



T4SE-XGB: Interpretable Sequence-Based Prediction of Type IV Secreted Effectors Using eXtreme Gradient Boosting Algorithm

Tianhang Chen¹, Xiangeng Wang^{1,2}, Yanyi Chu^{1,3}, Yanjing Wang¹, Mingming Jiang¹, Dong-Qing Wei^{1,3*} and Yi Xiong^{1*}

¹ State Key Laboratory of Microbial Metabolism, School of Life Sciences and Biotechnology, Shanghai Jiao Tong University, Shanghai, China, ² Department of Biomedical Sciences, City University of Hong Kong, Hong Kong, China, ³ Peng Cheng Laboratory, Shenzhen, China

OPEN ACCESS

Edited by:

Qi Zhao,
Liaoning University, China

Reviewed by:

Balachandran Manavalan,
Ajou University, South Korea
Xiaoyang Jing,
Toyota Technological Institute at
Chicago, United States

*Correspondence:

Dong-Qing Wei
dqwei@sjtu.edu.cn
Yi Xiong
xiongyi@sjtu.edu.cn

Specialty section:

This article was submitted to
Systems Microbiology,
a section of the journal
Frontiers in Microbiology

Received: 06 July 2020

Accepted: 21 August 2020

Published: 24 September 2020

Citation:

Chen T, Wang X, Chu Y, Wang Y,
Jiang M, Wei D-Q and Xiong Y (2020)
T4SE-XGB: Interpretable
Sequence-Based Prediction of Type IV
Secreted Effectors Using eXtreme
Gradient Boosting Algorithm.
Front. Microbiol. 11:580382.
doi: 10.3389/fmicb.2020.580382

Type IV secreted effectors (T4SEs) can be translocated into the cytosol of host cells via type IV secretion system (T4SS) and cause diseases. However, experimental approaches to identify T4SEs are time- and resource-consuming, and the existing computational tools based on machine learning techniques have some obvious limitations such as the lack of interpretability in the prediction models. In this study, we proposed a new model, T4SE-XGB, which uses the eXtreme gradient boosting (XGBoost) algorithm for accurate identification of type IV effectors based on optimal features based on protein sequences. After trying 20 different types of features, the best performance was achieved when all features were fed into XGBoost by the 5-fold cross validation in comparison with other machine learning methods. Then, the ReliefF algorithm was adopted to get the optimal feature set on our dataset, which further improved the model performance. T4SE-XGB exhibited highest predictive performance on the independent test set and outperformed other published prediction tools. Furthermore, the SHAP method was used to interpret the contribution of features to model predictions. The identification of key features can contribute to improved understanding of multifactorial contributors to host-pathogen interactions and bacterial pathogenesis. In addition to type IV effector prediction, we believe that the proposed framework can provide instructive guidance for similar studies to construct prediction methods on related biological problems. The data and source code of this study can be freely accessed at <https://github.com/CT001002/T4SE-XGB>.

Keywords: type IV secreted effector, feature selection, extreme gradient boosting, interpretable analysis, SHAP (SHapley additive exPlanations)

INTRODUCTION

Different secretion systems have been found in bacteria that secrete proteins into the extracellular environment. Gram-negative bacterial secretion can be categorized into eight types (from type I to type VIII), and the secreted proteins (also called effectors) play a vital role in bacterial pathogenesis and bacterium-host interactions. Some databases or web resource have been developed to store the experimentally validated effectors of Type III, IV, and VI secretion systems (Bi et al., 2013; Li et al., 2015; Eichinger et al., 2016; An et al., 2017). Type IV secretion system (T4SS) are protein complexes found in various species that deliver proteins into the cytoplasm of host cell and thus

cause infection, such as whooping cough (Dorji et al., 2018), gastritis, peptic ulcer, and crown-gall tumor (Kuzmanovic et al., 2018). Therefore, the identification of type IV secreted effector proteins (T4SEs) is a fundamental step toward understanding of the pathogenic mechanism of T4SS.

There are a variety of experimental methods for identifying new T4SEs such as immunoblot analysis and pull-down assay (Cunha et al., 2015). However, they are limited by both a *priori* knowledge about biological mechanisms and the sophisticated implementation of molecular experiments (Zeng and Zou, 2019). Furthermore, these experimental approaches are quite time-consuming and expensive. Instead, a large number of computational methods have been developed for prediction of T4SEs in the last decade, which successfully speed up the process in terms of time and efficiency. These computational approaches can be categorized into two main groups: the first group of approaches infer new effectors based on sequence similarity with currently known effectors (Chen et al., 2010; Lockwood et al., 2011; Marchesini et al., 2011; Meyer et al., 2013; Sankarasubramanian et al., 2016; Noroy et al., 2019) or phylogenetic profiling analysis (Zalguizuri et al., 2019), and the second group of approaches involve learning the patterns of known secreted effectors that distinguish them from non-secreted proteins based on machine learning and deep learning techniques (Burstein et al., 2009; Lifshitz et al., 2013; Zou et al., 2013; Wang et al., 2014; Ashari et al., 2017; Wang Y. et al., 2017; Esna Ashari et al., 2018, 2019a,b; Guo et al., 2018; Xiong et al., 2018; Xue et al., 2018; Acici et al., 2019; Chao et al., 2019; Hong et al., 2019; Wang J. et al., 2019; Li J. et al., 2020; Yan et al., 2020). In the latter group of methods, Burstein et al. (2009) worked on *Legionella pneumophila* to identify T4SEs and validated 40 novel effectors which were predicted by machine learning algorithms. Several features such as genomic organization, evolutionary based attributes, regulatory network attributes, and attributes specific to the *L. pneumophila* pathogenesis system were applied as input of the different machine learning algorithms: naïve Bayes, Bayesian networks, support vector machine (SVM), neural network and a voting classifier based on these four algorithms. Then, Zou et al. (2013) built the tool called T4EffPred based on the SVM algorithm with features such as amino acid composition (AAC), dipeptide composition (DPC), position specific scoring matrix composition (PSSM), auto covariance transformation of PSSM to identify T4SEs. Wang et al. (2014) constructed an effective inter-species T4SS effector prediction tool named T4SEpre, based on SVM by using C-terminal sequential and position-specific amino acid compositions, possible motifs, and structural features. Later, Xiong et al. (2018) used the same dataset as that of the previous study (Wang Y. et al., 2017) and developed a stacked ensemble classifier PredT4SE-Stack using various machine learning algorithms, such as SVM, gradient boosting machine, and extremely randomized trees. Wang J. et al. (2019) developed an ensemble classifier called Bastion4 which serves as an online T4SS effector predictor. They calculated 10 types of sequence-derived features. Then, the naïve Bayes, *k*-nearest neighbor, logistic regression, random forest, SVM and multilayer perceptron were trained and compared. Significantly improved predictive performance was achieved when they used

the majority voting strategy based on the six classifiers where the PSSM-based features were used as input vectors. Recently, Esna Ashari et al. developed the software called OPT4e (Esna Ashari et al., 2019a), which assembled various features used in prior studies to predict a set of candidate effectors for *A. phagocytophilum*. This tool yielded reasonable candidate effector predictions for most T4SS bacteria from the *Alphaproteobacteria* and *Gammaproteobacteria* classes.

Besides the traditional machine learning methods, deep learning is a new technology based on neural network architecture and has been successfully applied in wide range of applications in recent years (Yu et al., 2018; Lv et al., 2019; Ren et al., 2019; Wu et al., 2019; Deng et al., 2020). Some researchers explored deep learning techniques to identify T4SEs based on protein sequences. Xue et al. (2018) proposed a deep learning method to identify T4SEs from protein sequences. The model called DeepT4 utilized a convolutional neural network (CNN) to extract T4SEs-related features from 50 N-terminal and 100 C-terminal residues of the proteins. This work provided the original idea about using the deep learning method. However, only few information of protein sequences can be extracted, which showed a slightly weaker performance compared with the Bastion4 (Wang J. et al., 2019). Later, Acici et al. (2019) developed the CNN-based model based on the conversion from protein sequences to images using AAC, DPC and PSSM feature extraction methods. More recently, Hong et al. (2019) developed the new tool CNN-T4SE based on CNN, which integrated three encoding strategies: PSSM, protein secondary structure & solvent accessibility (PSSSA) and one-hot encoding scheme (Onehot), respectively. Compared with other machine learning methods, CNN-T4SE outperform all other state-of-the-art sequence-based T4SEs predictors. However, the less-than-optimal features analysis causes the limited deep learning for protein data and it is not straightforward to understand which features extracted from a given protein sequence drive the final prediction.

In this study, we proposed T4SE-XGB to predict type IV effectors using XGBoost based on sequence-derived features. To overcome the limitations of existing methods, we selectively summarized the features covered in previous studies and added some new features. The main strength of our method hinges on two aspects. On one hand, T4SE-XGB trained with features selected by the ReliefF algorithm significantly improved the overall performance on the benchmark dataset. On the other hand, T4SE-XGB uses a *post-hoc* interpretation technique: the SHAP (SHapley Additive exPlanations) method to demystify and explain specific features that led to deeper understanding of “black box” models.

MATERIALS AND METHODS

The overall workflow of T4SE-XGB is shown in **Figure 1**, which is composed of five stages: Dataset Collection, Feature Extraction, Feature Selection, Model Construction, and Model Interpretation. The detailed steps are described in the following sections.

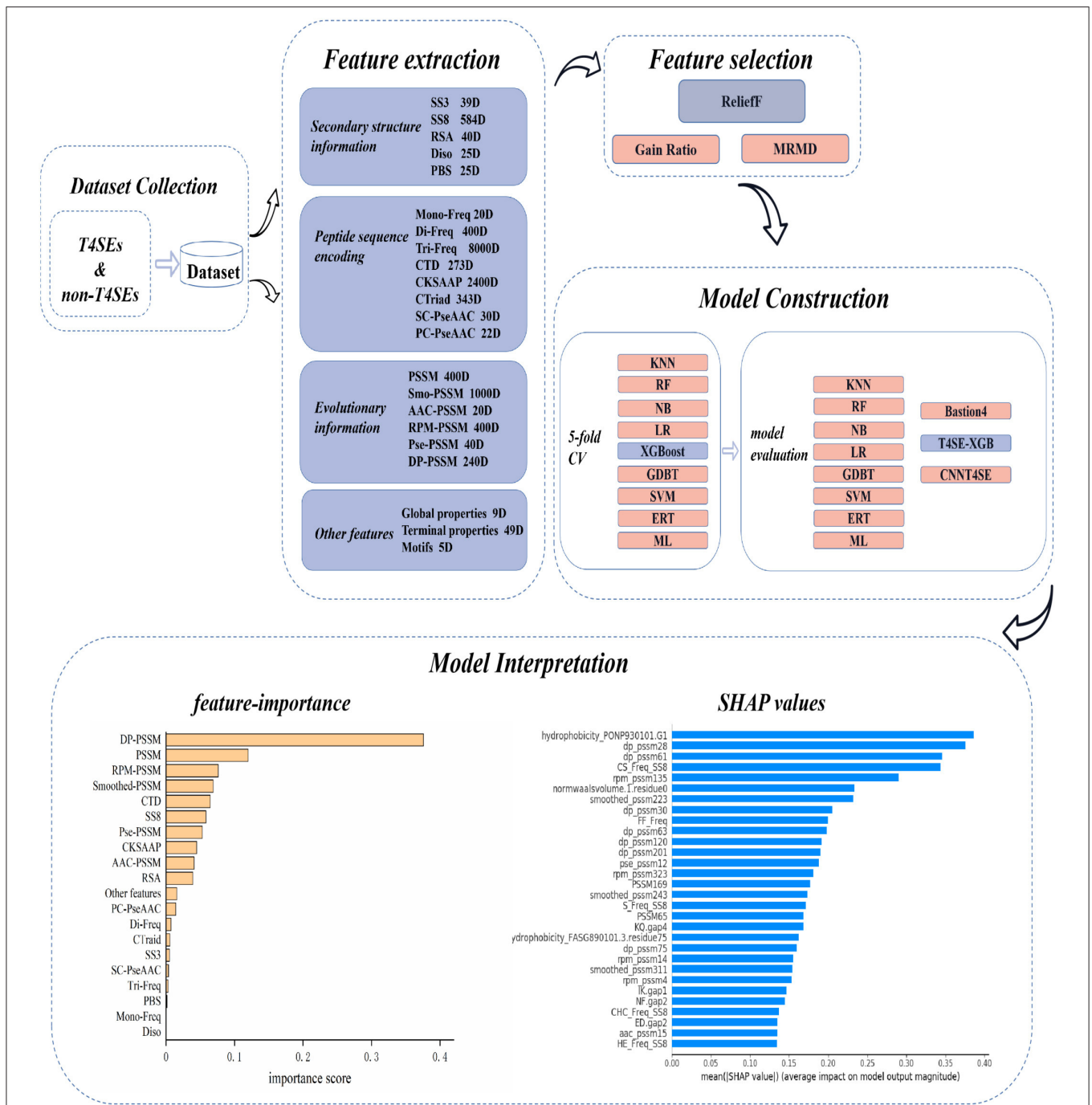


FIGURE 1 | Overview workflow of T4SE-XGB. First, the benchmark dataset was collected. Next, 20 types of features were used to extract information from original protein sequences. Then, the ReliefF algorithm was employed to select optimal features. Five-fold cross validation test and independent test were set to verify the validation of the model. Finally, we not only used the vanilla XGBoost method to get feature importance, but also got SHAP values to realize the model interpretation.

Dataset

In our study, type IV secreted effectors and non-effectors were selected to build the benchmark dataset to construct the machine learning-based model for prediction of T4SEs. Our dataset was directly obtained from the recently published study (Wang

J. et al., 2019), which contained 420 T4SEs and 1262 non-T4SEs. The protein dataset was passed through a filter of >30% sequence identity to reduce sequence redundancy by the CD-HIT tool (Huang et al., 2010). In the end, we got the final training dataset consisting of 365 T4SE and 1106 non-T4SEs,

and the independent test dataset consisting of 20 T4SEs and 139 non-T4SEs.

Feature Extraction

In this work, we took full advantage of features derived from protein sequences that former researchers have used and also added some novel features, which were used in other large-scale protein function prediction problems. We utilized the following four aspects of features to characterize protein sequences: secondary structure information (Zhu X. J. et al., 2019), peptide sequence encoding (Yang et al., 2019), evolutionary information and other features. Details about feature extraction are listed as below.

Secondary Structure Information

(i) First, we used SCRATCH (Cheng et al., 2005) to predict 3- and 8-state secondary structure (SS) information of amino acids of sequences and then *mono*- (1 state i.e., turn, strand, or coil), *di*- (two consecutive states) and *tri*-state (three consecutive states) frequencies from a given protein sequence were extracted. (ii) The fraction of exposed residues (FER) with 20 different relative solvent accessibility (RSA) cutoffs (0–95% cutoffs at 5% intervals) and the FER by the average hydrophobicity of these exposed residues at different RSA cutoffs were calculated. (iii) DISOPRED (Ward et al., 2004) can predict precise disordered region with annotated protein-binding activity. In the former study by Elbasir et al. (2019), they used DISOPRED to get 25 disordered features and 25 features of protein binding sites (PBS) in disordered regions.

Peptide Sequence Encoding

(i) Frequencies of 20 amino acids, 400 di-peptides, 8,000 tri-peptides were extracted from the protein sequences. (ii) The Composition, Transition, and Distribution (CTD) feature represents the amino acid distribution patterns of a specific structural or physicochemical property in a protein or peptide sequence. Various types of physicochemical properties such as hydrophobicity, normalized Van der Waals Volume, polarity, polarizability, charge, secondary structures, and solvent accessibility have been used for constituting the final feature vectors. (iii) The composition of *k*-spaced amino acid pairs (CKSAAP) feature calculates frequencies of amino acid pairs separated by any *k* (ranging from 0 to 5) residues. We use the default maximum value of *k* which is 5, and got a 2,400-dimensional feature vector for one protein sequence. (iv) The Conjoint Triad descriptor (CTriad) considers the properties of one amino acid and its vicinal amino acids by regarding any three continuous amino acids as a single unit (Shen et al., 2007). The occurrence that each triad appearing in the protein sequence is used to constitute a 343-dimensional vector after the amino acids are categorized into seven classes. (v) Pseudo amino acid composition analyses protein sequences about the physicochemical properties of the constituent amino acids. The final feature vectors include the global or long-range sequence order information. Series correlation pseudo amino acid composition (SC-PseAAC) (Chou, 2004) is a variant of PseAAC, which generates protein feature vectors by combining the amino acid composition and global sequence-order effects

via series correlation. Parallel correlation pseudo amino acid composition (PC-PseAAC) (Chou, 2001), derived from PseAAC, incorporates the contiguous local sequence-order information and the global sequence-order information into feature vectors of protein sequences.

The iFeature (Chen et al., 2018) sever is capable of calculating and extracting different sequence-based, structural, and physicochemical features derived from protein sequences. The BioSeq-Analysis2.0 (Liu B. et al., 2019) sever was employed to generate modes of pseudo amino acid compositions (such as SC-PseAAC and PC-PseAAC) for protein sequences.

Evolutionary Information

(i) PSSM of a protein sequence can be obtained in the form of $L \times 20$ matrix (L is the amino acid length). PSSM represents the evolutionary, residue, and sequence information features of input proteins. In our study, we got 400 feature vectors from the original PSSM profile by summing rows corresponding to the same amino acid residue. (ii) Smoothed-PSSM (Cheng et al., 2008) transformed from the standard PSSM encodes the correlation or dependency from surrounding residues, which significantly enhanced the performance of RNA-binding site prediction in proteins. The Smoothed-PSSM profile considered the first 50 amino acids starting from the protein's N-terminus to form a vector with the dimension 1,000. (iii) AAC-PSSM (Liu et al., 2010) represents the correlation of evolutionary conservation of the 20 residues between two positions separated by a predefined distance along the sequence and successfully converts a protein into a fixed length feature vector with dimension 20. It reveals the possibility of the amino acid residues mutated to different types during the evolution process. (iv) RPM-PSSM (Jeong et al., 2011) filters all entities with values of < 0 from the PSSM matrix by using the residue probing method, in which each amino acid is regarded as a probe corresponds to a particular column in the PSSM profiles, and the negative values were set to 0. For this method, original PSSM matrix finally transformed into the 20×20 matrix and can be constructed into a 400-dimensional vector. (v) Pse-PSSM (Chou and Shen, 2007) is similar to PseAAC and encodes the PSSM of proteins with different lengths using a uniform length matrix. (iv) DP-PSSM (Juan et al., 2009), a protein descriptor based on similarity, gets the hidden sequential order information and can avoid cancellation of positive or negative terms in the average process. As a result, we obtained a 400-dimensional vector for each sequence.

These PSSM-based features were achieved using the bioinformatics tool called POSSUM (Wang J. et al., 2017), including the original PSSM profiles, smoothed-PSSM, AAC-PSSM, RPM-PSSM, Pse-PSSM, and DP-PSSM. All PSSM-based features used default parameters the website provided: smoothing window=7 and sliding window=50 for smoothed-PSSM, $\xi = 1$ for Pse-PSSM, and $\alpha = 5$ for DP-PSSM.

Other Features

(i) Global properties of the protein were calculated such as sequence length, molecular weight, and total hydrophathy et al. and the list is shown in the **Supplementary Table 3**. (ii) Terminal properties like the frequencies of 20 amino acid types in the 50

amino acids at the C-terminus or N-terminus used in previous studies were also calculated (Wang Y. et al., 2017; An et al., 2018; Wang J. et al., 2019; Zeng and Zou, 2019). The frequencies of di-peptides at the C-terminus, like SS, KE, EE, EK, AA, AG, and LL involved in former studies have shown variances between effectors and the non-effectors were also calculated (Zou et al., 2013; Zou and Chen, 2016). (iii) We also searched for several types of protein motifs including nuclear localization signals (NLS), E-Block (EEXXE motif), conserved EPIYA motifs (EPIYA_CON), hypothetical EPIYA motifs (EPIYA_HYS), and Prenylation Domain (CaaX motif) that have been proposed and extracted before (Esna Ashari et al., 2018, 2019a; Noroy et al., 2019).

Feature Normalization

Normalization is a scaling technique in which values are shifted and rescaled so that they fall into the same numeric interval. Having features on a similar scale helps the gradient descent converge more quickly toward the minima. The following formula can be used to normalize all feature values and end up ranging between 0 and 1, which is known as Min-Max scaling:

$$X' = \frac{X - X_{\min}}{X_{\max} - X_{\min}} \quad (1)$$

Here, X_{\max} and X_{\min} are the maximum and the minimum values of the feature, respectively. We imported the MinMaxScaler package from the python scikit learn library to calculate the normalized values.

Feature Selection

The purpose of dimensionality reduction or feature selection is to reduce the computational time and complexity of the prediction model, and also to provide more insights into the data abundance (Xiong et al., 2012, 2019; He et al., 2018; Manavalan et al., 2018; Tang et al., 2018; Jing et al., 2019; Kang et al., 2019; Liu Y. et al., 2019; Shi et al., 2019; Basith et al., 2020; Govindaraj et al., 2020; Li K. et al., 2020; Su et al., 2020; Zhang S. et al., 2020). It is indispensable to reduce dimensionality to remove redundant features so that we can reserve the important ones.

Gain Ratio

The gain ratio algorithm based on information theory can be used to deal with oversized feature sets (Shannon, 1948). We used the gain ratio function from the R package named FSelector. The algorithm finds weights of discrete attributes based on their correlation with continuous class attribute.

ReliefF

The ReliefF algorithm is an improvement of Kononenko's standard Relief algorithm (Kira and Rendell, 1992). In this work, the ReliefF algorithm was implemented by the ReliefExpRank function in the attrEval method from R package named CORElearn (Yu et al., 2019). Rank of nearest instance is determined by the increasing (Manhattan) distance from the selected instance and the k nearest instances have weight exponentially decreasing with increasing rank. This is a default choice for methods taking conditional dependencies among the attributes into account.

The ReliefF algorithm fully considers the correlation between features and labels, in order to effectively remove unnecessary features after updating the feature weights according to the degree of correlation. The higher the weight value, the stronger the classification ability of the feature. The weight W of each feature is defined as:

$$W = W - \frac{\sum_{x \in X} d(f, m_i, x_i)}{nd} + \frac{\sum_{c \neq \text{class}(m_i)} \left[\frac{P(C)}{1 - P(\text{class}(m_i))} \sum_{x \in Y(c_i)} d(f, m_i, y_i) \right]}{nd} \quad (2)$$

where d means samples with the nearest distance m_i from each category selected by the ReliefF algorithm first, f means a certain feature, n means the number of samples, and $d(f, m_i, x)$ means the distance between sample X and sample Y for a certain feature f .

Maximum-Relevance-Maximum-Distance

The Maximum-Relevance-Maximum-Distance (MRMD) algorithm uses the Pearson's correlation coefficient to measure the relevance between features in a subset. The Pearson correlation coefficient shows the degree of relationship between features and labels. Besides, Euclidean distance, cosine similarity, and Tanimoto coefficient are utilized to calculate the redundancy between features in a subset. In the end, the MRMD algorithm selects features which have strong correlation with labels and have lowest redundancy between features (Zou et al., 2016).

Extreme Gradient Boosting

Extreme gradient boosting also named XGBoost (Chen and Guestrin, 2016) is an optimized distributed gradient boosting algorithm designed to be highly efficient, flexible, and portable (Wang et al., 2019a). XGBoost based on decision tree ensembles consists of a set of classification or regression trees. It uses the training data (with multiple features) x_i to predict a target variable y_i .

To begin with, the objective function is defined as:

$$obj = \sum_{i=1}^n l(y_i, \hat{y}_i^{(t)}) + \sum_{i=1}^t \Omega(f_i) \quad (3)$$

where n is the number of trees, l is the training loss function, Ω is the regularization term.

Then, the XGBoost takes the Taylor expansion of the loss function up to the second order and removes all the constants, so the specific objective at step t becomes:

$$L^{(t)} = \sum_{i=1}^n \left[g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i) \right] + \Omega(f_t) \quad (4)$$

where the g_i and h_i are defined as

$$\begin{cases} g_i = \partial_{\hat{y}_i^{(t-1)}} l(y_i, \hat{y}_i^{(t-1)}) \\ h_i = \partial_{\hat{y}_i^{(t-1)}}^2 l(y_i, \hat{y}_i^{(t-1)}) \end{cases} \quad (5)$$

The value of the objective function only depends on g_i and h_i . It can optimize every loss function, including logistic regression and pairwise ranking.

TABLE 1 | Comparison of predictive power of different features on the training data set by 5-fold cross validation test.

Feature types	ACC (%)	SE (%)	PRE (%)	F-score	MCC
ss3	79.94	43.84	64.52	0.5209	0.4125
ss8	81.44	48.77	68.57	0.5640	0.4650
RSA	84.23	58.90	72.80	0.6497	0.5556
Diso	79.81	35.62	68.16	0.4650	0.3856
PBS	79.47	36.16	65.74	0.4643	0.3760
Mono-Freq	85.18	60.00	75.54	0.6669	0.5809
Di-Freq	84.36	53.97	76.16	0.6302	0.5486
Tri-Freq	80.01	30.41	73.85	0.4301	0.3822
PSSM	92.11	78.08	88.99	0.8305	0.7833
smoothed-PSSM	88.51	70.14	81.03	0.7512	0.6806
AAC-PSSM	90.69	74.52	86.27	0.7976	0.7425
RPM-PSSM	91.98	77.26	89.32	0.8262	0.7797
Pse-PSSM	92.59	81.37	88.13	0.8446	0.7984
DP-PSSM	92.79	81.37	88.95	0.8486	0.8039
CKSAAP	84.02	51.78	76.31	0.6153	0.5359
CTD	87.70	67.12	80.24	0.7296	0.6562
CTraid	82.53	49.86	71.04	0.5816	0.4909
SC-PseAAC	85.66	61.92	75.96	0.6811	0.5958
PC-PseAAC	85.52	61.10	76.05	0.6769	0.5913
Other features	84.09	54.52	74.62	0.6288	0.5422
All features	93.95	81.92	93.04	0.8698	0.8346

The traditional treatment of tree learning only emphasized the improved impurity, while the complexity control was left to heuristics. Chen and Guestrin (2016) formally defined the complexity of the tree $\Omega(f)$ to obtain regularization, and the loss function in the t -th tree finally can be rewritten as:

$$L^{(t)} = -\frac{1}{2} \sum_{j=1}^T \frac{G_j^2}{H_j + \lambda} + \gamma T \quad (6)$$

where the G_j and H_j are defined as

$$\begin{cases} G_j = \sum_{i \in I_j} g_i \\ H_j = \sum_{i \in I_j} h_i \end{cases} \quad (7)$$

I_j is the sample set divided into the j -th leaf node according to the decision rules for a given tree. The formula (6) can be used as the score value to evaluate the quality of a tree. They also defined the score it gains when a leaf is split into two leaves:

$$Gain = \frac{1}{2} \left[\frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} \right] - \gamma \quad (8)$$

This formula is composed of the score on the new left leaf, the score on the new right leaf, and the score on the original leaf and regularization on the additional leaf. We can find the best split efficiently by the maximum value of *Gain* through a scan from left to right to get all possible split solutions.

XGBoost with many optimization techniques is able to solve problems using far fewer resources. It is simple to parallel and can greatly enhance the program efficiency with a fast model exploration. More details about XGBoost are given in (Chen and Guestrin, 2016).

Performance Evaluation

In this work, confusion matrix obtained after prediction contains four units: true positive (TP), false positive (FP), false negative (FN), and true negative (TN). In order to evaluate the overall predictive performance of different classification models, we used metrics such as Sensitivity (SE), Specificity (SP), Precision (PRE), Accuracy (ACC), F-score, and Matthew's correlation coefficient (MCC) to evaluate the model. They have been widely used in previous studies (Jing and Dong, 2017; Hu et al., 2018; Zhao et al., 2018a,b; Al-Ajlan and El Allali, 2019; Chu et al., 2019; Lin et al., 2019; Manavalan et al., 2019; Zhang et al., 2019a,b; Zhu X. et al., 2019; Cheng et al., 2020; Hasan et al., 2020; Liu et al., 2020; Yue et al., 2020; Zhang Y.-F. et al., 2020), with a higher value indicating better performances. The performance metrics can be defined as follows:

$$\left\{ \begin{array}{l} \text{Sensitivity} = \frac{TP}{TP+FN} \\ \text{Specificity} = \frac{TN}{TN+FP} \\ \text{Precision} = \frac{TP}{TP+FP} \\ \text{Accuracy} = \frac{TP+TN}{TP+TN+FP+FN} \\ \text{F-score} = \frac{2TP}{2TP+FP+FN} \\ \text{MCC} = \frac{TP \times TN - (FP \times FN)}{\sqrt{(TP+FP) \times (TP+FN) \times (TN+FP) \times (TN+FN)}} \end{array} \right. \quad (9)$$

RESULTS AND DISCUSSION

Comparison of Different Features and Their Combinations on the Training Data Set

In this section, we evaluated the predictive power of the individual types of features and their combinations using the XGBoost classification algorithm by the 5-fold cross validation (CV) method on the training data set. In 5-fold CV, the training data set was randomly divided into five subsets. XGBoost were trained by four subsets and the remaining one was used to evaluate the performance of the model. All steps were repeated five times. The average of the performance measures such as ACC and SE of the training set were calculated and the results are shown in **Table 1**. It can be seen that some individual feature types based on PSSM have higher overall prediction power on the training data set. This observation indicates that the features based on PSSM have better performance in the prediction of T4SE when compared with other types of features.

The combination of different features could depict protein sequences in a more comprehensive manner (Wang et al., 2019b). As illustrated in **Table 1**, using the combined features yield the ACC of 93.95% and the MCC of 0.8346, which are both higher than other PSSM-based features. In summary, compared with single feature-based models, the combination of all features achieved consistently better performance.

Comparison of Three Feature Selection Methods on the Training Data Set

In this section, three kinds of feature selection methods were compared on the training data set by the 5-fold cross validation test. They are gain ratio algorithm (Shannon, 1948),

TABLE 2 | Comparison of different feature selection methods with different dimensions of the selected features on the training data set by 5-fold cross validation test.

Method	500	600	700	800	900	1,000	1,100	1,200	1,300	1,400	1,500
GainRatio	92.79	92.73	93.34	93.34	93.41	93.47	93.34	93.54	93.60	93.41	93.81
MRMD	92.18	92.45	92.32	92.52	92.18	92.66	92.39	92.93	92.86	92.25	92.93
ReliefF	93.74	93.95	93.61	93.74	94.36	94.42	94.22	94.09	93.95	94.02	94.08

maximum relevance–maximum distance (MRMD) (Zou et al., 2016) and ReliefF algorithm (Kira and Rendell, 1992). The ACC of different number of dimensions were obtained and compared by using different feature selection algorithms to select most useful features. As shown in **Table 2**, when the MRMD algorithm was used for dimensionality reduction on the training data set, the highest ACC value was 92.93%. The gain ratio algorithm achieved ACC of 93.81% on the training data set. By comparing the prediction accuracy of three methods in different dimensions, it can be found that the ReliefF algorithm achieved the highest ACC value, 94.42% when the dimension was 1,000, obviously better than the models using all original features.

Therefore, the ReliefF algorithm can effectively eliminate redundant variables and improve prediction accuracy. In the following sections, the ReliefF algorithm was used for dimensionality reduction.

Comparison of Different Classification Algorithms on the Training Data Set

In order to objectively validate the prediction power of the XGBoost algorithm, we compared the performance of this algorithm with other classification algorithms by using the 5-fold cross validation on the training data set. Based on the optimal set of features, other classification algorithms such as Random Forests (RF) (Zhang Y.-F. et al., 2020), naïve Bayes (NB), Logistic Regression (LR), Gradient Boost (GDBT), support vector machine, *k*-nearest neighbor (KNN), Extremely randomized trees (ERT), and Multi-layer Perceptron (MLP) were all trained and compared. The grid search method was employed in this work to optimize hyper-parameters for each classifier (Shan et al., 2019), and the searching ranges of these parameters are shown in the **Supplementary Table 1**. For each ML classifier, we obtained the best hyper-parameter combination based on the highest accuracy by the 5-fold cross validation. The optimal parameters are shown in the **Supplementary Table 2**. **Table 3** shows the comparison results of XGBoost with other classification algorithms on the training data set by 5-fold cross validation.

As shown in **Table 3**, the ACCs of different classifiers were falling within the range from 90.89 to 94.42%, and their MCCs were ranging from 0.76 to 0.84 on the training data set. The results showed that XGBoost achieved the best performance, where the ACC, F-score and MCC were significantly higher than the other classifiers. All in all, the XGBoost algorithm performs better than the other machine learning-based methods when applied on the training data set.

TABLE 3 | Comparison of different classification algorithms on the training data set by 5-fold cross validation test.

	ACC (%)	SE (%)	PRE (%)	F-score	MCC
NB	90.89	84.11	80.76	0.8207	0.7631
ML	92.32	82.47	86.10	0.8409	0.7920
LR	93.00	83.01	88.28	0.8539	0.8101
KNN	93.20	80.82	91.20	0.8544	0.8148
RF	93.27	80.27	91.94	0.8554	0.8163
ERT	93.54	80.55	92.76	0.8604	0.8235
GDBT	93.81	84.11	90.55	0.8710	0.8323
SVM	94.36	83.56	93.28	0.8794	0.8466
XGB	94.42	83.01	94.02	0.8803	0.8481

Comparison With Other Classification Algorithms and Existing Methods on the Independent Test Data Set

To further validate the performance of the proposed model in the real test, we compared the performance of our T4SE-XGB model with other classification algorithms and several state-of-the-art methods on the independent data set. The performance results of these methods are provided in **Table 4**, and the ROC curves are shown in the **Supplementary Figure 1**. To make a fair comparison, the same independent data set, which consists of 20 T4SEs and 139 non-T4SEs, was used for all models.

Among these machine learning-based methods, the results showed that our T4SE-XGB model achieved the overall best performance with an ACC of 97.48%, *F*-value of 90.48% and MCC of 0.8916, followed by the state-of-the-art machine learning model called Bastion4 (Wang J. et al., 2019), which achieved 96.23% on ACC, 86.96% on *F*-value and 0.8579 on MCC. Moreover, the T4SE-XGB trained by fewer training samples also gets more stable prediction performance than the deep learning-based method named CNN-T4SE (VOTE 2/3), which takes the majority votes of the three best-performing convolutional neural network-based models (CNN-PSSM, CNNPSSSA, and CNN-Onehot). The CNN-PSSM, a deep learning-based model based on PSSM features, achieved the best results. However, it gets two less false positive and one less false negative when compared with our model.

In summary, there is a consistent observation (from the results obtained from the 5-fold cross validation test and independent test) that our T4SE-XGB model achieved higher performance in terms of sensitivity, specificity, accuracy, and MCC on both the training data set and independent data set.

TABLE 4 | Comparison of different classification algorithms and the state-of-the-art methods on the independent data set.

Model	TP	FN	TN	FP	ACC (%)	SE (%)	SP (%)	PRE (%)	F-score	MCC
SVM	19	1	134	5	96.23	95.00	96.40	79.17	0.8636	0.8467
LR	19	1	131	8	94.34	95.00	94.24	70.37	0.8085	0.7882
NB	19	1	126	13	91.19	95.00	90.65	59.38	0.7308	0.7084
GDBT	19	1	131	8	94.34	95.00	94.24	70.37	0.8085	0.7882
RF	19	1	132	7	94.97	94.96	95.68	73.08	0.8261	0.8066
ERT	19	1	134	5	96.23	95.00	96.40	79.17	0.8636	0.8467
KNN	20	0	128	11	93.08	100.0	92.09	64.52	0.7843	0.7708
ML	18	2	129	10	92.45	90.0	92.81	64.29	0.7500	0.7209
Bastion4	20	0	133	6	96.23	100.0	95.68	76.92	0.8696	0.8579
CNNT4SE(PSSSA)	14	6	138	1	95.60	70.00	99.28	93.33	0.8000	0.7860
CNNT4SE(Onehot)	14	6	139	0	96.23	70.00	100.0	100.0	0.8235	0.8192
CNNT4SE(PSSM)	20	0	138	1	99.37	100.0	99.28	95.24	0.9756	0.9724
CNNT4SE(VOTE 2/3)	16	4	139	0	97.48	80.00	100.0	100.0	0.8889	0.8818
T4SE-XGB	19	1	136	3	97.48	95.00	97.84	86.36	0.9048	0.8916

Model Interpretation

Estimation of Feature Importance by XGBoost

As a tree-based non-linear machine learning technique, XGBoost can exploit the interactions between the engineered features. In contrast to black-box modeling techniques such as SVM, ANN, CNN, the XGBoost algorithm can easily obtain feature importance scores for all input features. XGBoost can also obtain the importance score efficiently based on the frequency of a feature which is used to split data or according to the average gain a feature brings when it was used during node splitting across all established trees. For the optimal set of features constructed on the benchmark dataset, the importance of each feature during training is available in the **Supplementary Table 4**, which is the sum of information gained when used for splits (tree branching).

The total feature importance contribution of all features according to their feature types are shown in **Table 5** and **Figure 2**. We can see that the DP-PSSM feature gets the maximum value of the importance score which is 0.3758. This may mean that the DP-PSSM feature is more important. Besides, the PSSM feature which incorporated evolutionary information has the importance score of 0.1199, followed by other features based on the transformation of the standard PSSM profile, such as RPM-PSSM and Smoothed-PSSM. There are also other features showing high importance. For example, CTD accounts for 6.46% of all feature importance score. SS8 makes up 5.84% of the total variable importance.

Model Interpretation by SHAP

SHAP, a unified framework for interpreting predictions, assigns each feature an importance value for a particular prediction (Lundberg and Lee, 2017) and improves the interpretability of tree-based models such as random forests, decision trees, and gradient boosted trees (Lundberg et al., 2018, 2020). SHAP is based on the game theoretically optimal Shapley values that can be calculated as below (Lipovetsky and Conklin, 2001):

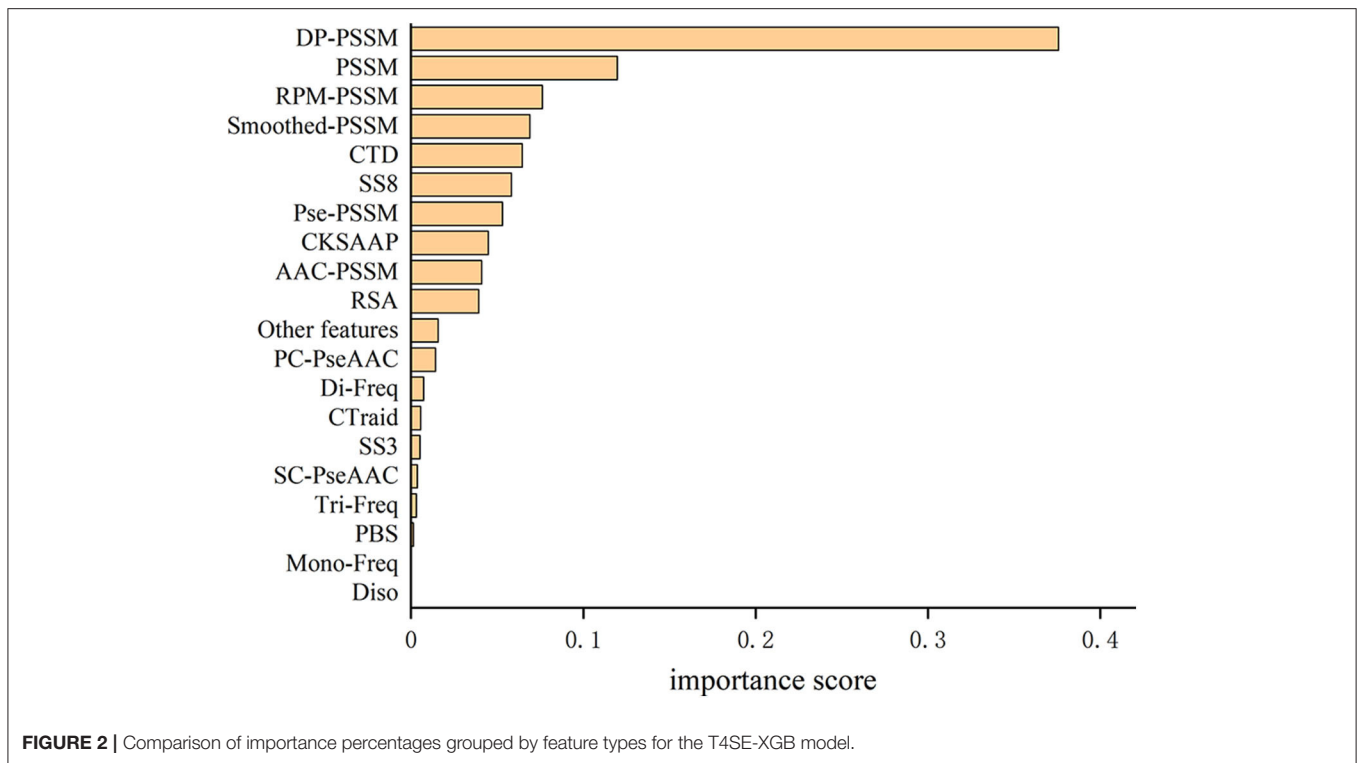
$$\phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|!(|F|-|S|-1)!}{|F|!} [f_{S \cup \{i\}}(x_{S \cup \{i\}}) - f_S(x_S)] \quad (10)$$

TABLE 5 | Importance percentages grouped by feature types for the T4SE-XGB model.

Feature name	Importance score
DP-PSSM	0.3758
PSSM	0.1199
RPM-PSSM	0.0764
Smoothed-PSSM	0.0690
CTD	0.0646
SS8	0.0584
Pse-PSSM	0.0532
CKSAAP	0.0449
AAC-PSSM	0.0411
RSA	0.0394
Other features	0.0158
PC-PseAAC	0.0143
Di-Freq	0.0075
CTraid	0.0057
SS3	0.0053
SC-PseAAC	0.0038
Tri-Freq	0.0032
PBS	0.0015
Diso	0
Mono-Freq	0

where F is the set of all features, S is a subset of the features used in the model without the i th feature, x is the feature vector of the instance to be explained. SHAP comes with many global interpretation methods based on aggregations of Shapley values. More detailed description of the SHAP method is available in (Lundberg and Lee, 2017).

The SHAP method has the ability to provide interpretable predictions and also overcomes limitation that the feature importance scores obtained from XGBoost model, which is in lack of directivity, and is unable to correspond to specific

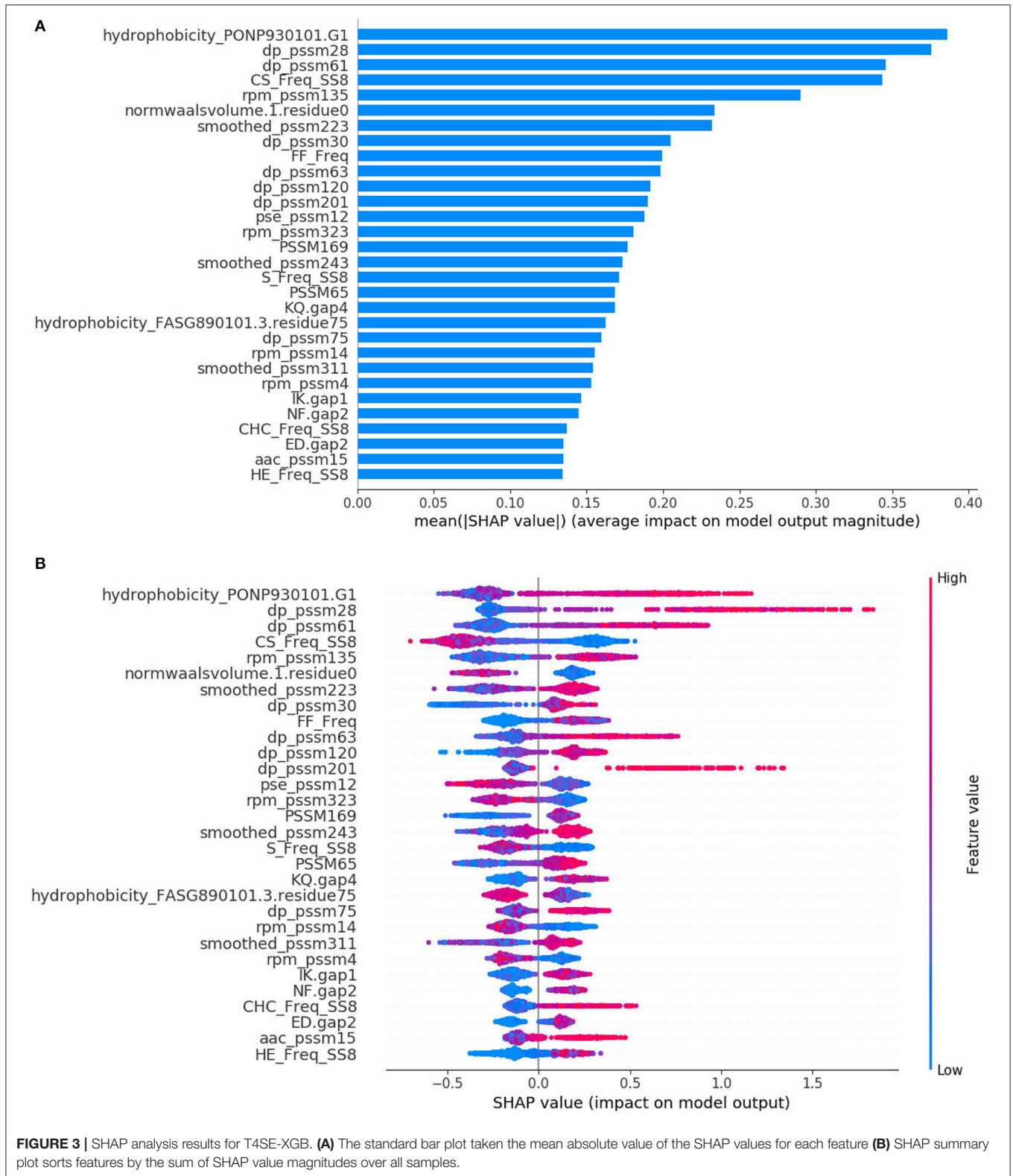


eigenvalues. **Figure 3A** is the standard bar-chart based on the average magnitude of the Shapley values over all training samples. The higher value indicates higher feature importance. It can be seen that DP-PSSM has the largest number of important features, accounting for 7, among the 30 most important features. Meanwhile, other features based on PSSM also form the majority. Among them, the hydrophobicity_PONP930101.G1 coming from the feature unit of CTD can be obviously identified as the most important. Hydrophobicity_PONP930101 is one physicochemical attribute based on the main clusters of the amino acid indices of Tomii and Kanehisa (1996). The hydrophobicity_PONP93-0101.G1= $\{N(r)/N, r \in \{KPDESNT\}\}$ represents the global compositions (percentage) of polar residues of the protein under the hydrophobicity_PONP930101 attribute (Chen et al., 2018). Several studies have suggested that type IV effector proteins exhibited some specificities in regard to amino acid frequency. Zou et al. (2013) calculated the ACC and the variance in their dataset. They found that Asn (N), Glu (E), and Lys (K) have higher compositions in type IVB effectors than non-effectors, and Ala (A), Glu (E), and Ser (S) have higher compositions in type IVA effectors than non-effectors. Some polar amino acids, such as Asp (D), Cys (C), and His (H), have small differences between secreted proteins and non-secreted proteins. Similarly, The Mann-Whitney *U*-test and the permutation test on amino acid frequencies were conducted by An et al. (2018). It was showed that Ala (A), Gly (G), Met (M), Arg (R), Val (V), occurred less frequently in type IV effectors than in cytoplasmic proteins. Meanwhile, Phe (F), Ile (I), Lys (K), Asn (N), Ser (S), Tyr (Y), Thr (T) occurred more frequently in type IV effectors than in cytoplasmic proteins. Since different

benchmark datasets were used, the final results are debatable and incomplete. However, this is the first time to pay attention to the feature named hydrophobicity_PONP930101.G1, which not only corresponds to the amino acid frequency, but also represents the corresponding hydrophilicity. The SHAP summary plot from TreeExplainer (Lundberg et al., 2020) succinctly displays the magnitude, prevalence, and direction of a feature's effect. Each dot in **Figure 3B** corresponds to a protein sample in the study. The position of the dot on the x-axis is the impact that feature has on the model's prediction for that protein. For example, the higher value of hydrophobicity_PONP930101.G1 has higher contribution on predicting a protein being an effector. In contrast, when the values of top features such as CS_Freq_SS8 and normwaalsvolume.1.residue0 are high, the corresponding Shapley values are negative driving the model prediction toward non-effector class. Besides, there are many long tails mean features with a low global importance which can yet be extremely important for specific samples. From the analysis above, it is necessary and effective to consider many characteristics at the same time.

CONCLUSION

In this study, we have presented T4SE-XGB, a predictor developed for accurate identification of T4SE proteins based on the XGBoost algorithm. Especially, we have achieved the state-of-the-art performance compared with previous predictors on the benchmark dataset. There are three major conclusions can be drawn. First, compared with different



algorithm, the XGBoost algorithm gives more stable and accurate prediction performance for prediction of T4SEs. Second, the feature selection method called ReliefF was

utilized to optimize feature vectors, which extracted important features from a large number of candidate features and improved the model performance. Furthermore, unlike other

sequence-based T4SEs predictors, T4SE-XGB can provide meaningful explanation based on samples provided using the feature importance and the SHAP method. It gives us the details about how some features, such as DP-PSSM features and hydrophobicity_PONP930101.G1 from CTD contributed to the final direction of prediction. Meanwhile, it explains the reason why it is essential to pay attention to some certain identities, and also considers a variety of features at the same time.

The final results showed that T4SE-XGB achieved satisfying and promising performance which is stable and credible. However, the model is still constrained by the quantity of T4SE proteins which need to be further updated and the characteristics of T4SEs which need to be discovered. Besides, some potential relationships between features need to be explored. In the future, we plan to find and extract as many features as possible from a large amount of collected data to discriminate type IV secreted effectors from non-effectors.

DATA AVAILABILITY STATEMENT

All datasets generated for this study are included in the article/**Supplementary Material**.

AUTHOR CONTRIBUTIONS

YX and D-QW conceived the study. TC conducted the experiments. TC, YX, XW, YC, YW, and MJ wrote paper.

REFERENCES

- Acici, K., Asuroglu, T., Erdas, C. B., and Ogul, H. (2019). T4SS effector protein prediction with deep learning. *Data* 4:45. doi: 10.3390/data4010045
- Al-Ajlan, A., and El Allali, A. (2019). CNN-MGP: convolutional neural networks for metagenomics gene prediction. *Interdiscip. Sci.* 11, 628–635. doi: 10.1007/s12539-018-0313-4
- An, Y., Wang, J., Li, C., Leier, A., Marquez-Lago, T., Wilksch, J., et al. (2018). Comprehensive assessment and performance improvement of effector protein predictors for bacterial secretion systems III, IV and VI. *Brief Bioinform.* 19, 148–161. doi: 10.1093/bib/bbw100
- An, Y., Wang, J., Li, C., Revote, J., Zhang, Y., Naderer, T., et al. (2017). SecretEPDB: a comprehensive web-based resource for secreted effector proteins of the bacterial types III, IV and VI secretion systems. *Sci. Rep.* 7:41031. doi: 10.1038/srep41031
- Ashari, Z. E., Brayton, K. A., and Broschat, S. L. (2017). “Determining optimal features for predicting type IV secretion system effector proteins for *Coxiella burnetii*. ACM-Bcb’ 2017,” in *Proceedings of the 8th Acm International Conference on Bioinformatics, Computational Biology, and Health Informatics* (Boston, MA), 346–351.
- Basith, S., Manavalan, B., Hwan Shin, T., and Lee, G. (2020). Machine intelligence in peptide therapeutics: a next-generation tool for rapid disease screening. *Med. Res. Rev.* 40, 1276–1314. doi: 10.1002/med.21658
- Bi, D., Liu, L., Tai, C., Deng, Z., Rajakumar, K., and Ou, H. Y. (2013). SecReT4: a web-based bacterial type IV secretion system resource. *Nucleic Acids Res.* 41, D660–D665. doi: 10.1093/nar/gks1248
- Burstein, D., Zusman, T., Degtyar, E., Viner, R., Segal, G., and Pupko, T. (2009). Genome-scale identification of *Legionella pneumophila* effectors using a machine learning approach. *PLoS Pathog* 5:e1000508. doi: 10.1371/journal.ppat.1000508

All authors contributed to the article and approved the submitted version.

FUNDING

This work was supported by the grants from the 20th college students’ innovative practice project of Shanghai Jiao Tong University, the National Natural Science Foundation of China (Contract Nos. 61872094 and 61832019), the Key Research Area Grant No. 2016YFA0501703 of the Ministry of Science and Technology of China, the Science and Technology Commission of Shanghai Municipality (Grant No. 19430750600), and the Natural Science Foundation of Henan Province (162300410060), as well as SJTU JiRLMDS Joint Research Fund and Joint Research Funds for Medical and Engineering and Scientific Research at Shanghai Jiao Tong University (YG2017ZD14, ZH2018QNA41, YG2019GD01, and YG2019ZDA12).

ACKNOWLEDGMENTS

This manuscript has been released as a pre-print at bioRxiv (Chen et al., 2020).

SUPPLEMENTARY MATERIAL

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fmicb.2020.580382/full#supplementary-material>

- Chao, L., Wei, L., and Zou, Q. (2019). SecProMTB: a SVM-based classifier for secretory proteins of *Mycobacterium tuberculosis* with imbalanced data set. *Proteomics* 19:e1900007. doi: 10.1002/pmic.201900007
- Chen, C., Banga, S., Mertens, K., Weber, M. M., Gorbaslieva, I., Tan, Y., et al. (2010). Large-scale identification and translocation of type IV secretion substrates by *Coxiella burnetii*. *Proc. Natl. Acad. Sci. U.S.A.* 107, 21755–21760. doi: 10.1073/pnas.1010485107
- Chen, T., and Guestrin, C. (2016). “XGBoost: a scalable tree boosting system,” in *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (San Francisco, CA), 785–794.
- Chen, T., Wang, X., Chu, Y., Wei, D.-Q., and Xiong, Y. (2020). T4SE-XGB: interpretable sequence-based prediction of type IV secreted effectors using extreme gradient boosting algorithm. *bioRxiv [Preprint]*. doi: 10.1101/2020.06.18.158253
- Chen, Z., Zhao, P., Li, F., Leier, A., Marquez-Lago, T. T., Wang, Y., et al. (2018). iFeature: a Python package and web server for features extraction and selection from protein and peptide sequences. *Bioinformatics* 34, 2499–2502. doi: 10.1093/bioinformatics/bty140
- Cheng, C.-W., Su, E. C.-Y., Hwang, J.-K., Sung, T.-Y., and Hsu, W.-L. (2008). Predicting RNA-binding sites of proteins using support vector machines and evolutionary information. *BMC Bioinform.* 9 (Suppl. 12):S6. doi: 10.1186/1471-2105-9-S12-S6
- Cheng, J., Randall, A. Z., Sweredoski, M. J., and Baldi, P. (2005). SCRATCH: a protein structure and structural feature prediction server. *Nucleic Acids Res.* 33, W72–W76. doi: 10.1093/nar/gki396
- Cheng, N., Li, M., Zhao, L., Zhang, B., Yang, Y., Zheng, C. H., et al. (2020). Comparison and integration of computational methods for deleterious synonymous mutation prediction. *Brief Bioinform.* 21, 970–981. doi: 10.1093/bib/bbz047
- Chou, K.-C. (2001). Prediction of protein cellular attributes using pseudo-amino acid composition. *Proteins* 43, 246–255. doi: 10.1002/prot.1035

- Chou, K.-C. (2004). Using amphiphilic pseudo amino acid composition to predict enzyme subfamily classes. *Bioinformatics* 21, 10–19. doi: 10.1093/bioinformatics/bth466
- Chou, K.-C., and Shen, H.-B. (2007). MemType-2L: a web server for predicting membrane proteins and their types by incorporating evolution information through Pse-PSSM. *Biochem. Biophys. Res. Commun.* 360, 339–345. doi: 10.1016/j.bbrc.2007.06.027
- Chu, Y., Kaushik, A. C., Wang, X., Wang, W., Zhang, Y., Shan, X., et al. (2019). DTI-CDF: a cascade deep forest model towards the prediction of drug-target interactions based on hybrid features. *Brief Bioinform.* doi: 10.1093/bib/bbz152. [Epub ahead of print].
- Cunha, L. D., Ribeiro, J. M., Fernandes, T. D., Massis, L. M., Khoo, C. A., Moffatt, J. H., et al. (2015). Inhibition of inflammasome activation by *Coxiella burnetii* type IV secretion system effector IcaA. *Nat. Commun.* 6:10205. doi: 10.1038/ncomms10205
- Deng, Y., Xu, X., Qiu, Y., Xia, J., Zhang, W., and Liu, S. (2020). A multimodal deep learning framework for predicting drug-drug interaction events. *Bioinformatics*. doi: 10.1093/bioinformatics/btaa501. [Epub ahead of print].
- Dorji, D., Mooi, F., Yantorno, O., Deora, R., Graham, R. M., and Mukkur, T. K. (2018). *Bordetella pertussis* virulence factors in the continuing evolution of whooping cough vaccines for improved performance. *Med. Microbiol. Immunol.* 207, 3–26. doi: 10.1007/s00430-017-0524-z
- Eichinger, V., Nussbaumer, T., Platzer, A., Jehl, M. A., Arnold, R., and Rattei, T. (2016). EffectiveDB—updates and novel features for a better annotation of bacterial secreted proteins and type III, IV, VI secretion systems. *Nucleic Acids Res.* 44, D669–674. doi: 10.1093/nar/gkv1269
- Elbasir, A., Mall, R., Kunji, K., Rawi, R., Islam, Z., Chuang, G.-Y., et al. (2019). BCystal: an interpretable sequence-based protein crystallization predictor. *Bioinformatics* 36, 1429–1438. doi: 10.1093/bioinformatics/btz762
- Esna Ashari, Z., Brayton, K. A., and Broschat, S. L. (2019a). Prediction of T4SS effector proteins for *Anaplasma phagocytophilum* using OPT4e, a new software tool. *Front. Microbiol.* 10:1391. doi: 10.3389/fmicb.2019.01391
- Esna Ashari, Z., Brayton, K. A., and Broschat, S. L. (2019b). Using an optimal set of features with a machine learning-based approach to predict effector proteins for *Legionella pneumophila*. *PLoS ONE* 14:e0202312. doi: 10.1371/journal.pone.0202312
- Esna Ashari, Z., Dasgupta, N., Brayton, K. A., and Broschat, S. L. (2018). An optimal set of features for predicting type IV secretion system effector proteins for a subset of species based on a multi-level feature selection approach. *PLoS ONE* 13:e0197041. doi: 10.1371/journal.pone.0197041
- Govindaraj, R. G., Subramaniam, S., and Manavalan, B. (2020). Extremely-randomized-tree-based prediction of N(6)-methyladenosine sites in *Saccharomyces cerevisiae*. *Curr. Genomics* 21, 26–33. doi: 10.2174/1389202921666200219125625
- Guo, Z., Cheng, X., Hui, X., Shu, X., White, A. P., Hu, Y., et al. (2018). Prediction of new bacterial type III secreted effectors with a recursive hidden markov model profile-alignment strategy. *Curr. Bioinform.* 13, 280–289. doi: 10.2174/1574893612666170725122633
- Hasan, M. M., Schaduagrat, N., Basith, S., Lee, G., Shoombatong, W., and Manavalan, B. (2020). HLPpred-Fuse: improved and robust prediction of hemolytic peptide and its activity by fusing multiple feature representation. *Bioinformatics* 36, 3350–3356. doi: 10.1093/bioinformatics/btaa160
- He, J., Fang, T., Zhang, Z., Huang, B., Zhu, X., and Xiong, Y. (2018). PseUI: pseudouridine sites identification based on RNA sequence information. *BMC Bioinform.* 19:306. doi: 10.1186/s12859-018-2321-0
- Hong, J., Luo, Y., Mou, M., Fu, J., Zhang, Y., Xue, W., et al. (2019). Convolutional neural network-based annotation of bacterial type IV secretion system effectors with enhanced accuracy and reduced false discovery. *Brief Bioinform.* doi: 10.1093/bib/bbz120. [Epub ahead of print].
- Hu, H., Zhang, L., Ai, H., Zhang, H., Fan, Y., Zhao, Q., et al. (2018). HLPi-ensemble: prediction of human lncRNA-protein interactions based on ensemble strategy. *RNA Biol.* 15, 797–806. doi: 10.1080/15476286.2018.1457935
- Huang, Y., Niu, B., Gao, Y., Fu, L., and Li, W. (2010). CD-HIT Suite: a web server for clustering and comparing biological sequences. *Bioinformatics* 26, 680–682. doi: 10.1093/bioinformatics/btq003
- Jeong, J. C., Lin, X., and Chen, X. (2011). On position-specific scoring matrix for protein function prediction. *IEEE/ACM Trans. Comput. Biol. Bioinform.* 8, 308–315. doi: 10.1109/TCBB.2010.93
- Jing, X., and Dong, Q. (2017). MQAPRank: improved global protein model quality assessment by learning-to-rank. *BMC Bioinform.* 18:275. doi: 10.1186/s12859-017-1691-z
- Jing, X., Dong, Q., Lu, R., and Dong, Q. (2019). Protein inter-residue contacts prediction: methods, performances and applications. *Curr. Bioinform.* 14, 178–189. doi: 10.2174/1574893613666181109130430
- Juan, E. Y. T., Li, W. J., Jhang, J. H., and Chiu, C. H. (2009). “Predicting protein subcellular localizations for gram-negative bacteria using DP-PSSM and support vector machines,” in: *2009 International Conference on Complex, Intelligent and Software Intensive Systems* (Fukuoka) 836–841.
- Kang, J., Fang, Y., Yao, P., Li, N., Tang, Q., and Huang, J. (2019). NeuroPP: a tool for the prediction of neuropeptide precursors based on optimal sequence composition. *Interdiscip. Sci.* 11, 108–114. doi: 10.1007/s12539-018-0287-2
- Kira, K., and Rendell, L. A. (1992). “The feature selection problem: traditional methods and a new algorithm. AAAI-92,” in *Proceedings Tenth National Conference on Artificial Intelligence*, 129–134.
- Kuzmanovic, N., Pulawska, J., Hao, L., and Burr, T. J. (2018). The ecology of agrobacterium vitis and management of crown gall disease in vineyards. *Curr. Top. Microbiol. Immunol.* 418, 15–53. doi: 10.1007/82_2018_85
- Li, J., Wei, L., Guo, F., and Zou, Q. (2020). EP3: An ensemble predictor that accurately identifies type III secreted effectors. *Brief Bioinform.* doi: 10.1093/bib/bbaa008. [Epub ahead of print].
- Li, J., Yao, Y., Xu, H. H., Hao, L., Deng, Z., Rajakumar, K., et al. (2015). SecReT6: a web-based resource for type VI secretion systems found in bacteria. *Environ. Microbiol.* 17, 2196–2202. doi: 10.1111/1462-2920.12794
- Li, K., Zhang, S., Yan, D., Bin, Y., Xia, J., and Mitchell, J. C. (2020). Prediction of hot spots in protein-DNA binding interfaces based on supervised isometric feature mapping and extreme gradient boosting. *BMC Bioinform.* 21(Suppl 13):381. doi: 10.1186/s12859-020-03683-3
- Lifshitz, Z., Burstein, D., Peeri, M., Zusman, T., Schwartz, K., Shuman, H. A., et al. (2013). Computational modeling and experimental validation of the *Legionella* and *Coxiella* virulence-related type-IVB secretion signal. *Proc. Natl. Acad. Sci. U.S.A.* 110, E707–715. doi: 10.1073/pnas.1215278110
- Lin, H., Liang, Z. Y., Tang, H., and Chen, W. (2019). Identifying sigma70 promoters with novel pseudo nucleotide composition. *IEEE/ACM Trans. Comput. Biol. Bioinform.* 16, 1316–1321. doi: 10.1109/TCBB.2017.2666141
- Lipovetsky, S., and Conklin, M. (2001). Analysis of regression in game theory approach. *Appl. Stoch. Model. Bus.* 17, 319–330. doi: 10.1002/asmb.446
- Liu, B., Gao, X., and Zhang, H. (2019). BioSeq-Analysis2.0: an updated platform for analyzing DNA, RNA and protein sequences at sequence level and residue level based on machine learning approaches. *Nucleic Acids Res.* 47:e127. doi: 10.1093/nar/gkz740
- Liu, H., Ren, G., Chen, H., Liu, Q., Yang, Y., and Zhao, Q. (2020). Predicting lncRNA-miRNA interactions based on logistic matrix factorization with neighborhood regularized. *Knowl. Based Syst.* 191:105261. doi: 10.1016/j.knosys.2019.105261
- Liu, T., Zheng, X., and Wang, J. (2010). Prediction of protein structural class for low-similarity sequences using support vector machine and PSI-BLAST profile. *Biochimie* 92, 1330–1334. doi: 10.1016/j.biochi.2010.06.013
- Liu, Y., Guo, Y., Wu, W., Xiong, Y., Sun, C., Yuan, L., et al. (2019). A machine learning-based QSAR model for benzimidazole derivatives as corrosion inhibitors by incorporating comprehensive feature selection. *Interdiscip. Sci.* 11, 738–747. doi: 10.1007/s12539-019-00346-7
- Lockwood, S., Voth, D. E., Brayton, K. A., Beare, P. A., Brown, W. C., Heinzen, R. A., et al. (2011). Identification of anaplasma marginale type IV secretion system effector proteins. *PLoS ONE* 6:e27724. doi: 10.1371/journal.pone.0027724
- Lundberg, S., and Lee, S.-I. (2017). “A unified approach to interpreting model predictions,” *Advances in Neural Information Processing Systems 30 (NIPS 2017)*, 4765–4774.
- Lundberg, S. M., Erion, G., Chen, H., DeGrave, A., Prutkin, J. M., Nair, B., et al. (2020). From local explanations to global understanding with explainable AI for trees. *Nat. Mach. Intell.* 2, 56–67. doi: 10.1038/s42256-019-0138-9

- Lundberg, S. M., Nair, B., Vavilala, M. S., Horibe, M., Eisses, M. J., Adams, T., et al. (2018). Explainable machine-learning predictions for the prevention of hypoxaemia during surgery. *Nat. Biomed. Eng.* 2, 749–760. doi: 10.1038/s41551-018-0304-0
- Lv, Z. B., Ao, C. Y., and Zou, Q. (2019). Protein function prediction: from traditional classifier to deep learning. *Proteomics* 19:1900119. doi: 10.1002/pmic.201900119
- Manavalan, B., Basith, S., Shin, T. H., Wei, L., and Lee, G. (2019). AtbPpred: a robust sequence-based prediction of anti-tubercular peptides using extremely randomized trees. *Comput. Struct. Biotechnol. J.* 17, 972–981. doi: 10.1016/j.csbj.2019.06.024
- Manavalan, B., Subramaniam, S., Shin, T. H., Kim, M. O., and Lee, G. (2018). Machine-learning-based prediction of cell-penetrating peptides and their uptake efficiency with improved accuracy. *J. Proteome Res.* 17, 2715–2726. doi: 10.1021/acs.jproteome.8b00148
- Marchesini, M. I., Herrmann, C. K., Salcedo, S. P., Gorvel, J. P., and Comerchi, D. J. (2011). In search of *Brucella abortus* type IV secretion substrates: screening and identification of four proteins translocated into host cells through VirB system. *Cell Microbiol.* 13, 1261–1274. doi: 10.1111/j.1462-5822.2011.01618.x
- Meyer, D. F., Noroy, C., Moumene, A., Raffaele, S., Albina, E., and Vachieri, N. (2013). Searching algorithm for type IV secretion system effectors 1.0: a tool for predicting type IV effectors and exploring their genomic context. *Nucleic Acids Res.* 41, 9218–9229. doi: 10.1093/nar/gkt718
- Noroy, C., Lefrançois, T., and Meyer, D. F. (2019). Searching algorithm for type IV effector proteins (S4TE) 2.0: improved tools for Type IV effector prediction, analysis and comparison in proteobacteria. *PLoS Comput. Biol.* 15:e1006847. doi: 10.1371/journal.pcbi.1006847
- Ren, F., Yang, C., Qiu, Q., Zeng, N., Cai, C., Hou, C., et al. (2019). Exploiting discriminative regions of brain slices based on 2D CNNs for Alzheimer's disease classification. *IEEE Access* 7:181423. doi: 10.1109/ACCESS.2019.2920241
- Sankarasubramanian, J., Vishnu, U. S., Dinakaran, V., Sridhar, J., Gunasekaran, P., and Rajendhran, J. (2016). Computational prediction of secretion systems and secretomes of *Brucella*: identification of novel type IV effectors and their interaction with the host. *Mol. Biosyst.* 12, 178–190. doi: 10.1039/C5MB00607D
- Shan, X., Wang, X., Li, C.-D., Chu, Y., Zhang, Y., Xiong, Y., et al. (2019). Prediction of CYP450 enzyme–substrate selectivity based on the network-based label space division method. *J. Chem. Inf. Model* 59, 4577–4586. doi: 10.1021/acs.jcim.9b00749
- Shannon, C. E. (1948). A mathematical theory of communication. *Bell Syst. Tech. J.* 27, 379–423. doi: 10.1002/j.1538-7305.1948.tb01338.x
- Shen, J., Zhang, J., Luo, X., Zhu, W., Yu, K., Chen, K., et al. (2007). Predicting protein-protein interactions based only on sequences information. *Proc. Natl. Acad. Sci. U.S.A.* 104, 4337–4341. doi: 10.1073/pnas.0607879104
- Shi, F., Yao, Y., Bin, Y., Zheng, C. H., and Xia, J. (2019). Computational identification of deleterious synonymous variants in human genomes using a feature-based approach. *BMC Med. Genomics* 12 (Suppl. 1):12. doi: 10.1186/s12920-018-0455-6
- Su, R., Hu, J., Zou, Q., Manavalan, B., and Wei, L. (2020). Empirical comparison and analysis of web-based cell-penetrating peptide prediction tools. *Brief Bioinform.* 21, 408–420. doi: 10.1093/bib/bby124
- Tang, H., Zhao, Y. W., Zou, P., Zhang, C. M., Chen, R., Huang, P., et al. (2018). HBPred: a tool to identify growth hormone-binding proteins. *Int. J. Biol. Sci.* 14, 957–964. doi: 10.7150/ijbs.24174
- Tomii, K., and Kanehisa, M. (1996). Analysis of amino acid indices and mutation matrices for sequence comparison and structure prediction of proteins. *Protein Eng.* 9, 27–36. doi: 10.1093/protein/9.1.27
- Wang, J., Yang, B., An, Y., Marquez-Lago, T., Leier, A., Wilksch, J., et al. (2019). Systematic analysis and prediction of type IV secreted effector proteins by machine learning approaches. *Brief Bioinform.* 20, 931–951. doi: 10.1093/bib/bbx164
- Wang, J., Yang, B., Revote, J., Leier, A., Marquez-Lago, T. T., Webb, G., et al. (2017). POSSUM: a bioinformatics toolkit for generating numerical sequence feature descriptors based on PSSM profiles. *Bioinformatics* 33, 2756–2758. doi: 10.1093/bioinformatics/btx302
- Wang, X., Wang, Y., Xu, Z., Xiong, Y., and Wei, D. (2019a). ATC-NLSP: prediction of the classes of anatomical therapeutic chemicals using a network-based label space partition method. *Front. Pharmacol.* 10:971. doi: 10.3389/fphar.2019.00971
- Wang, X., Zhu, X., Ye, M., Wang, Y., Li, C.-D., Xiong, Y., et al. (2019b). STS-NLSP: a network-based label space partition method for predicting the specificity of membrane transporter substrates using a hybrid feature of structural and semantic similarity. *Front. Bioeng. Biotechnol.* 7:306. doi: 10.3389/fbioe.2019.00306
- Wang, Y., Guo, Y., Pu, X., and Li, M. (2017). Effective prediction of bacterial type IV secreted effectors by combined features of both C-termini and N-termini. *J. Comput. Aided Mol. Des.* 31, 1029–1038. doi: 10.1007/s10822-017-0080-z
- Wang, Y., Wei, X., Bao, H., and Liu, S. L. (2014). Prediction of bacterial type IV secreted effectors by C-terminal features. *BMC Genomics* 15:50. doi: 10.1186/1471-2164-15-50
- Ward, J. J., McGuffin, L. J., Bryson, K., Buxton, B. F., and Jones, D. T. (2004). The DISOPRED server for the prediction of protein disorder. *Bioinformatics* 20, 2138–2139. doi: 10.1093/bioinformatics/bth195
- Wu, B., Zhang, H., Lin, L., Wang, H., Gao, Y., Zhao, L., et al. (2019). A similarity searching system for biological phenotype images using deep convolutional encoder-decoder architecture. *Curr. Bioinform.* 14, 628–639. doi: 10.2174/1574893614666190204150109
- Xiong, Y., Liu, J., Zhang, W., and Zeng, T. (2012). Prediction of heme binding residues from protein sequences with integrative sequence profiles. *Proteome Sci.* 10 (Suppl. 1):S20. doi: 10.1186/1477-5956-10-S1-S20
- Xiong, Y., Qiao, Y., Kihara, D., Zhang, H. Y., Zhu, X., and Wei, D. Q. (2019). Survey of machine learning techniques for prediction of the isoform specificity of cytochrome P450 Substrates. *Curr. Drug Metab.* 20, 229–235. doi: 10.2174/1389200219666181019094526
- Xiong, Y., Wang, Q., Yang, J., Zhu, X., and Wei, D. Q. (2018). PredT4SE-stack: prediction of bacterial type IV secreted effectors from protein sequences using a stacked ensemble method. *Front. Microbiol.* 9:2571. doi: 10.3389/fmicb.2018.02571
- Xue, L., Tang, B., Chen, W., and Luo, J. S. (2018). A deep learning framework for sequence-based bacteria type IV secreted effectors prediction. *Chemometr. Intell. Lab. Syst.* 183, 134–139. doi: 10.1016/j.chemolab.2018.11.002
- Yan, Z. H., Chen, D., Teng, Z. X., Wang, D. H., and Li, Y. J. (2020). SMOPredT4SE: an effective prediction of bacterial type IV secreted effectors using SVM training with SMO. *IEEE Access* 8, 25570–25578. doi: 10.1109/ACCESS.2020.2971091
- Yang, W., Zhu, X. J., Huang, J., Ding, H., and Lin, H. (2019). A brief survey of machine learning methods in protein sub-golgi localization. *Curr. Bioinform.* 14, 234–240. doi: 10.2174/1574893613666181113131415
- Yu, B., Qiu, W., Chen, C., Ma, A., Jiang, J., Zhou, H., et al. (2019). SubMito-XGBoost: predicting protein submitochondrial localization by fusing multiple feature information and eXtreme gradient boosting. *Bioinformatics* 36, 1074–1081. doi: 10.1093/bioinformatics/btz734
- Yu, L., Sun, X., Tian, S. W., Shi, X. Y., and Yan, Y. L. (2018). Drug and nondrug classification based on deep learning with various feature selection strategies. *Curr. Bioinform.* 13, 253–259. doi: 10.2174/1574893612666170125124538
- Yue, Z., Chu, X., and Xia, J. (2020). PredCID: prediction of driver frameshift indels in human cancer. *Brief Bioinform.* doi: 10.1093/bib/bbaa119. [Epub ahead of print].
- Zalguizuri, A., Caetano-Anolles, G., and Lepek, V. C. (2019). Phylogenetic profiling, an untapped resource for the prediction of secreted proteins and its complementation with sequence-based classifiers in bacterial type III, IV and VI secretion systems. *Brief Bioinform.* 20, 1395–1402. doi: 10.1093/bib/bby009
- Zeng, C., and Zou, L. (2019). An account of *in silico* identification tools of secreted effector proteins in bacteria and future challenges. *Brief Bioinform.* 20, 110–129. doi: 10.1093/bib/bbx078
- Zhang, S., Zhao, L., Zheng, C. H., and Xia, J. (2020). A feature-based approach to predict hot spots in protein-DNA binding interfaces. *Brief Bioinform.* 21, 1038–1046. doi: 10.1093/bib/bbz037
- Zhang, W., Jing, K., Huang, F., Chen, Y., Li, B., Li, J., et al. (2019a). SFLN: a sparse feature learning ensemble method with linear neighborhood regularization for predicting drug–drug interactions. *Inf. Sci.* 497, 189–201. doi: 10.1016/j.ins.2019.05.017

- Zhang, W., Li, Z., Guo, W., Yang, W., and Huang, F. (2019b). "A fast linear neighborhood similarity-based network link inference method to predict microRNA-disease associations," *IEEE/ACM Transactions on Computation Biology and Bioinformatics*.
- Zhang, Y.-F., Wang, X., Kaushik, A. C., Chu, Y., Shan, X., Zhao, M.-Z., et al. (2020). SPVec: a Word2vec-inspired feature representation method for drug-target interaction prediction. *Front. Chem.* 7:895. doi: 10.3389/fchem.2019.00895
- Zhao, Q., Yu, H., Ming, Z., Hu, H., Ren, G., and Liu, H. (2018a). The bipartite network projection-recommended algorithm for predicting long non-coding RNA-protein interactions. *Mol. Ther. Nucleic Acids* 13, 464–471. doi: 10.1016/j.omtn.2018.09.020
- Zhao, Q., Zhang, Y., Hu, H., Ren, G., Zhang, W., and Liu, H. (2018b). IRWNRLPI: integrating random walk and neighborhood regularized logistic matrix factorization for lncRNA-protein interaction prediction. *Front. Genet.* 9:239. doi: 10.3389/fgene.2018.00239
- Zhu, X., He, J., Zhao, S., Tao, W., Xiong, Y., and Bi, S. (2019). A comprehensive comparison and analysis of computational predictors for RNA N6-methyladenosine sites of *Saccharomyces cerevisiae*. *Brief. Funct. Genomics* 18, 367–376. doi: 10.1093/bfgp/elz018
- Zhu, X. J., Feng, C. Q., Lai, H. Y., Chen, W., and Lin, H. (2019). Predicting protein structural classes for low-similarity sequences by evaluating different features. *Knowl-Based Syst.* 163, 787–793. doi: 10.1016/j.knosys.2018.10.007
- Zou, L., and Chen, K. (2016). "Computational prediction of bacterial type IV-B effectors using C-terminal signals and machine learning algorithms," in *2016 IEEE Conference on Computational Intelligence in Bioinformatics and Computational Biology (CIBCB)* (Chiang Mai), 1–5. doi: 10.1109/CIBCB.2016.7758098
- Zou, L., Nan, C., and Hu, F. (2013). Accurate prediction of bacterial type IV secreted effectors using amino acid composition and PSSM profiles. *Bioinformatics* 29, 3135–3142. doi: 10.1093/bioinformatics/btt554
- Zou, Q., Zeng, J., Cao, L., and Ji, R. (2016). A novel features ranking metric with application to scalable visual and bioinformatics data classification. *Neurocomputing* 173, 346–354. doi: 10.1016/j.neucom.2014.12.123

Conflict of Interest: The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Copyright © 2020 Chen, Wang, Chu, Wang, Jiang, Wei and Xiong. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.