervised variational autoencoder model geneVAE and recti-⁹
¹⁰fied junction tree variational autoencoder (JTVAE). GeneVAE and JTVAE were¹⁰
¹¹found to have strong robustness in drug response prediction of breast cancer cell¹¹
¹²lines and whole cancer cell lines [8]. Su, Ran et al. used genetic information, chem-¹²
¹³ical characteristics and biological context with the ensemble optimization strategy,¹³
¹⁴and combined with the weighted model META-GDBP to predict drug response,¹⁴
¹⁵which found a high correlation between predicted drug response and observed drug¹⁵
¹⁶response [9]. Sharifi-Noghabi Hossein et al. proposed a deep neural network MOLI¹⁶
¹⁷algorithm, which took somatic mutation, copy number variation and gene expression¹⁷
¹⁸as input data and used a combination of multi-omics methods and clinical data to¹⁸
¹⁹predict drug response. Compared with the latest single-omics and early integrated¹⁹
²⁰multi-omics methods, their proposed method had a significant improvement in pre-²⁰
²¹diction performance [10]. Similarly, Szalai Bence et al. conducted a model prediction²¹
²²analysis based on the correlation between the differentially expressed genes mea-²²
²³sured in the cell lines and the drug sensitivity under the action of the the drug at a²³
²⁴specific concentration, and found that the cell line response was correlated with the²⁴
²⁵drug concentration and time. However, the model achieved low accuracy and poor²⁵
²⁶fitting in the prediction process because it ignored the non-linear characteristics²⁶
²⁷between differentially expressed genes and the drug sensitivity [11].²⁷

²⁸ In this study, we developed the WRFEN-XGBoost algorithm to predict the cell²⁸
²⁹viability under the drug induction using LINCS-L1000 perturbation signatures.²⁹
³⁰Firstly, we screened and matched the three data sets, including perturbation tran-³⁰
³¹scriptomics signatures (LINCS-L1000), cancer treatment response portal (CTRP)³¹
³²and cancer dependence map database (Achilles), and divided them into nine data³²
³³subsets. Secondly, we proposed a weighted fusion algorithm based on random forest³³

¹and elastic nets to effectively extract non-linear features between differentially ex-¹
²pressed genes and cell viability, and completed the selection of key genes. Then, we²
³used the XGBoost algorithm to predict the cell viability and analyzed the apoptosis³
⁴response under the action of drug toxicity and gene silencing. At the same time,⁴
⁵in order to avoid the problem of tedious parameter adjustment, we introduced the⁵
⁶FEBPSO algorithm into the XGBoost learning algorithm. Finally, in order to mea-⁶
⁷sure the feasibility of our method, we completed cross-dataset validation between⁷
⁸compounds and shRNAs at different perturbation times. In addition, we validated⁸
⁹the drug sensitivity inference on the two benchmark data sets of CCLE and NCI60.⁹
10 10

¹¹Methods 11

¹²Dataset collection 12

¹³We used five datasets in this study, including the perturbation transcriptomics¹³
¹⁴signatures (LINCS-L1000), the Cancer Therapeutics Response Portal (CTRP),¹⁴
¹⁵the Cancer Dependence Map Database (Achilles), the Cancer Cell Line Ency-¹⁵
¹⁶clopedia (CCLE) and NCI-60 dataset. LINCS adopted L1000 technology to de-¹⁶
¹⁷tect the transcriptome expression data in human cancer cell lines under vari-¹⁷
¹⁸ous external stimulation. The expression of the whole genome was extrapolated¹⁸
¹⁹by detecting the expression levels of 978 genes [12],[13]. The differentially ex-¹⁹
²⁰pressed signatures corresponding to level five in the LINCS project were cho-²⁰
²¹sen as the training data set, and the data could be obtained from the website²¹
²²<https://www.ncbi.nlm.nih.gov/geo/>. To analyze the cellular response of the can-²²
²³cer cell lines to specific therapeutic drugs, we used the Cancer Treatment Re-²³
²⁴sponse Portal (CTRP), which covered the link between compound sensitivity and²⁴
²⁵genetic or lineage characteristics in 70,000 cancer cell lines. We selected post-quality-²⁵
²⁶control cell viability values as a target for our modeling, which could be downloaded²⁶
²⁷from the website <https://ocg.cancer.gov/programs/ctd2/data-portal> [14]. The third²⁷
²⁸dataset, Cancer Dependence Map Database, could be obtained from the website²⁸
²⁹<https://portals.Broadinstitute.org/achilles> and we selected the log fold scores of²⁹
³⁰effects change before and after shRNA treatment for our model analysis [15]. 30

³¹To verify the effectiveness of our prediction model, we used the NCI-60³¹
³²dataset and the Cancer Cell Line Encyclopedia (CCLE) as validation datasets³²
³³in the end, respectively. The NCI-60 dataset could be downloaded from website³³

¹https://dtp.cancer.gov/discovery_development/nci-60, and we set GI50 value as¹
²the evaluation standard for drug sensitivity [16]. The last dataset was the CCLE²
³dataset, which consisted of the responses of more than 400 cell lines and 24 com-³
⁴pounds at eight concentration points, as well as the expression data of 18,926⁴
⁵genes for each cell line. The CCLE dataset could be downloaded from the website⁵
⁶<https://portals.broadinstitute.org/ccle>, and we used the active area of the drug as⁶
⁷the evaluation standard for drug sensitivity [17].⁷

8

8

⁹Dataset preprocessing⁹

¹⁰We first merged the two-stage perturbation screens LINCS-L1000-PhaseI and¹⁰
¹¹LINCS-L1000-PhaseII, and obtained the genome-wide gene expression levels un-¹¹
¹²der various perturbations in LINCS-L1000. To further analyze the cell viability of¹²
¹³different cell lines under the compound perturbation, we correlated it with the cell¹³
¹⁴viability data after drug treatment in CTRP. We matched the sample instances¹⁴
¹⁵based on the same cell line and the drug identification number provided by the¹⁵
¹⁶Broad Institute. We referred to (1) to match samples with similar concentrations.¹⁶
¹⁷For different experimental batches, we took the average value of the cell viability¹⁷
¹⁸which was measured in the same concentration.¹⁸

19

19

$$\text{doseDiff} = |\log_{10}(Cdose) - \log_{10}(Ldose)| \leq 0.2 \quad (1)$$

21

21

²²where $Cdose$ was the concentration value corresponding to the cancer treatment²²
²³drug in CTRP, and $Ldose$ was the concentration value corresponding to the per-²³
²⁴turbation signatures in LINCS-L1000.²⁴

²⁵In the course of the research, in order to enable our training model to be tested²⁵
²⁶independently on other datasets to verify the effectiveness of the model, we at-²⁶
²⁷tempted to use similar phenotypic information to the cancer treatment response²⁷
²⁸portal CTRP for further research. We associated the merged two-stage LINCS-²⁸
²⁹L1000 perturbation screen data with the Achilles project, the cancer dependency²⁹
³⁰map database, to investigate the effect of single gene knockdown or knockout on³⁰
³¹apoptosis or proliferation of cancer cells under the action of shRNA. Since the num-³¹
³²ber of cell survival after drug treatment or shRNA treatment was proportional to³²
³³the evaluation indicators in the CTRP project or the Achilles project, for simplicity,³³

we referred to the cell phenotypic information in the above two data sets as cell viability. The specific process above was shown in Fig.1.

Model establishment

The research framework of this study is shown in Fig.2. In the first place, we completed the selection of differentially expressed genes and predictive analysis of cell viability on the perturbation transcriptomics signatures LINCS-L1000 and the cancer treatment response portal CTRP dataset. To derive the model's performance across the datasets, we then performed independent screen tests on the cancer dependency map database Achilles (only the test process of the CTRP-L1000 model on the data set Achilles-L1000 is presented here, and vice versa). At the same time, we conducted the model validation based on the active area value in the Cancer Cell Line Encyclopedia CCLE dataset and the drug sensitivity index in the NCI-60 dataset.

Feature extraction based on random forest and elastic net

Random forest, as a typical representative of the Bagging method in ensemble learning, can guarantee the improvement of the regression accuracy and search for a large number of non-linear features [18]. It is considered as one of the most successful algorithms to describe the correlation between key genes and cell phenotype studies [19]. In this study, the sample space is randomly divided into different parts by bootstrapping method. For each node of the decision tree, several genes are randomly selected from the M -dimensional differentially expressed gene space $\mathbf{OriDEGs} = (g_1, g_2, g_3, \dots, g_M)$ and then form the Z -dimensional gene subspace $\mathbf{SubGenes} = (i_1, i_2, i_3, \dots, i_Z)$. Then we select the best split node and get the result of the sample by the weak decision tree. To obtain the final results, prediction of each weak decision tree is averaged. After obtaining the prediction results, we used the Pearson correlation coefficient to evaluate the performance of the random forest to prepare for the feature-weighted fusion. We arranged each attribute in descending order according to the importance of the genes. The non-contributing genes were removed and the number of remaining genes were recorded after sorting.

Elastic network regression, as a combination of ridge regression and lasso regression, can not only reduce the prediction variance but also achieve the purpose of coefficient shrinkage and variable selection [20]. Therefore, we use elastic net regres-

1 sion to select the key genes. We used the Pearson correlation on the validation set¹
 2 to select the appropriate parameter settings for the model. We evaluated the con-²
 3 tribution of each characteristic gene in the model and ranked them in descending³
 4 order of gene contribution. 4

5 In order to screen out effective differentially expressed genes (DEGs), we used⁵
 6 a weighted fusion algorithm of the random forest and elastic network (referred as⁶
 7 WRFEN) to select key genes. 7

$$8 \quad W(DEGs)_{Rank} = \frac{e^{RFPearson} * (DEGs)_{Rank}^{RF} + e^{ENPearson} * (DEGs)_{Rank}^{EN}}{e^{RFPearson} + e^{ENPearson}} \quad (2)^9$$

10 10

11 where $RFPearson$ and $ENPearson$ are the Pearson correlation on the valida-¹¹
 12 tion set using random forest and elastic network algorithms. $(DEGs)_{Rank}^{RF}$ and¹²
 13 $(DEGs)_{Rank}^{EN}$ are the feature importance order of the differentially expressed gene¹³
 14 DEGs and the number of genes selected in the random forest and elastic network¹⁴
 15 algorithm, respectively. 15

16 We ranked the key genes in the random forest and elastic network respectively,¹⁶
 17 and use (2) to perform weighted summation. Finally, we ranked the result and the¹⁷
 18 optimal number of genes in order of gene contribution. The algorithm flowchart¹⁸
 19 was shown in Supplementary Figure 1. More precisely, it was a feature selection¹⁹
 20 method based on the combination of random forest and elastic net. It calculated²⁰
 21 the order of each gene in two methods and the performance of the two methods²¹
 22 in the prediction performance (Pearson correlation) was used as the weight. If the²²
 23 prediction performance of the model was better, the more weight it occupied in²³
 24 gene ranking and the higher the genes in the final ranking. 24

25 25

26 *Cell viability prediction algorithm based on XGBoost and FEBPSO* 26

27 XGBoost is one of the most competitive prediction algorithm in machine learning. 27

28 It improves the integration of the gradient boosting algorithm and has high per-²⁸
 29 formance in solving both classification and regression problems [21]. We used the²⁹
 30 XGBoost algorithm to predict cell viability and obtained a prediction score on the³⁰
 31 leaf node of each decision tree based on the differential expression of genes in each³¹
 32 sample. Multiple weak estimators are constructed one by one through multiple it-³²
 33 erations. The cell viability prediction result is defined as the sum of the prediction 33

1 scores of all the trees as follows. 1

$$2 \quad c\hat{v}_i = \sum_{k=1}^K f_k(\text{sample}_i[\text{DEGs}]) \quad (3) \quad 3$$

4

5 where $f_k(\text{sample}_i[\text{DEGs}])$ represents the prediction score on the k -th decision tree 5
 6 for the i -th sample on the selected differentially expressed gene set DEGs. K is the 6
 7 number of decision trees. Then during the t -th iteration of the sample, the model's 7
 8 predicted value $c\hat{v}_i$ can be described as follows: 8

$$9 \quad c\hat{v}_i^{(t)} = c\hat{v}_i^{(t-1)} + f_t(\text{sample}_i[\text{DEGs}]) \quad (4) \quad 10$$

11

12 In this study, in order to improve the prediction accuracy of cell viability and 12
 13 reduce the prediction bias, we used the discrete binary particle swarm optimization 13
 14 with flexible weights algorithm FEBPSO to adaptively adjust the parameters of XG- 14
 15 Boost. As a typical representative of swarm intelligence algorithms, particle swarm 15
 16 optimization can effectively solve nonlinear continuous optimization problems [22]. 16
 17 Meanwhile, it solves the problem of too long training time due to a large amount of 17
 18 adjustment parameters [23]. In the prediction process of FEBPSO-XGBoost, we first 18
 19 initialized the binary particle swarm, encoded each parameter as a binary number 19
 20 and transformed the parameter optimization into a discrete combinatorial optimiza- 20
 21 tion problem. During each iteration, the parameters were converted into decimal 21
 22 numbers within the specified range in a group of six. At this time, we calculated 22
 23 the Pearson correlation coefficient of each individual particle running in XGBoost 23
 24 algorithm and evaluated the fitness of each individual particle. For each particle, 24
 25 we compared the current fitness value with the individual's historical best position 25
 26 or global best position. If the current fitness value was higher, the historical best 26
 27 position and global best position would be updated with the current position of the 27
 28 particle. At the same time, the particle speed and position information would be 28
 29 updated to enter the next iteration until the termination condition has been met. 29
 30 Finally, the global optimal value and the best parameter settings would be output 30
 31 at this time. The particle speed is updated as follows: 31

$$32 \quad v_i^{k+1} = wv_i^k + c_1r_1(x_{pbest,i}^k - x_i^k) + c_2r_2(x_{gbest}^k - x_i^k) \quad (5) \quad 33$$

¹where v_i^k represents the velocity vector of particle i during the k -th iteration, x_i^k
²represents the position vector of particle i during the k -th iteration, c_1 and c_2 are
³the acceleration constant, r_1, r_2 are the random number, w is the inertial weight,
⁴ $x_{pbest,i}^k$ denotes the best position of the individual particle and x_{gbest}^k denotes the
⁵best position of the global particle.

⁶In order to overcome the shortcomings of premature convergence and falling into
⁷local extremes of particle swarm optimization, we used the formula shown below to
⁸update the weights [24].

$$w(k) = \alpha_1 e^{-\frac{\psi * k}{T}} + \alpha_2 e^{\frac{\psi * k}{T}} \quad (6)$$

¹¹where $\alpha_1 = \frac{w_2 e^{\psi} - w_1 e^{2\psi}}{1 - e^{2\psi}}$, $\alpha_2 = \frac{w_1 - w_2 e^{\psi}}{1 - e^{2\psi}}$, T denotes the maximum number of itera-
¹²tions, k is the current number of iterations, w_1, w_2 are the minimum inertia weight
¹³and maximum inertia weight greater than zero, respectively.

¹⁵We used WRFEN for core gene selection and FEBPSO-XGBoost for predictive
¹⁶analysis. Through this, we formed a complete prediction model and explained the
¹⁷complete apoptotic levels observed in cell lines with specific drugs and concentra-
¹⁸tions.

Results

²⁰Based on the latest transcriptomic perturbation screens in LINCS-L1000, we con-
²¹ducted the study with the cell viability after the drug treatment in CTRP and
²²the effect change score before and after the treatment with shRNA in the Achilles
²³project, respectively. From the perspective of gene regulation, we examined the
²⁴relationship between key genes and drug response. At the same time, the FEBPSO-
²⁵XGBoost machine learning algorithm was used to predict the cell viability of differ-
²⁶ent cell lines with the treatment of various drugs or shRNA by using the expression
²⁷levels of characteristic genes under the action of different perturbation times and
²⁸different drug concentrations.

Analysis of feature selection

³¹In the feature selection process, we firstly selected 40 trees for the establishment of
³²a random forest, and the results were ranked according to the variable contribution.

³³Secondly, the ratio of the lasso penalty term was set to 0.1, 0.2, 0.5, 0.7, 0.95, 1 and

¹the coefficient penalty term was controlled to from 0.1 to 1.0 by step 0.1 in the¹
²elastic net. The best combination of the parameters was decided on the validation²
³set. Then, we sorted the variables according to their contribution and deleted the³
⁴non-contributing genes. Finally, we calculated the selected characteristic genes ac-⁴
⁵cording to Formula (2), and obtained the final genes. The feature genes selected⁵
⁶on each subset (subset names were shown in Supplementary Table 1) was ranked⁶
⁷according to their contribution. We listed the number of feature genes selected and⁷
⁸the contribution ranking of the fifteen key genes in each subset in Supplementary⁸
⁹Table 1. 9

¹⁰ Taking the LINCS-L1000-CTRP-24h dataset as an example, we compared the¹⁰
¹¹WRFEN with the existing traditional methods FTest [25], MI [26], RFFS [27] and¹¹
¹²LRFS [28], and tested it on multiple predictors at the same time (Supplementary¹²
¹³Figure 2). The results showed that the results of the gene selection algorithm in¹³
¹⁴this paper were better than the existing single algorithms. It could also be observed¹⁴
¹⁵that the prediction performance of the model would be gradually stabilize as the¹⁵
¹⁶number of selected feature genes increases. 16

¹⁷ In order to further understand the biological functions performed by the selected¹⁷
¹⁸characteristic genes, we took the subsets of CTRP-L1000-24h and Achilles-L1000-¹⁸
¹⁹96h as examples to perform analysis on the extracted characteristic genes. We could¹⁹
²⁰find that they were all closely related to the apoptotic process from Fig. 3 and²⁰
²¹Supplementary Figure 3. The most significantly enriched pathways, r-has-1640170²¹
²²and GO:0007346, were involved in the regulation of cell cycle and apoptosis, which²²
²³also confirmed that the differentially expressed genes selected in this study after²³
²⁴treatment with drugs or shRNA constituted the pathway of apoptosis. 24

²⁵ 25

²⁶Prediction and analysis of drug induced cell viability 26

²⁷ We updated and adjusted the parameter combination of XGBoost with the binary²⁷
²⁸discrete particle swarm optimization with flexible weight. We set the number of²⁸
²⁹swarm particles to be 25, the dimension of the particles to be 48, the maximum²⁹
³⁰number of iterations in CTRP-L1000 and Achilles-L1000 series models to be 50 and³⁰
³¹20 respectively, the acceleration constants to be 1.5, the maximum and minimum³¹
³²values of inertia weight to be 0.8 and 0.4 respectively, the maximum and minimum³²
³³values of velocity to be 10 and -10 respectively and weight updating formula of³³

¹parameter ψ to be 2.6. The correlation coefficient between the observed value and¹
²the predicted value was used as the model evaluation index and the fitness function.²
³At the beginning of the particle swarm optimization algorithm, the population was³
⁴generated randomly. When the iteration reached a certain number, the optimal⁴
⁵solution or approximate optimal solution would be found with a high probability.⁵
⁶The experimental results of parameter optimization in XGBoost by using FEBPSO⁶
⁷algorithm were shown in Fig. 4. 7

⁸8
⁹9
¹⁰ From the above experimental results, it was obvious that the measurement of¹⁰
¹¹cell viability in CTRP required a long perturbation time. With the increasement of¹¹
¹²the perturbation time, the reliability of the forecast also continued to rise, and the¹²
¹³prediction results of the 24-hour perturbation time was more reliable. When the con-¹³
¹⁴centration factor was added in the prediction of the CTRP dataset, the prediction¹⁴
¹⁵accuracy of the model could be improved, which indicated that the cell viability¹⁵
¹⁶depended on the concentration of the drug to some extent. In the LINCS-L1000¹⁶
¹⁷perturbation screens and cancer dependency map database Achilles, the model pro-¹⁷
¹⁸duced by the 96-hours perturbation time had the most significant prediction effect.¹⁸
¹⁹It could be seen from the results that the disturbance time was not necessarily as¹⁹
²⁰long as possible. 20

²¹21
²² In the optimization process of the CTRP-L1000 series model, when the number of²²
²³iterations reached about 20 rounds, the prediction performance of the model gradu-²³
²⁴ally tended to be stable. In the process of Achilles-L1000 series model optimization,²⁴
²⁵when the number of iterations reached about 8 rounds, the prediction performance²⁵
²⁶of the model also gradually tended to be stable. After we used FEBPSO to adjust²⁶
²⁷the parameters of the XGBoost model, the optimal parameter combinations and²⁷
²⁸default values of each parameter were shown in Table 1 and Supplementary Table2²⁸
²⁹below. It could be seen that this experiment fully proves the effectiveness of the²⁹
³⁰parameter optimization algorithm proposed by this research. Compared with the³⁰
³¹traditional default parameters, using the FEBPSO algorithm to optimize the pa-³¹
³²rameters of the XGBoost model had significantly improved the accuracy of model³²
³³prediction. 33

¹Independent dataset validation on CTRP-L1000 and Achilles-L1000 1

²In order to verify the reliability of the model predictions, we used independent²
³datasets to verify the model's prediction capabilities. We had implemented the in-³
⁴teractive test in the CTRP-L1000 series model and the Achilles-L1000 series model.⁴
⁵The Fig. 5 showed the experimental results. From the figure above, it could be found⁵
⁶that the 24-hour perturbation time was the best in the CTRP-L1000 data set. The⁶
⁷Pearson correlation of the model on this data set was 0.8321, which was better⁷
⁸than the 3-hour and 6-hour perturbation times. In the Achilles-L1000 dataset, the⁸
⁹96-hour perturbation time was considered to be the best. The performance of the⁹
¹⁰model on this data set is better than the perturbation time of 120 hours and 144¹⁰
¹¹hours with 0.5893 Pearson correlation. Similarly, in terms of independent set valida-¹¹
¹²tion, the CTRP-L1000-6h model, CTRP-L1000-24h model and Achilles-L1000-96h¹²
¹³model was superior to other models in CTRP-L1000-24h screen with 0.7416, 0.8321¹³
¹⁴and 0.7319 Pearson correlation, respectively. Therefore, we further confirmed that¹⁴
¹⁵the drug could achieve excellent predictive performance after a longer perturbation¹⁵
¹⁶time. 16

17

17

¹⁸Model validation on the NCI60 dataset 18

¹⁹In order to validate the model across the NCI60 dataset, we used the GI50 value¹⁹
²⁰as the indicator of drug sensitivity evaluation and binarized the GI50 value (50%²⁰
²¹growth inhibition). In the NCI60 dataset, when the efficacy was within the range of²¹
²²50% growth inhibition concentration, it corresponded to the GI50 value in the drug²²
²³sensitivity evaluation index. When the efficacy was not effective within the 50%²³
²⁴growth inhibition concentration range, it was recorded as the highest concentration²⁴
²⁵value. In this study, we would define the drug concentration difference variable,²⁵
²⁶which portrayed the efficacy of the drugs and was calculated as shown in Formula²⁶
²⁷(7). In other words, when the value of the drug concentration difference was less than²⁷
²⁸zero, it meant that the drug was an effective drug, otherwise it was an ineffective²⁸
²⁹drug. 29

30

30

$$\Delta drug_conc(dr, cl) = drug_sens(dr, cl) - test_max_conc(dr, cl) \quad (7) \quad 31$$

³²where, $\Delta drug_conc(dr, cl)$ was the difference in drug concentration when the cell³²
³³lines cl under the treatment of the specific drug dr . $drug_sens(dr, cl)$ was the drug³³

¹sensitivity value GI50 for cl treated by dr . $test_max_conc(dr, cl)$ was the maximum¹
²tested drug concentration used in the treatment of cell line cl with the drug dr .²

³ In this study, ROC curve and PR curve were used to measure the contribution³
⁴of the algorithm in evaluating the drug effectiveness. By observing the ROC curve⁴
⁵shown in Fig. 6(a), we could find that the prediction made by the Achilles-L1000-96h⁵
⁶model is the most accurate in the LINCS-L1000-NCI60-24h dataset. When using⁶
⁷this model for prediction, the AUC area under the ROC curve reached 0.80, the 95%⁷
⁸confidence interval ranged from 0.769 to 0.822, and the significance level was less⁸
⁹than 0.0001. The other two models also had good performance. Among them, the⁹
¹⁰AUC area under the ROC curve of the CTRP-L1000-24h model reached 0.76, and¹⁰
¹¹the area under the ROC curve of the CTRP-L1000-6h model reached 0.74. In the¹¹
¹²accuracy-recall evaluation curve shown in Fig. 6(b), the Achilles-L1000-96h model¹²
¹³still surpassed other models with the area under the curve $AUC = 0.94$. Through¹³
¹⁴the above analysis, we further confirmed that the Achilles-L1000-96h model was¹⁴
¹⁵effective during the prediction process of the LINCS-L1000-NCI60-24h data set,¹⁵
¹⁶and it could be further used for the effectiveness testing of other drugs.¹⁶

¹⁷ Furthermore, we also matched and correlated the LINCS-L1000 perturbation¹⁷
¹⁸screens, CTRP data and NCI60 data according to the matching method described¹⁸
¹⁹above. The drug with the perturbation time of 24 hours was recorded as LINCS-¹⁹
²⁰L1000-CTRP-NCI60-24h. In this experiment, we used CTRP-L1000-6h, CTRP-²⁰
²¹L1000-24h and Achilles-L1000-96h models to predict the cell viability in three major²¹
²²data sets, which had drugs and cell lines in common. We also binarized the drug²²
²³sensitivity data in NCI60.²³

²⁴ Finally, we used the ROC curve and PR curve to discuss and analyze the exper-²⁴
²⁵imental results. As shown in Supplementary Figure 4, when we used the Achilles-²⁵
²⁶L1000-96h model, the CTRP-L1000-24h model and the CTRP-L1000-6h model to²⁶
²⁷predict the effectiveness of the drug, the area under the ROC curve achieved 0.78,²⁷
²⁸0.80 and 0.72, respectively, and the area under the PR curve achieved 0.98, 0.98 and²⁸
²⁹0.97, respectively. The above results indicated the superior prediction performance²⁹
³⁰of the Achilles-L1000-96h model and the CTRP-L1000-24h model.³⁰

³¹ While predicting the effectiveness of the drugs, we required that the predictors³¹
³²used in this study could make effective predictions. In addition, whether the appro-³²
³³priate features could be selected during the feature selection stage directly affected³³

¹the predictive performance of the predictors. To do this, we correlated the selected¹
²feature genes with the effectiveness of the drug. We observed whether the differential²
³expression levels of selected characteristic genes have significantly different expres-³
⁴sion patterns under the action of effective or ineffective drugs. For this reason, we⁴
⁵mapped the differential expression levels of the first 15 differentially expressed genes⁵
⁶selected in the feature selection stage under the treatment of effective drugs and⁶
⁷ineffective drugs. Fig. 7(a) was the result of the LINCS-L1000-NCI60-24h dataset⁷
⁸and Fig. 7(b) was the result of the LINCS-L1000-CTRP-NCI60-24h dataset. By⁸
⁹comparison, in the effective drug group, we could find that the expression level⁹
¹⁰of differentially expressed genes had significantly up-regulated or down-regulated.¹⁰
¹¹However, in the ineffective drug group, there was no significant change in the ex-¹¹
¹²pression level of differentially expressed genes. Therefore, we further demonstrated¹²
¹³the validity of selected feature genes. 13

¹⁴ So far, we had completed inferring the effectiveness of the drug from the predicted¹⁴
¹⁵cell viability of each model. To further examine whether there was a significant¹⁵
¹⁶difference between the effective and ineffective drugs on the cell viability, we used a¹⁶
¹⁷non-parametric Mann Whitney test to analyze the cell viability prediction results,¹⁷
¹⁸as shown in Fig. 8. Different models were predicted on LINCS-L1000-NCI60-24h¹⁸
¹⁹screen and LINCS-L1000-CTRP-NCI60-24h screen respectively. The results found¹⁹
²⁰that using the Achilles-L1000-96h model to discriminate between effective drugs²⁰
²¹and ineffective drugs had a significant difference in the mean value, the significance²¹
²²levels were $P \leq 0.0001$ and $P = 0.0004$, respectively. In addition, similar results²²
²³were obtained in the use of CTRP-L1000-24h model for inferring drug effectiveness,²³
²⁴the significance levels were $P \leq 0.0001$ and $P = 0.0002$, respectively. 24

²⁵

25

²⁶Model validation on the CCLE dataset 26

²⁷Our model was also verified on CCLE, and we used the active area as the evaluation 27
²⁸criterion of drug sensitivity. In order to achieve binarization of drug sensitivity on 28
²⁹the CCLE data set, we first normalized the active area in CCLE to zero mean. 29
³⁰Meanwhile, we defined the active area with 0.8 variance above the mean as an ef- 30
³¹fective drug, and the active area with 0.8 variance below the mean as an ineffective 31
³²drug. We then searched for common combination pairs of cell lines and drugs in the 32
³³LINCS-L1000 perturbation screen. Since there were only a small number of 24 drugs 33

¹in the CCLE data set, we used the PubChem database to find synonymous drugs.¹
²We marked the data after matching as LINCS-L1000-CCLE. Similarly, we screened²
³the drugs corresponding to the perturbation time of 24 hours, which were included³
⁴in LINCS-L1000-CCLE-24h. At the same time, we selected the drugs whose concen-⁴
⁵tration was greater than or equal to 10 micromoles. In addition, when multiple drug⁵
⁶perturbation signatures were presented, we choose the lowest cell viability value. ⁶
⁷ We used the ROC curve shown in Supplementary Figure 5(a) and the PR curve⁷
⁸shown in Supplementary Figure 5(b) to measure the results of the algorithm. By⁸
⁹observing the experimental results, we found that when we used the drug sensitivity⁹
¹⁰data in CCLE to evaluate the predicted cell viability values, the Achilles-L1000-96h¹⁰
¹¹model also showed excellent performance in cross-dataset validation. When we used¹¹
¹²Achilles-L1000-96h model to predict the effectiveness of the drug, the area under¹²
¹³the ROC curve achieved 0.84 and the area under the PR curve achieved 0.88. The¹³
¹⁴differential expression on effective and ineffective drugs was shown in Supplementary¹⁴
¹⁵Figure 6. We could see that the LINCS-L1000-CCLE-24h dataset still showed the¹⁵
¹⁶same gene expression pattern as the LINCS-L1000-NCI60-24h dataset. That was to¹⁶
¹⁷say, the differentially expressed genes in the effective drug group were significantly¹⁷
¹⁸up-regulated and down-regulated. ¹⁸

¹⁹

¹⁹

²⁰**Discussion** ²⁰

²¹In order to evaluate the effectiveness of the algorithm in this paper, we analyzed²¹
²²and compared our algorithm with other existing methods including PCA-Lasso,²²
²³PCA-SVR, FTest-RF, MI-KNN, VAE [8] and DAE-NN [29]. The Principal Compo-²³
²⁴nents Analysis (PCA), Ftest and Mutual Information (MI) were used to extract the²⁴
²⁵features, and the Lasso, Support Vector Regression (SVR), Random Forest (RF)²⁵
²⁶and k-nearest neighbor (KNN) were used for the final prediction. VAE and DAE-²⁶
²⁷NN are proposed by the recent literature in drug response prediction. VAE used²⁷
²⁸the variational autoencoder to predict the response of different anti-cancer drugs.²⁸
²⁹DAE-NN used a deep autoencoder to extract the features and the neural network²⁹
³⁰was for the final prediction. ³⁰

³¹ In the present paper, we used the Pearson correlation coefficient, coefficient of³¹
³²determination (R^2) and mean squared error of the predicted and actual values to³²
³³measure the prediction performance of the model. In the training process of VAE³³

¹and DAE-NN algorithms, we used grid search to select the best training parameters¹
²for the learning rate [0.001, 0.005, 0.01, 0.05, 0.1] and iteration period [30, 90,²
³150, 220, 300]. The detailed experimental results of these seven algorithms were³
⁴shown in Supplementary Tables 3-5. Taking the CTRP-L1000-24h(S1) dataset as an⁴
⁵example, the predicted results were shown in Table 2. Our algorithm outperformed⁵
⁶other algorithms with the maximum correlation coefficient 0.8321, the maximum⁶
⁷coefficient of determination 0.6922 and the minimum mean squared error 0.025. ⁷

⁸ ⁸
⁹ Compared with PCA-Lasso, PCA-SVR, FTest-RF, MI-KNN, VAE and DAE-NN⁹
¹⁰algorithms, Pearson correlation coefficient of our method increased by 5.50%, 5.33%,¹⁰
¹¹4.77%, 3.32%, 0.39%, 3.59% and R^2 increased by 11.45%, 11.45%, 9.80%, 7.92%,¹¹
¹²1.45% and 12.12%. In terms of the mean squared error, our method decreased from¹²
¹³3.85% to 21.88% comparing with the other six algorithms above. The experimental¹³
¹⁴results showed that the prediction performance of the proposed algorithm have¹⁴
¹⁵been further improved. For the CTRP-1000-3h, CTRP-L1000-6h, CTRP-L1000-¹⁵
¹⁶24h, Achilles-L1000-96h, Achilles-L1000-120h and Achilles-L1000-144h datasets ,¹⁶
¹⁷the evaluation results of other models were shown in Supplementary Tables 3-5. ¹⁷

¹⁸ ¹⁸
¹⁹ In addition to reliably and effectively inferring cell viability through the pre-¹⁹
²⁰dictive models, we also needed to correlate our results with the literature on cell²⁰
²¹viability, as shown in Fig. 9. As a member of tumor necrosis factor receptor super-²¹
²²family, high affinity nerve growth factor receptor p75NTR could induce apoptosis²²
²³and inhibit the growth of prostate epithelial cells. Azacitidine-mediated p75NTR²³
²⁴had anti-tumor effects on androgen-independent prostate cancer cells 22Rv1 and²⁴
²⁵PC3 [30]. After Bortezomib treatment, the cells with suppressed C/EBPbeta levels²⁵
²⁶showed delayed autophagy activation. The growth of the PC3 cells and xenografts²⁶
²⁷has been decreased with the C/EBbeta gene knockdown, which could make PC3²⁷
²⁸cells sensitive to Bortezomib [31]. Another study has tested the effects on three²⁸
²⁹related human glioma cell lines treated by the new epidermal growth factor recep-²⁹
³⁰tor (EGFR) tyrosine kinase Tyrphostin-AG-1478, and found that AG-1478 was the³⁰
³¹relatively specific inhibitor of truncated EGFR. They had important medical sig-³¹
³²nificance because the truncated EGFR occurred frequently in glioblastoma, breast,³²
³³lung and ovarian cancer [32]. ³³

1 Conclusions 1

2 In this paper, we managed to predict the drug-induced cell viability from the differ-2
 3 ential gene expression data through the WRFEN-XGBoost algorithm. The study3
 4 of cell phenotype was firstly correlated with the drugs and shRNA perturbation4
 5 signatures. In addition, we have completed the selection of key genes based on5
 6 the WRFEN algorithm and proposed a novel FEBPSO-XGBoost machine learning6
 7 method to predict the cell viability. Through the connection between cell viabil-7
 8 ity and pharmacogenomics, the establishment of the prediction model trained from8
 9 perturbation transcriptomics signatures, cell phenotype and drug response data has9
 10 been completed. At the same time, the robustness and effectiveness of our proposed10
 11 modeling strategy in drug sensitivity analysis were verified on CCLE and NCI-11
 12 60 datasets. This study could provide help for the biomedical researchers in drug12
 13 screening and promote the analysis of anticancer drugs in pharmacogenomics. 13

14 However, in the clinical application of cancer cell lines and anticancer therapies,14
 15 it is urgent to identify the biomarkers that can distinguish between drug-sensitive15
 16 cell lines and drug-resistant cell lines. Firstly, besides gene expression, drug char-16
 17 acteristics can be integrated into the model to achieve better accuracy. Secondly,17
 18 a more appropriate supervised machine learning algorithm is hoped to be designed18
 19 to reveal the sensitivity between cancer cell lines and drug treatment. Finally, we19
 20 will continue to reveal new biomarkers that are sensitive and resistant to the can-20
 21 certherapies. It provides more opportunities for exploring the biological behavior21
 22 of cancer cell lines at the cellular level, and it is also the direction of our future22
 23 research. 23

25 Abbreviations 25

26 LINC: Library of integrated network-based cellular signatures; GEO: Gene Expression Omnibus; CTRP: Cancer
 27 treatment response portal; NCI: National cancer institute; CCLE: Cancer Cell Line Encyclopedia; WRFEN: Weighted
 28 fusion algorithm of the random forest and elastic network; FEBPSO: Binary particle swarm optimization with
 29 flexible inertia weight; XGBoost: Extreme Gradient Boosting; PCA: Principal Components Analysis; MI: Mutual
 30 Information; SVR: Support Vector Regression; RF: Random Forest; KNN: K-nearest neighbor; VAE: Variational
 31 Autoencoder; DAE: Deep autoencoder; NN: Neural network; ROC curve: Receiver operating characteristic curve;
 32 PR curve: Precision-Recall curve; AUC: Area Under Curve 30

31 Declarations 31

32

33 Acknowledgements 33

Not applicable.

| | | |
|-----------|---|-----------|
| 1 | Funding | 1 |
| 2 | This work was supported in part by Shanghai Science and Technology Innovation Plan Project (20dz1203800 to | 2 |
| | M.C.), the National Natural Science Foundation of China (61702325 to Y.Q.), National Key R&D Program Projects | |
| 3 | (2018YFD0701003 to M.C.) and Shanghai Science and Technology Innovation Action Plan (16391902900 to M.C.). ³ | 3 |
| 4 | Availability of data and materials | 4 |
| 5 | The screening data (LINCS-L1000,CTRP,Achilles,NCI60 and CCLE) used to train the machine learning model | 5 |
| | presented in this study are available at https://www.ncbi.nlm.nih.gov/geo/ , | |
| 6 | https://ocg.cancer.gov/programs/ctd2/data-portal , https://portals.broadinstitute.org/achilles , | 6 |
| 7 | https://dtp.cancer.gov/discovery-development/nci-60 and https://portals.broadinstitute.org/ccle , respectively. | 7 |
| | Source code is available at https://github.com/RuyiMz/SJZY.git . The results of the data analysis could be | |
| 8 | obtained from the additional files. | 8 |
| 9 | Author's contributions | 9 |
| 10 | JL conducted the experiments, performed the data analysis and wrote the paper. YQ designed the study and | 10 |
| | supervised the research. MC contributed critical review. All authors have read and approved of the final manuscript. | |
| 11 | | 11 |
| | Ethics approval and consent to participate | |
| 12 | Not applicable. | 12 |
| 13 | Consent for publication | 13 |
| 14 | Not applicable. | 14 |
| 15 | Competing interests | 15 |
| 16 | The authors declare that they have no competing interests. | 16 |
| 17 | References | 17 |
| 18 | 1. Samane M, Hossein K, Nafiseh E, Nahid E, Ilnaz R, Abbas R, et al. Producing Soluble Human Programmed | 18 |
| 19 | Cell Death Protein-1: A Natural Supporter for CD4+T Cell Cytotoxicity and Tumor Cells Apoptosis. Iranian | 19 |
| | journal of biotechnology. 2019;17(4):266–267. | |
| 20 | 2. Cubillos-Ruiz JR, Mohamed E, Rodriguez PC. Unfolding anti-tumor immunity: ER stress responses sculpt | 20 |
| | tolerogenic myeloid cells in cancer. J Immunother Cancer. 2017;5:5. | |
| 21 | 3. Mostaghimi H. Quantitative determination of tumor platinum concentration of patients with advanced Breast, | 21 |
| | lung, prostate, or colorectal cancers undergone platinum-based chemotherapy. J Cancer Res Ther. | |
| 22 | 2017;13(6):930–935. | 22 |
| 23 | 4. Yousefi MR, Datta A, Dougherty ER. Optimal Intervention in Markovian Gene Regulatory Networks With | 23 |
| | Random-Length Therapeutic Response to Antitumor Drug. IEEE transactions on bio-medical engineering. | |
| 24 | 2013;60(12):3542–3552. | 24 |
| 25 | 5. Menden MP, Iorio F, Garnett M, McDermott U, Benes CH, Ballester PJ, et al. Machine learning prediction of | 25 |
| | cancer cell sensitivity to drugs based on genomic and chemical properties. PLoS One. 2013;8(4):e61318. | |
| 26 | 6. Yongcui W, Jianwen F, Shilong C. Inferences of drug responses in cancer cells from cancer genomic features | 26 |
| | and compound chemical and therapeutic properties. Scientific reports. 2016;6(6):32679. | |
| 27 | 7. Akram E, Changiz E. DSPLMF: A Method for Cancer Drug Sensitivity Prediction Using a Novel Regularization | 27 |
| | Approach in Logistic Matrix Factorization. Frontiers in genetics. 2020;11. | |
| 28 | 8. Xie J, Dong H, Jing Z, Ren D. Variational Autoencoder for Anti-Cancer Drug Response Prediction. | 28 |
| 29 | Bioinformatics. 2020;Preprint at https://arxiv.org/abs/2008.09763?context=cs.LG (2020). | 29 |
| 30 | 9. Ran S, Xinyi L, Guobao X, Leyi W. Meta-GDBP: a high-level stacked regression model to improve anticancer | 30 |
| | drug response prediction. Briefings in Bioinformatics. 2020;21(3):996–1005. | |
| 31 | 10. Hossein SN, Olga Z, C CC, Martin E. MOLl: multi-omics late integration with deep neural networks for drug | 31 |
| | response prediction. Bioinformatics. 2019;35(14):i501–i509. | |
| 32 | 11. Szalai B, Subramanian V, Holland CH, Alfoldi R, Puskas LG, Saez-Rodriguez J. Signatures of cell death and | 32 |
| 33 | proliferation in perturbation transcriptomics data-from confounding factor to effective prediction. Nucleic Acids | 33 |
| | Res. 2019;47(19):10010–10026. | |

- 1¹². Subramanian A, Narayan R, Corsello SM, Peck DD, Natoli TE, Lu X, et al. A Next Generation Connectivity 1
 2 Map: L1000 Platform and the First 1,000,000 Profiles. *Cell*. 2017;171(6):1437–1452 e17. 2
13. Liu C, Su J, Yang F, Wei K, Ma J, Zhou X. Compound signature detection on LINCS L1000 big data. *Mol*
 3 *Biosyst*. 2015;11(3):714–722. 3
- 4¹⁴. Seashore-Ludlow B, Rees MG, Cheah JH, Cokol M, Price EV, Coletti ME, et al. Harnessing Connectivity in a 4
 Large-Scale Small-Molecule Sensitivity Dataset. *Cancer Discovery*. 2015;5(11):1210–1223. 4
- 5¹⁵. Tsherniak A, Vazquez F, Montgomery PG, Weir BA, Kryukov G, Cowley GS, et al. Defining a Cancer 5
 Dependency Map. *Cell*. 2017;170(3):564–576 e16. 5
- 6¹⁶. Shoemaker RH. The NCI60 human tumour cell line anticancer drug screen. *Nature Reviews Cancer*. 6
 2006;6(10):813–823. 6
- 7¹⁷. Ghandi M, Huang FW, Jane-Valbuena J, Kryukov GV, Lo CC, McDonald ER, et al. Next-generation 7
 8 characterization of the Cancer Cell Line Encyclopedia. *Nature*. 2019;569(7757):503–508. 8
- 9¹⁸. Qi Y. Random Forest for Bioinformatics. *Ensemble Machine Learning*. 2012;p. 307–323. 9
- 10¹⁹. Rahman R, Haider S, Ghosh S, Pal R. Design of Probabilistic Random Forests with Applications to Anticancer 10
 Drug Sensitivity Prediction. *Cancer informatics*. 2016;15(Suppl. 5):57–73. 10
- 11²⁰. Soomro BN, Xiao L, Huang L, Soomro SH, Molaei M. Bilayer Elastic Net Regression Model for Supervised 11
 Spectral-Spatial Hyperspectral Image Classification. *IEEE Journal of Selected Topics in Applied Earth*
 12 *Observations and Remote Sensing*. 2017;9(9):4102–4116. 12
- 13²¹. Li W, Yin Y, Quan X, Zhang H. Gene Expression Value Prediction Based on XGBoost Algorithm. *Frontiers in*
 13 *Genetics*. 2019;10:1077–1077. 13
- 14²². Gong YJ, Li JJ, Zhou Y, Li Y, Chung SH, Shi YH, et al. Genetic Learning Particle Swarm Optimization. *IEEE* 14
Transactions on Cybernetics. 2017;46(10):2277–2290. 14
- 15²³. Mizuho N, Mitsuo N, Osamu S, Ryosuke K, Masahiro Y, Tomohiro K, et al. Computer-aided diagnosis of lung 15
 nodule using gradient tree boosting and Bayesian optimization. *Plos One*. 2018;13(4):e0195875. 16
- 16²⁴. Javad AM, Mousa S, Hossein SM, Deng Y. A Novel Flexible Inertia Weight Particle Swarm Optimization 16
 Algorithm. *Plos One*. 2016;11(8):e0161558. 17
- 17²⁵. Dhanya R, Paul IR, Akula SS, Madhumathi Sivakumar JJN. F-test feature selection in Stacking ensemble 17
 model for breast cancer prediction. *Procedia Computer Science*. 2020;171:1561–1570. 18
- 18²⁶. B B, M CFE. Feature selection, mutual information, and the classification of high-dimensional patterns. 18
Pattern Analysis and Applications. 2008;11(3-4):309–319. 19
- 19²⁷. S MR, E WR. A method for simultaneous variable selection and outlier identification in linear regression. 19
Computational Statistics and Data Analysis. 1996;22(12):251–270. 20
- 20²⁸. M S, V AK. Feature selection and classification of leukocytes using random forest. *Medical and Biological*
 21 *Engineering and Computing*. 2014;52(12):1041–1052. 21
- 21²⁹. Amarbayasgalan T, Lee JY, Kim KR, Ryu KH, editors. Deep Autoencoder Based Neural Networks for Coronary 21
 Heart Disease Risk Prediction. Los Angeles: Plos One; 2019. 22
- 22³⁰. Marampon F, Sanita P, Mancini A, Colapietro A, Scarsella L, Jitariuc A, et al. Increased expression and activity 22
 of p75NTR are crucial events in azacitidine-induced cell death in prostate cancer. *Oncology Reports*. 23
 2016;36(1):125–130. 23
- 23³¹. Barakat DJ, Mendonca J, Barberi T, Zhang J, Kachhap SK, Paz-Priel I, et al. C/EBP beta regulates sensitivity 23
 to bortezomib in prostate cancer cells by inducing REDD1 and autophagosome-lysosome fusion. *Cancer letters*. 24
 2016;375(1):152–161. 24
- 24³². Han YC, Caday CG, Nanda A, Cavenee WK, Huang HJS. Tyrphostin AG 1478 Preferentially Inhibits Human 24
 Glioma Cells Expressing Truncated Rather than Wild-Type Epidermal Growth Factor Receptors. *Cancer*
 25 *Research*. 1996;56(17):3859–3861. 25

30 Figure Legends 30

31

32

33

Figure 1 LINCS-L1000 and CTRP, Achilles data association diagram. The process of data association consisted of two parts: perturbation signatures and cell phenotypic information. The LINCS-L1000-PhaseI and LINCS-L1000-PhaseII were combined and renamed LINCS-L1000. The compound perturbation signatures and shRNA perturbation signatures involved in LINCS-L1000 were respectively associated with CTRP and Achilles datasets according to relevant conditions, which were named CTRP-L1000 and Achilles-L1000. The datasets were divided into CTRP-L1000-3h, CTRP-L1000-6h, CTRP-L1000-24h, Achilles-L1000-96h, Achilles-L1000-120h and Achilles-L1000-144h according to different perturbation time. CTRP-L1000-3h, CTRP-L1000-6h, and CTRP-L1000-24h were divided into six subsets according to the concentration factor was considered(S2) or not considered(S1).

Figure 2 Framework diagram of cell viability prediction.

Figure 3 Enrichment analysis of differentially expressed genes in the CTRP-L1000-24h dataset.

Figure 4 Iterative process of FEBPSO in XGBoost algorithm. a, CTRP-L1000 Optimization. b, Achilles-L1000 Optimization.

Figure 5 Independent dataset validation. Using the Achilles-L1000 series model to predict cell viability in CTRP-L1000 data and vice versa.

Figure 6 ROC curve and PR curve of the model evaluation on LINCS-L1000-NCI60-24h dataset. a, The graph of Receiver Operating Characteristic. b, The graph of Precision-Recall.

Figure 7 Heat map of the first fifteen genes. a, LINCS-L1000-NCI60-24h. b, LINCS-L1000-CTRP-NCI60-24h.

Figure 8 Box plot. comparison of the effective drug group and the ineffective drug group.

Figure 9 The predicted cell viability for different drugs and cell lines. (a-c) showed the cell viability of the drugs Vorinostat, Bardoxolone-methyl and Tyrphostin-AG-1478 in different cell lines. (d-f) showed the cell viability of the cell lines HUES3, MCF7, PC3 in different drugs.

Tables

28
29
30
31
32
33

Table 1 XGBoost parameters and best parameter combinations (CTRP-L1000 Series Model).

| Parameter Name | L1000-CTRP-3h (S1/S2) | L1000-CTRP-6h (S1/S2) | L1000-CTRP-24h (S1/S2) |
|------------------|--------------------------|--------------------------|---------------------------|
| learning rate | 0.0225/0.0476 | 0.01/0.0225 | 0.01/0.035 |
| gamma | 0/0.0317 | 0.1587/0 | 0/0 |
| max depth | 6/3 | 5/5 | 6/5 |
| min child weight | 4/5 | 3/13 | 8/10 |
| subsample | 0.5757/0.7957 | 0.4343/0.5129 | 0.2457/0.6700 |
| colsample_bytree | 0.1111/0.0794 | 0.4286/0.1270 | 0.4762/0.8095 |
| lambda | 0.01/1.1156 | 1.2103/0.3259 | 1.4946/0.7997 |
| Iteration times | 4174/1476 | 5841/3460 | 4968/4492 |

Table 2 Comparison of the algorithm in this paper with other algorithms (Taking the CTRP-L1000-24h(S1) dataset as an example).

| Methods | Pearson Correlation | R^2 | Mean Squared Error |
|-----------|---------------------|--------|--------------------|
| Our model | 0.8321 | 0.6922 | 0.025 |
| PCA-Lasso | 0.7887 | 0.6211 | 0.031 |
| PCA-SVR | 0.7900 | 0.6211 | 0.031 |
| FTest-RF | 0.7942 | 0.6304 | 0.030 |
| MI-KNN | 0.8054 | 0.6414 | 0.030 |
| VAE | 0.8289 | 0.6823 | 0.026 |
| DAE-NN | 0.8033 | 0.6174 | 0.032 |

Additional Files

Additional file 1 — Supplementary Material

Supplementary Material contains supplementary figures and supplementary tables of the results in this study.

Figures

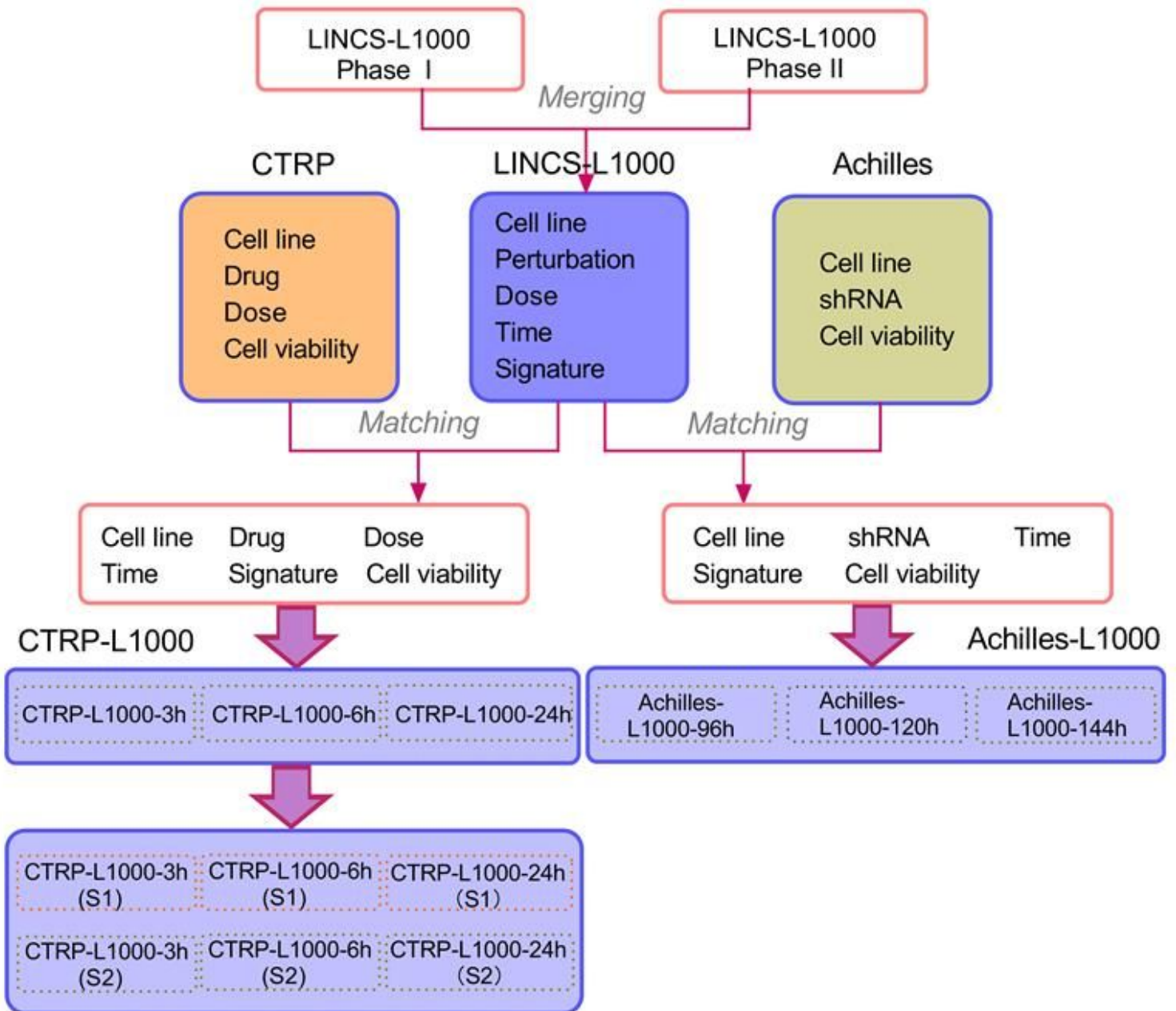


Figure 1

LINCS-L1000 and CTRP, Achilles data association diagram. The process of data association consisted of two parts: perturbation signatures and cell phenotypic information. The LINCS-L1000-Phase I and LINCS-L1000-Phase II were combined and renamed LINCS-L1000. The compound perturbation signatures and shRNA perturbation signatures involved M LINCS-L1000 were respectively associated with CTRP and Achilles datasets according to relevant conditions, which were named CTRP-L1000 and Achilles-L1000. The datasets were divided into CTRP-L1000-30 CTRP-L1000-60 CTRP-L1000-24h, Achilles-L1000-96h, Achilles-L1000-120h and Achilles-L1000-144h according to different perturbation time. CTRP-L1000-3h,

CTRP-L1000-611, and CTRP-L1000-24h were divided into six subsets according to the concentration factor was considered(52) or not considered(S1).

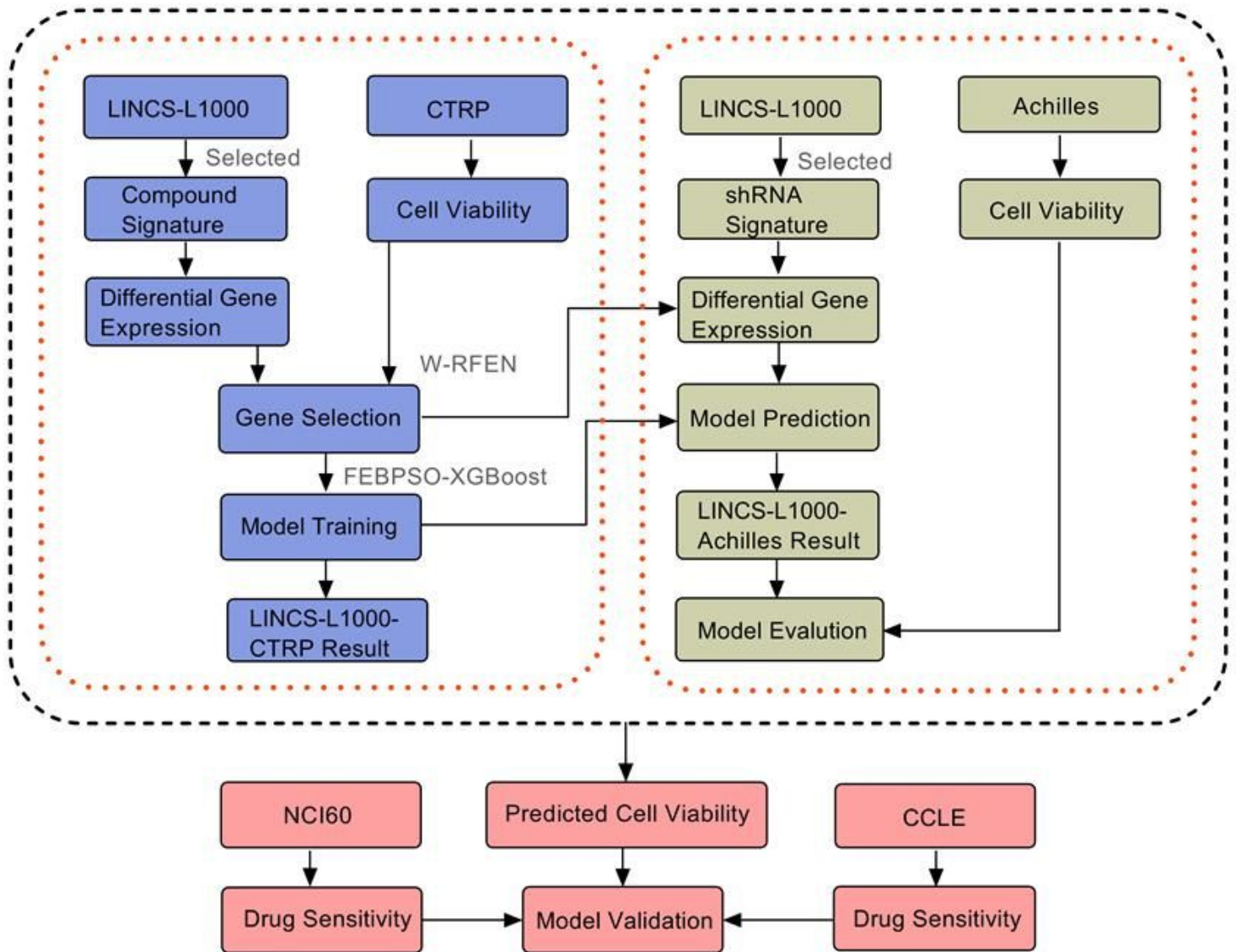


Figure 2

Framework diagram of cell viability prediction.

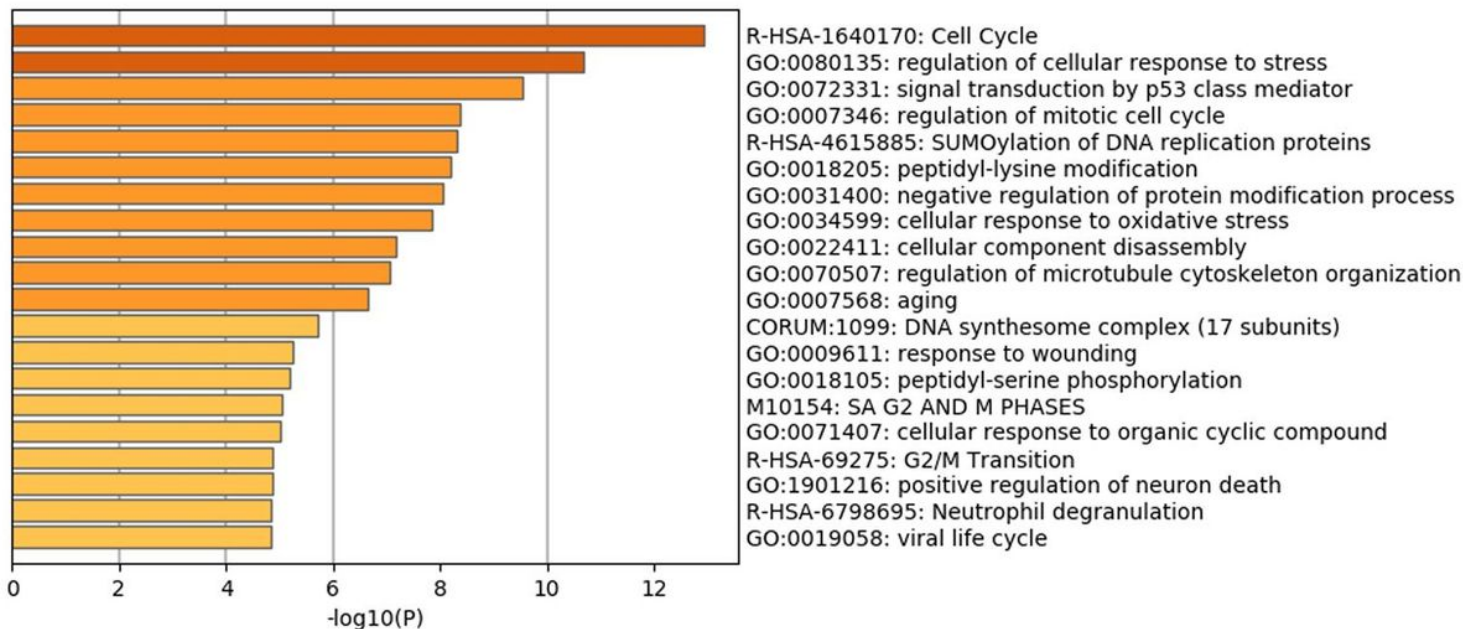


Figure 3

Enrichment analysis of differentially expressed genes in the CTRP-L1000-24h dataset.

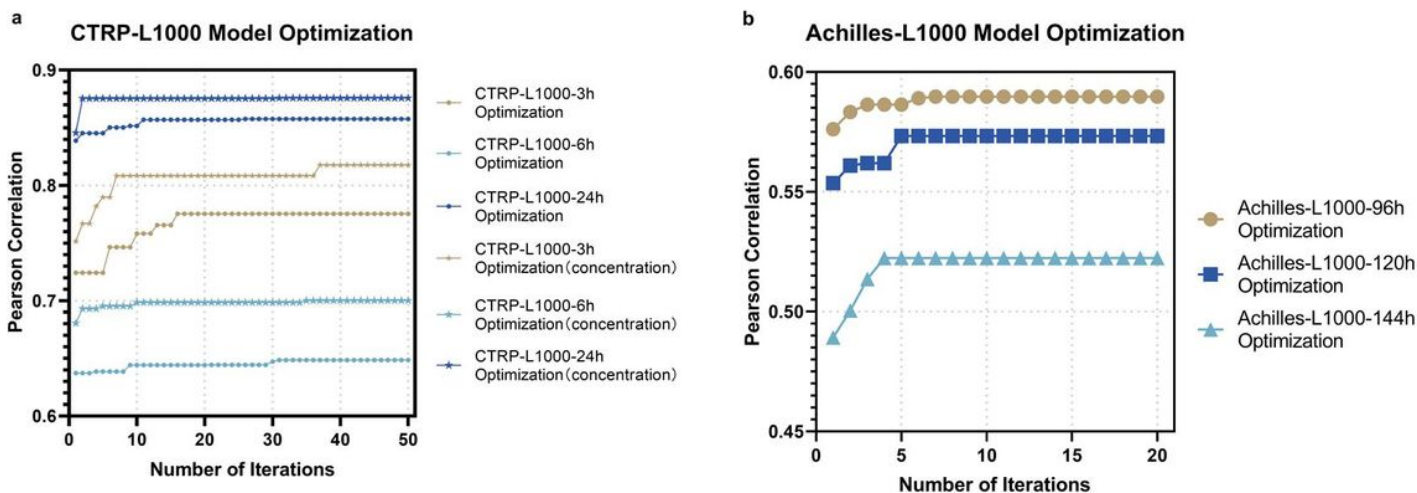


Figure 4

Iterative process of FEBPSO in XGBoost algorithm. a, CTRP-L1000 Optimization. b, Achilles-L1000 Optimization.

Across Screen Validation

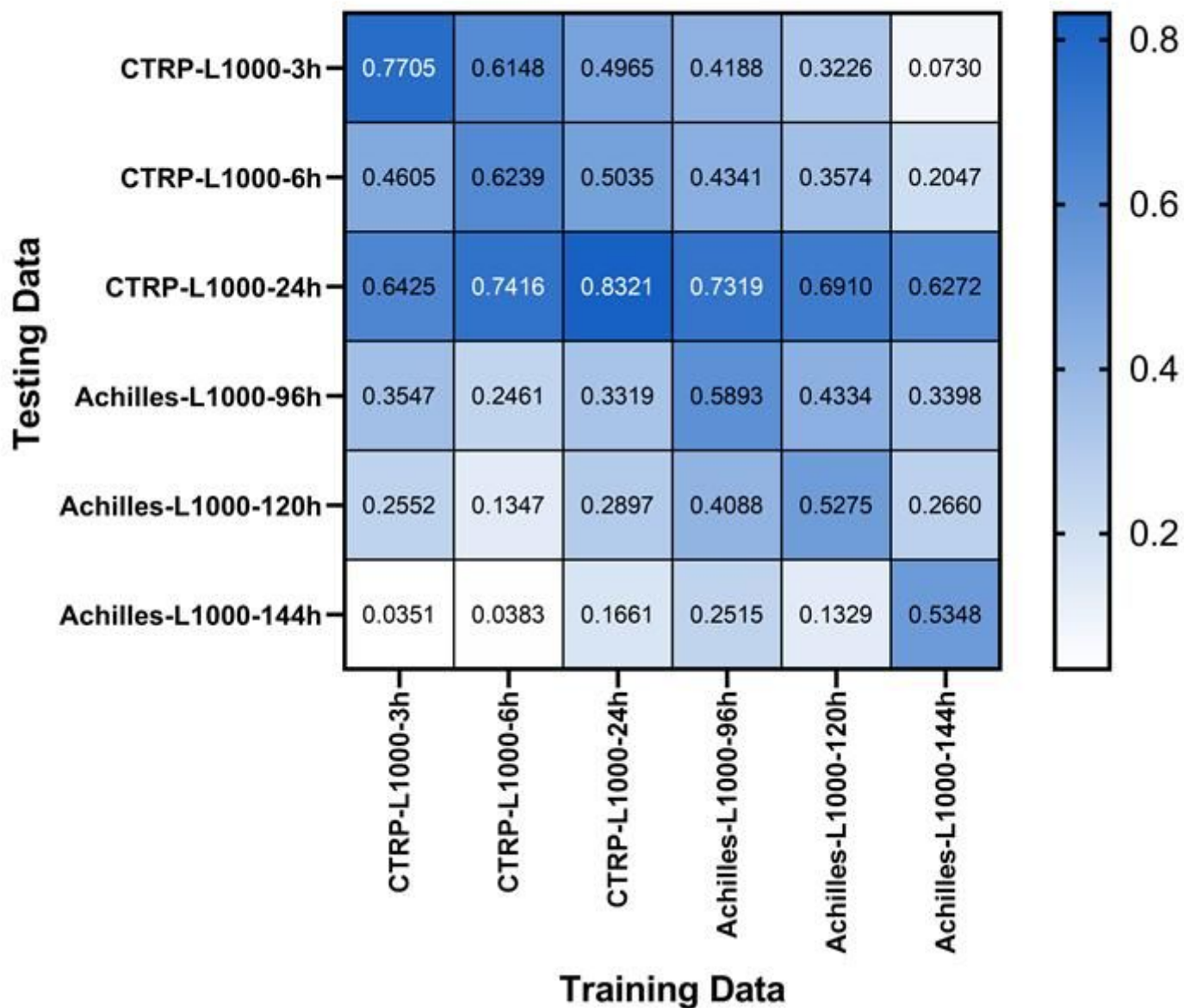


Figure 5

Independent dataset validation. Using the Achilles-L1000 series model to predict cell viability in CTRP-L1000 data and vice versa.

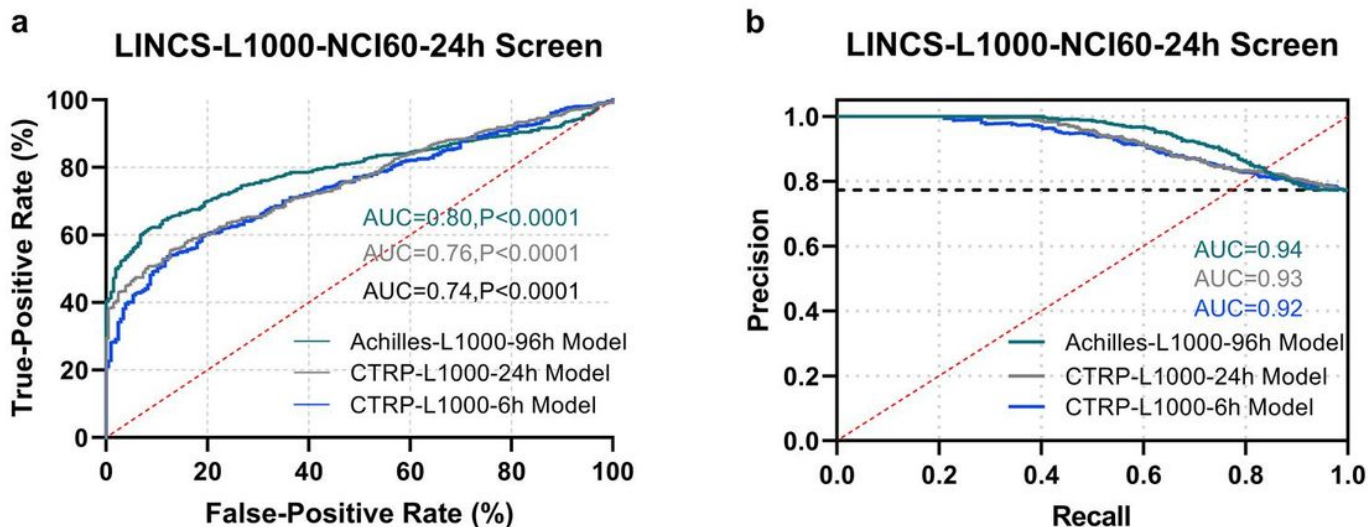


Figure 6

ROC curve and PR curve of the model evaluation on LINC-S-L1000-NCI60-24h dataset. a, The graph of Receiver Operating Characteristic. b, The graph of Precision-Recall.

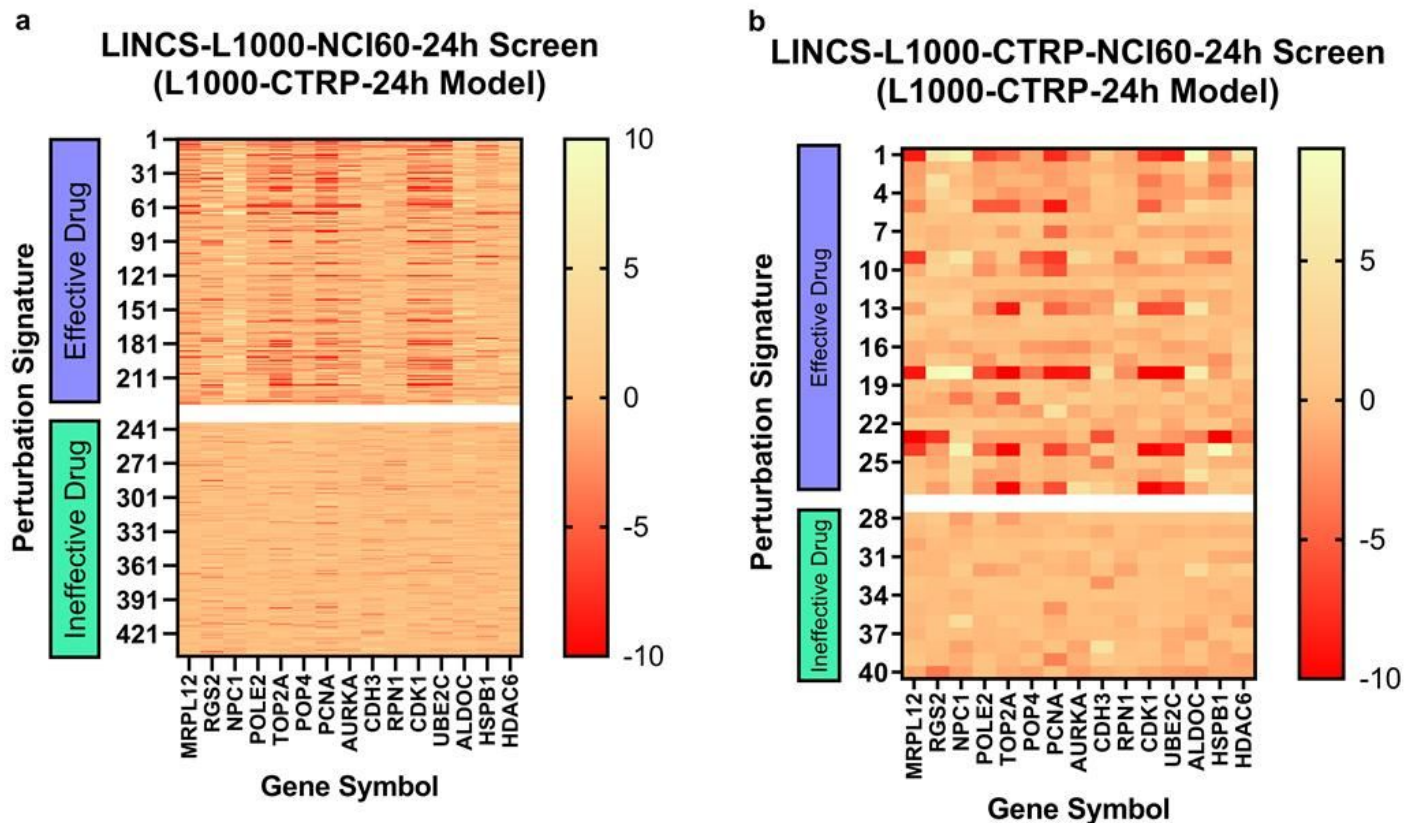


Figure 7

Heat map of the rst fteen genes. a, LINCS-L1000-NCI60-24h. b, LINCS-L1000-CTRP-NCI60-24h.

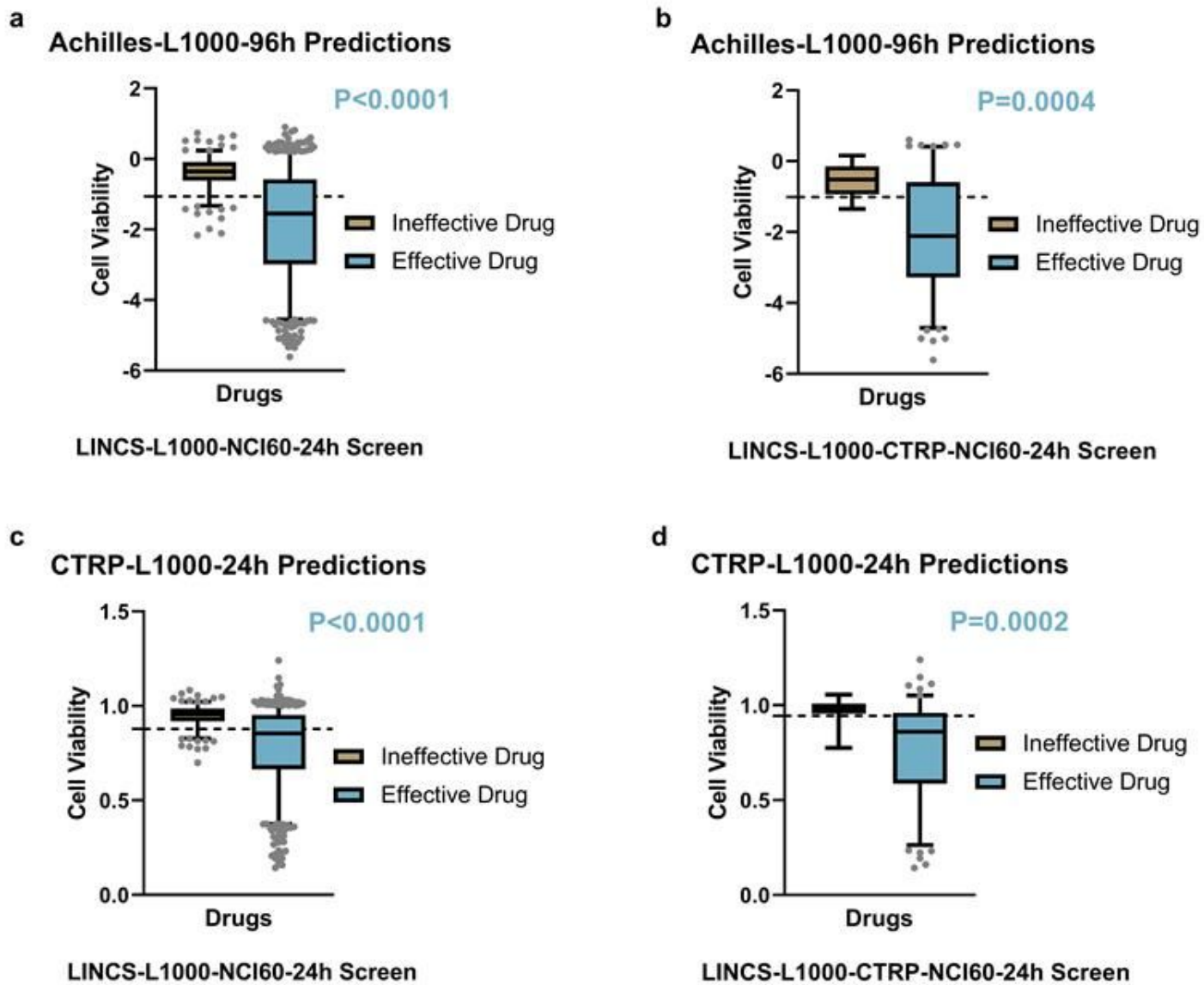


Figure 8

Box plot. comparison of the effective drug group and the ineffective drug group.

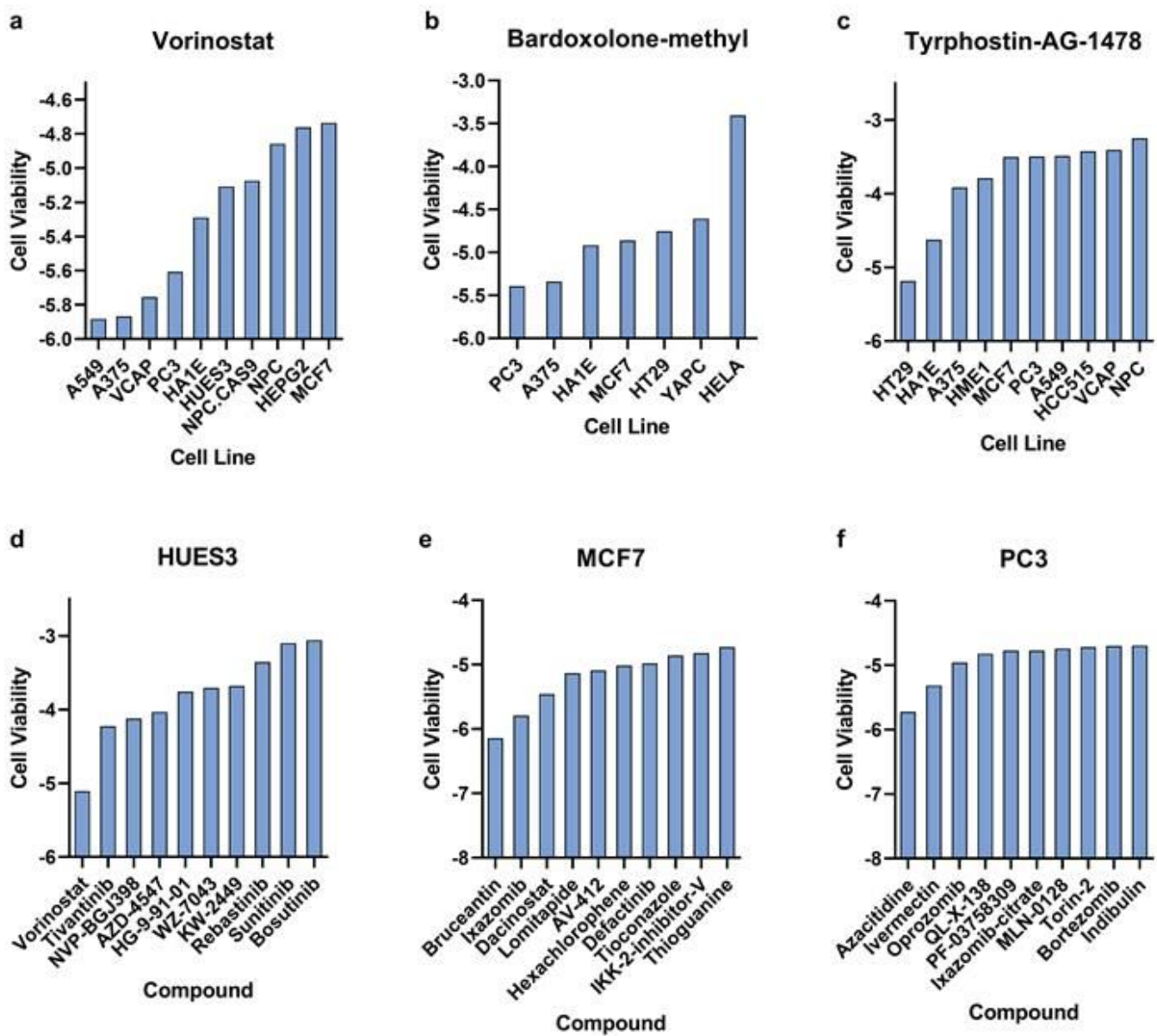


Figure 9

The predicted cell viability for different drugs and cell lines. (a-c) showed the cell viability of the drugs Vorinostat, Bardoxolone-methyl and Tyrphostin-AG-1478 in different cell lines. (d-f) showed the cell viability of the cell lines HUES3, MCF7, PC3 in different drugs.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementaryMaterial.pdf](#)