



SPVec: A Word2vec-Inspired Feature Representation Method for Drug-Target Interaction Prediction

Yu-Fang Zhang¹, Xiangeng Wang¹, Aman Chandra Kaushik^{1,2}, Yanyi Chu¹, Xiaoqi Shan¹, Ming-Zhu Zhao³, Qin Xu^{1*} and Dong-Qing Wei^{1,4*}

¹ State Key Laboratory of Microbial Metabolism, and SJTU-Yale Joint Center for Biostatistics and Data Science, School of Life Sciences and Biotechnology, and Joint Laboratory of International Cooperation in Metabolic and Developmental Sciences, Ministry of Education, Shanghai Jiao Tong University, Shanghai, China, ² Wuxi School of Medicine, Jiangnan University, Wuxi, China, ³ Instrumental Analysis Center, Shanghai Jiao Tong University, Shanghai, China, ⁴ Peng Cheng Laboratory, Shenzhen, China

OPEN ACCESS

Edited by

Zunnan Huang,
Guangdong Medical University, China

Reviewed by

Francesco Ortuso,
University of Catanzaro, Italy
Ling Wang,
South China University of
Technology, China

Correspondence

Qin Xu
xuqin523@sjtu.edu.cn
Dong-Qing Wei
dqwei@sjtu.edu.cn

Specialty section

This article was submitted to
Medicinal and Pharmaceutical
Chemistry,
a section of the journal
Frontiers in Chemistry

Received: 11 October 2019

Accepted: 12 December 2019

Published: 10 January 2020

Citation

Zhang Y-F, Wang X, Kaushik AC,
Chu Y, Shan X, Zhao M-Z, Xu Q and
Wei D-Q (2020) SPVec: A
Word2vec-Inspired Feature
Representation Method for
Drug-Target Interaction Prediction.
Front. Chem. 7:895.
doi: 10.3389/fchem.2019.00895

Drug discovery is an academical and commercial process of global importance. Accurate identification of drug-target interactions (DTIs) can significantly facilitate the drug discovery process. Compared to the costly, labor-intensive and time-consuming experimental methods, machine learning (ML) plays an ever-increasingly important role in effective, efficient and high-throughput identification of DTIs. However, upstream feature extraction methods require tremendous human resources and expert insights, which limits the application of ML approaches. Inspired by the unsupervised representation learning methods like Word2vec, we here proposed SPVec, a novel way to automatically represent raw data such as SMILES strings and protein sequences into continuous, information-rich and lower-dimensional vectors, so as to avoid the sparseness and bit collisions from the cumbersome manually extracted features. Visualization of SPVec nicely illustrated that the similar compounds or proteins occupy similar vector space, which indicated that SPVec not only encodes compound substructures or protein sequences efficiently, but also implicitly reveals some important biophysical and biochemical patterns. Compared with manually-designed features like MACCS fingerprints and amino acid composition (AAC), SPVec showed better performance with several state-of-art machine learning classifiers such as Gradient Boosting Decision Tree, Random Forest and Deep Neural Network on BindingDB. The performance and robustness of SPVec were also confirmed on independent test sets obtained from DrugBank database. Also, based on the whole DrugBank dataset, we predicted the possibilities of all unlabeled DTIs, where two of the top five predicted novel DTIs were supported by external evidences. These results indicated that SPVec can provide an effective and efficient way to discover reliable DTIs, which would be beneficial for drug reprofiling.

Keywords: drug-target interaction, representation learning, Word2vec, machine learning, feature embedding

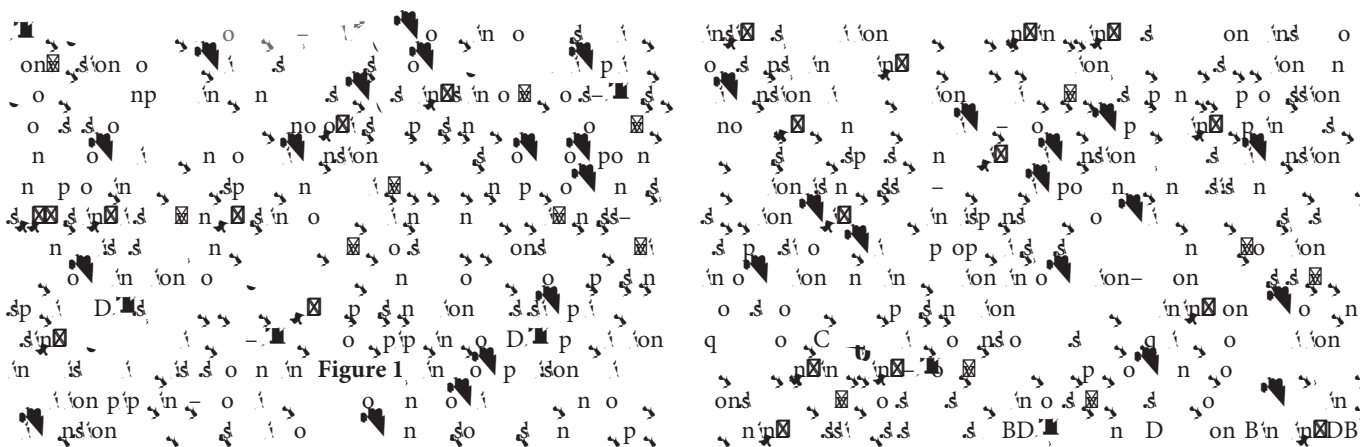


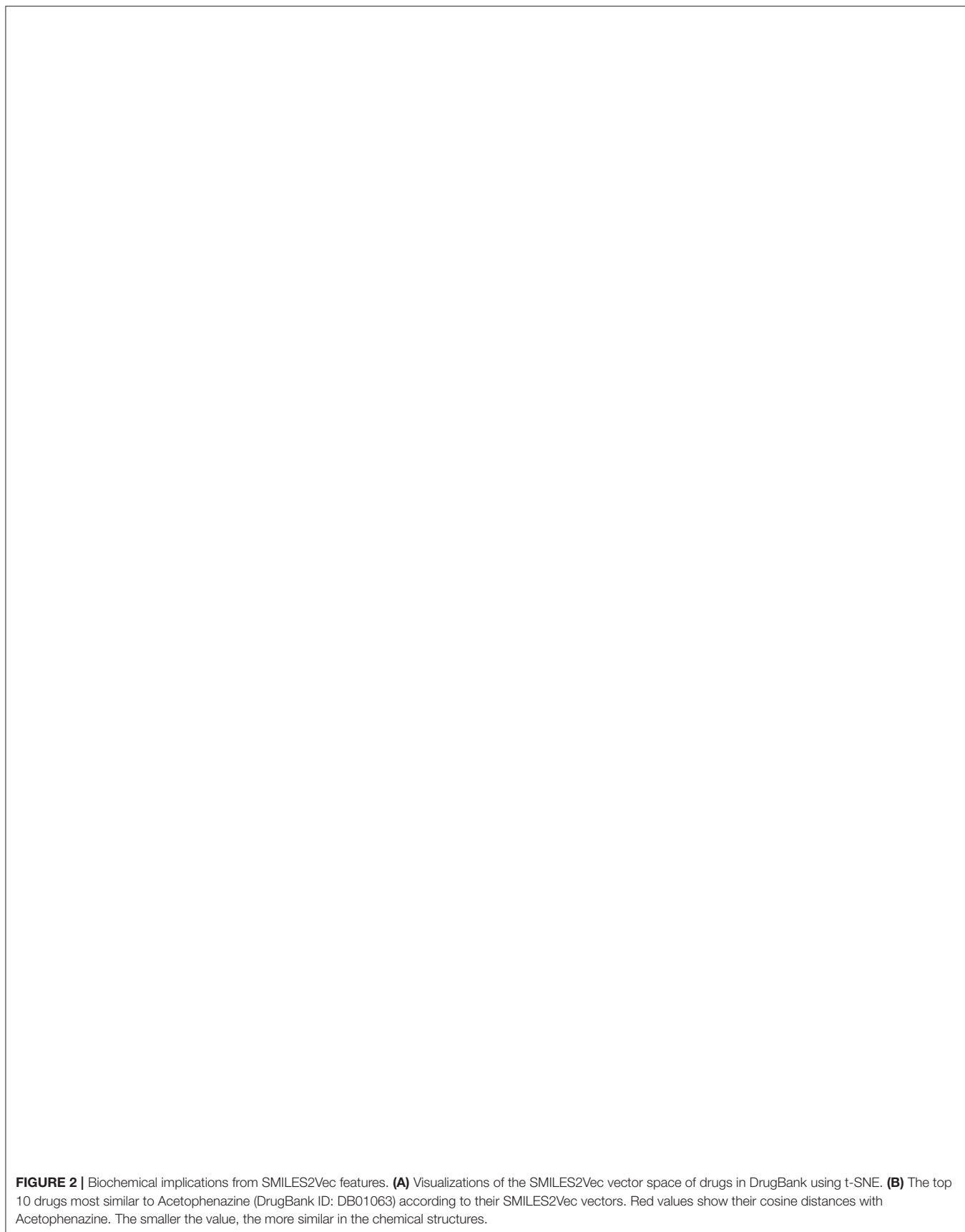
FIGURE 1 | Flowchart of the whole pipeline for DTI prediction in this article (**left**) in comparison to the traditional pipeline (**right**), with the procedures of feature representations squared in dashed lines.

son - p o n n o s n s o
 so on n n n on s n D B n
 -A so p possi o n D s
 n D B n o o op
 p no D s ppo n n n
 n n n s o n D s
 o n o p o n

METHOD

Datasets

Bin DB is p s s s o
 n n n n s o s n on n ons o
 p o n s n o n o B n DB
 on n n n o o p o n n
 o o p on -Cons n
 p o n s n n n n o n s n
 on on s n n > n n
 C n s n o > n n n
 o n n o n on D p s -A on p s
 no n s pas t D s n n n
 s n o n n / n n p s
 n o s n n -
 o D B n - o n n n n
 s p Ap as o
 ons n pas t s s o n n Table 1
 n on p s s o o n o n o
 n s n n n on p s s ons s o
 n o n n on p s s
 ons s o o n n n n on p s
 n on p s - n n o



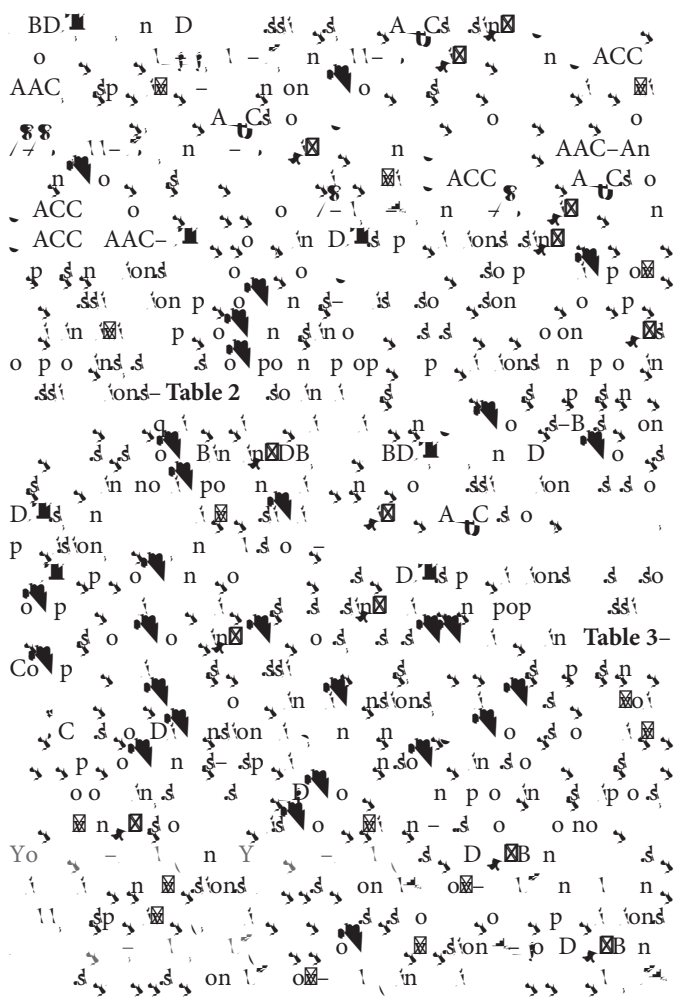
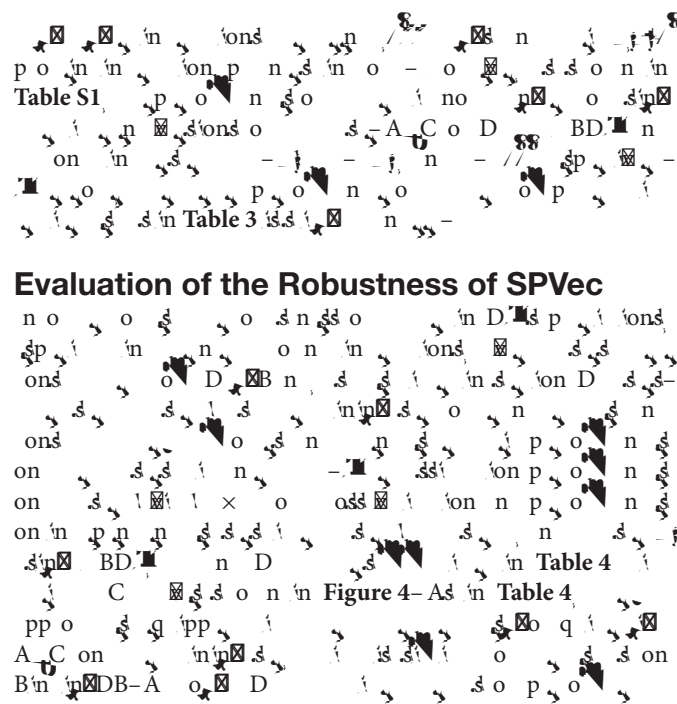


TABLE 1 | Number of entries of the five different datasets obtained from DrugBank dataset.

Datasets	Dataset_1	Dataset_2	Dataset_3	Dataset_4	Dataset_5
Drug	6,068	6,068	537	6,068	537
Target	3,839	3,839	3,839	160	160
Interactions	15,434	3,348	1,735	264	37



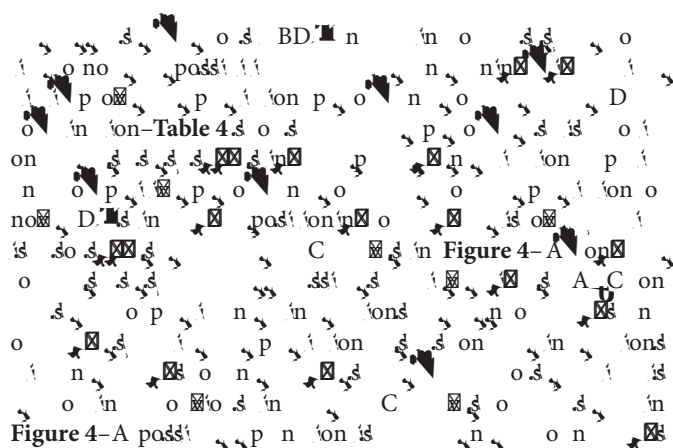


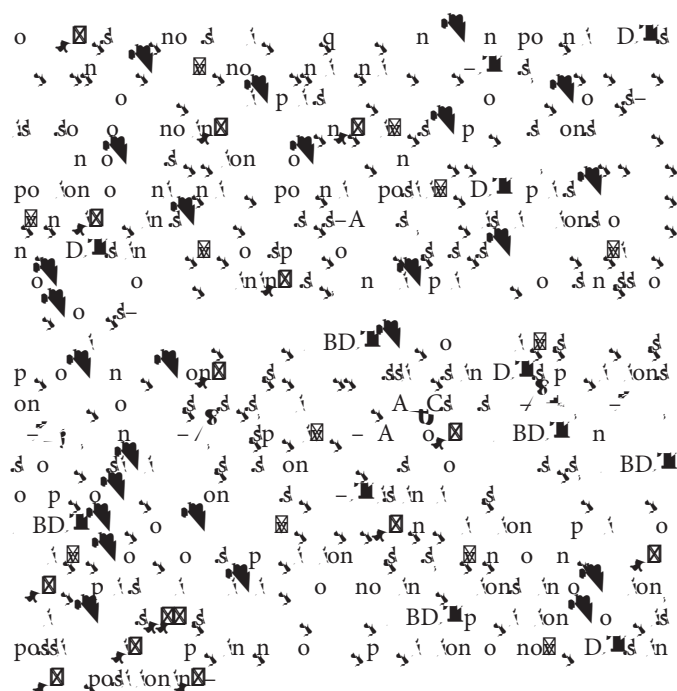
TABLE 2 | Results of classification performance of four feature combinations using three classifiers on BindingDB via 10 × 5-fold cross-validation, with the highest scores highlighted in the bold font.

Feature combinations	Model	AUC	Accuracy	Precision	Recall	F1-score
SPVec (SMILES2Vec-ProtVec)	GBDT	0.9923	0.9680	0.9695	0.9667	0.9681
	RF	0.9927	0.9675	0.9808	0.9540	0.9672
	DNN	0.9617	0.9332	0.9287	0.9248	0.9197
SMILES2Vec-AAC	GBDT	0.9037	0.8272	0.8563	0.7873	0.8204
	RF	0.8770	0.7974	0.8657	0.7050	0.7772
	DNN	0.8708	0.8124	0.7993	0.7879	0.7126
MACCS-ProtVec	GBDT	0.9479	0.8810	0.8908	0.8690	0.8798
	RF	0.9302	0.8542	0.8712	0.8322	0.8512
	DNN	0.9136	0.8034	0.8025	0.8097	0.8074
MACCS-AAC	GBDT	0.8588	0.7811	0.8077	0.7392	0.7719
	RF	0.8360	0.7468	0.8366	0.6150	0.7089
	DNN	0.8451	0.7832	0.7884	0.7726	0.7724

TABLE 3 | AUCs of SPVec and other models on DTI predictions using DrugBank.

Drug features	Drug dim.	Protein features	Protein dim.	ML method	AUC	References
Drug structure information	2,216	AAC, DC ^a and TC ^b	11,943	DNN	0.81	You et al., 2019 ^c
Constitutional, topological and molecular descriptors, 2D autocorrelations, topological charge indices, eigenvalue-based indices	1,664	AAC; DC ^a ; autocorrelation; Composition, Transition, Distribution descriptors; Quasi-sequence-order	1,080	RF	0.8950	Yu et al., 2012 ^c
Constitutional, topological and geometrical descriptors	193	AAC; DC ^a ; autocorrelation; composition, transition and distribution; quasi-sequence-order; amphiphilic pseudo-amino acid composition and total amino acid properties	1,260	DT RF	0.760 0.855	Ezzat et al., 2016
PubChem fingerprints indicating presence or absence of 881 known chemical substructures	881	Fingerprints of 876 different protein domains that are obtained from the Pfam database	876	EnsemDT	0.882	Ezzat et al., 2017
SMILES2Vec	100	ProtVec	100	RF	0.855	This work
				GBDT	0.9467	
				RF	0.9469	
				DNN	0.8637	

^aDC, dipeptide composition; ^bTC, tripeptide composition; ^cThese models are trained on different versions of DrugBank, whose AUCs are only as references.



Prediction and Validation on Unidentified DTIs

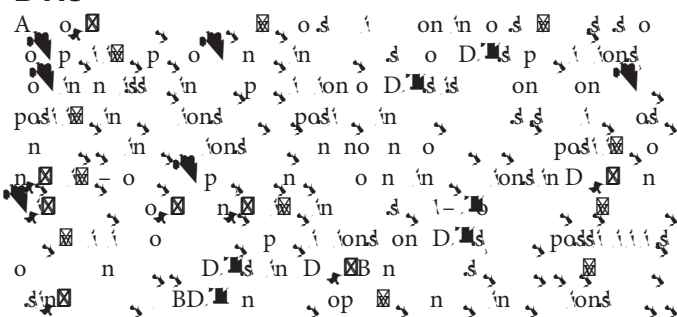


TABLE 4 | Results of classification performance using three classifiers on datasets obtained from DrugBank, with the highest scores highlighted in the bold font.

Dataset	Model	AUC	Accuracy	Precision	Recall	F1-score
Training set		10 × 5-fold cross-validation				
Dataset_1	GBDT	0.9506	0.9323	0.9456	0.9367	0.9343
	RF	0.9557	0.9234	0.9378	0.9369	0.9337
	DNN	0.8952	0.8732	0.8345	0.8437	0.8654
Test sets		Independent validation				
Dataset_2	GBDT	0.8945	0.8628	0.8747	0.8696	0.8637
	RF	0.8930	0.8753	0.8645	0.8467	0.8555
	DNN	0.8201	0.8026	0.8138	0.8199	0.8144
Dataset_3	GBDT	0.7502	0.7389	0.7340	0.7245	0.7333
	RF	0.7448	0.7299	0.7198	0.7243	0.7230
	DNN	0.6999	0.6922	0.6825	0.6798	0.6832
Dataset_4	GBDT	0.7356	0.7223	0.7167	0.7177	0.7201
	RF	0.7235	0.7034	0.7108	0.7078	0.71
	DNN	0.7173	0.6899	0.6884	0.6896	0.6866
Dataset_5	GBDT	0.68	0.6703	0.6679	0.6664	0.6688
	RF	0.5689	0.5605	0.5398	0.5321	0.5411
	DNN	0.6267	0.6098	0.607	0.6122	0.6114



Table 5 –

TABLE 5 | Top five novel DTIs predicted by SPVec-GBDT.

Drug ID	Target ID	Drug name	Target name	Validation source
DB11805	P07947	Saracatinib	The tyrosine-protein kinase Yes	Patel et al., 2013
DB09282	P42262	Molsidomine	Glutamate receptor 2	None
DB05524	Q99640	Pelitinib	Membrane-associated tyrosine and threonine-specific cdc2-inhibitory kinase	https://pubchem.ncbi.nlm.nih.gov/compound/6445562
DB03017	Q16620	Lauric acid	BDNF/NT-3 growth factors receptor	None
DB13165	P11362	Ripasudil	Fibroblast growth factor receptor 1	None

in D D DB W) s n
 o s -n - s so p in on
 n in D D DB n sso
 as n n on n sp in io
 D n on C
 n B is po n o o
 n s i n io o p i o
 sso as n n on n sp in io
 is n s o n o n ons
 p on o no D s

CONCLUSION

Co in n o n o ns
 in o n o n o ns on o s o
 o s o s i o as in o sp
 s n p i s so po n
 op s n o p ns-B on B n DB n
 D B n o s in o s s
 o in n n o s BD n D o
 in D s p on o s n B n DB
 s o n p opas o n

REFERENCES

A. ... - Con n o s s i p s n on
 o o p s q p s o p p o n n o s - PLoS ONE
 B. ... Y- Co A-C n in p p s n on n p
 o i - A -
 C. ... Y- X- X- X- n C o - p p o o in s o
 p on o p o n s o on in o p n q s s g n
 o - J. Cell. Biochem. / - o i - - -

p on p o n n n
 o n on o ACC n AAC-A so
 on D B n n n o pp o p
 BD n n o D s n n o n
 n s i n n o p o n
 A n D s n D B n
 p BD n o n o o op
 p no D s on n i n s o
 o o o o n on
 o s so n n o o n n o
 n on n sp p n n
 o o n s n s n p n
 o o n on o D p on s-

DATA AVAILABILITY STATEMENT

is s n on in ps
 -o- q - -

AUTHOR CONTRIBUTIONS

Y Z X n D on p on n s n
 s -Y Z n X o n o n i s s -Y
 Z A n Y C p o s n s i s - X n Y Z
 o n s p -X on o p o s o
 n s p - ZZ on o p o n s p -A
 o on o n s p s on n p p o
 s on-

FUNDING

is o s s p p o n s o on
 n on on o C n Con nos - V -
 n - A n Y A
 o n i s o n n n o o C n n o n
 n s o n s on n n n o n
 n s i Z V Z D A -

SUPPLEMENTARY MATERIAL

pp n o s n o n
 on n ps - on s n o
 - / - # s pp n

C. ... X- n Y n - D n n on p s on
 n o on s n o - Mol. Biosyst. /
 Co. ... C- B n n n n s n -
 Con on n o p s o p s i p o p
 p on - J. Chem. Inf. Model. /
 Co. ... n on - A n n on n
 p o s s p n n o s n n p - ACM /
 Co. ... n p - Co p s s s n o n p
 o n i s n - Science / - o i - - -

Co-CNN-patch-ppo on o s-Mach. Learn.

D n n on - / - s p -J. Mach. Learn. Res.

B -C n - o D p p n s o n s -J. Chem. Inf. Model.

A - X - n o C - D in on p on s s n n -BMC Bioinf.

A - X - n o C - D in on p on s s n n n n s on -Methods

n - n on pp o n o s p in Ann. Stat.

on A s A o o B n o A C n D C B s in -Nucleic Acids Res. D D

son - B B o n C on - B n DB n p o n is o p on is n s s p o -Nucleic Acids Res.

o -B- o s - C- n is n A- in p n p pas p n n o o p p op -arXiv [Preprint]. A on in ps

o -B- C- is n A- o s - n B C p on p n n o in is no P o n o P op A o -arXiv [Preprint]. A on in ps

o -B- C- is n A- o s - n B C p on p n n o in is no P o n o P op A o -arXiv [Preprint]. A on in ps

o -B- C- is n A- o s - n B C p on p n n o in is no P o n o P op A o -arXiv [Preprint]. A on in ps

o -B- C- is n A- o s - n B C p on p n n o in is no P o n o P op A o -arXiv [Preprint]. A on in ps

o -B- C- is n A- o s - n B C p on p n n o in is no P o n o P op A o -arXiv [Preprint]. A on in ps

o -B- C- is n A- o s - n B C p on p n n o in is no P o n o P op A o -arXiv [Preprint]. A on in ps

o -B- C- is n A- o s - n B C p on p n n o in is no P o n o P op A o -arXiv [Preprint]. A on in ps

o -B- C- is n A- o s - n B C p on p n n o in is no P o n o P op A o -arXiv [Preprint]. A on in ps

o -B- C- is n A- o s - n B C p on p n n o in is no P o n o P op A o -arXiv [Preprint]. A on in ps

o -B- C- is n A- o s - n B C p on p n n o in is no P o n o P op A o -arXiv [Preprint]. A on in ps

o -B- C- is n A- o s - n B C p on p n n o in is no P o n o P op A o -arXiv [Preprint]. A on in ps

o -B- C- is n A- o s - n B C p on p n n o in is no P o n o P op A o -arXiv [Preprint]. A on in ps

o -B- C- is n A- o s - n B C p on p n n o in is no P o n o P op A o -arXiv [Preprint]. A on in ps

o -B- C- is n A- o s - n B C p on p n n o in is no P o n o P op A o -arXiv [Preprint]. A on in ps

o -B- C- is n A- o s - n B C p on p n n o in is no P o n o P op A o -arXiv [Preprint]. A on in ps

o -B- C- is n A- o s - n B C p on p n n o in is no P o n o P op A o -arXiv [Preprint]. A on in ps

o -B- C- is n A- o s - n B C p on p n n o in is no P o n o P op A o -arXiv [Preprint]. A on in ps

o -B- C- is n A- o s - n B C p on p n n o in is no P o n o P op A o -arXiv [Preprint]. A on in ps

o n o s-Mach. Learn.

-J. Mach. Learn. Res.

-J. Chem. Inf. Model.

-BMC Bioinf.

-Methods

Ann. Stat.

-Nucleic Acids Res. D D

-Nucleic Acids Res.

-arXiv [Preprint]. A on in ps

-arXiv [Preprint]. A on in ps

-arXiv [Preprint]. A on in ps

-arXiv [Preprint]. A on in ps

-arXiv [Preprint]. A on in ps

-arXiv [Preprint]. A on in ps

-arXiv [Preprint]. A on in ps

-arXiv [Preprint]. A on in ps

-arXiv [Preprint]. A on in ps

-arXiv [Preprint]. A on in ps

-arXiv [Preprint]. A on in ps

-arXiv [Preprint]. A on in ps

-arXiv [Preprint]. A on in ps

-arXiv [Preprint]. A on in ps

-arXiv [Preprint]. A on in ps

-arXiv [Preprint]. A on in ps

-arXiv [Preprint]. A on in ps

-arXiv [Preprint]. A on in ps

-arXiv [Preprint]. A on in ps