Investigating the mechanisms of emerging and legacy contaminants bioconcentration in fish tissues using machine learning models



AT ARLINGTON

Background & Objective

- Bioconcentration of contaminants in aquatic organ poses environmental and public health risks as to accumulate in fish and transfer through the food
- In this study, we examined different chemical feat representations (i.e., physicochemical properties, fingerprints including ECFP and MACCS, and RD generated molecular descriptors) to classify fish bioconcentration mechanisms.



five repeats. All hyperparameters are fine-tuned usin training, 10% validation, and 20% test data.

Funding: College of Science Research Innovation Award; University of Texas Rising STARs award

Shashwat Dhayade^{1,2}, Dhruvilkumar Ashokbhai Chodvadiya³, Yihan Fan⁴, Feng Gao⁴, Yike Shen¹ ¹Department of Earth and Environmental Sciences, University of Texas at Arlington ²Division of Data Science, University of Texas at Arlington ³Department of Computer Science and Engineering, University of Texas at Arlington ⁴Department of Environmental Health Sciences, University of California Los Angeles ssd6515@mavs.uta.edu

nisms xins	GBDT n	nodel achieved	a high predictio	on perf	ormance	of 0.89	95 for phys	sicocher	nical features	s anc	
chain.					Prooie			abtod)			
ure	INIOGEI	INIQUEI ACCUIACY Recall Precision F1 (We Danal A: Danal A:						gnieu)	RDKit feature		
molecular			$\frac{1}{2} = \frac{1}{2} = \frac{1}$				- 0 0 0 0 1	features car		can a	
Kit-	GDUI	0.095 ± 0.022	0.020 ± 0.0	104	0.904 ± 0	0.033	$0.093 \pm$	0.023	a one-sto	op sc	
		0.900 ± 0.026	0.824 ± 0.0)49)40	0.910 ± 0	0.037	$0.897 \pm$	0.027	prediction	n	
		0.849 ± 0.025	0.801 ± 0.0)40)50	0.764 ± 0	0.024	0.857 ± 0.700	0.025	_		
	KF	0.784 ± 0.039	0.031 ± 0.0	158 2000r0	0.830 ± 0.830	0.059	0.769 ±	0.043	logKOW	V -	
				genera			0 666 1	0 0 2 7	logBCF	F -	
	GDD1 SV/C	0.002 ± 0.033	0.334 ± 0.0		0.093 ± 0	0.007	0.000 ± 0.001	0.037			
		0.599 ± 0.035	0.300 ± 0.0	10 17	0.003 ± 0	0.100	$0.401 \pm$	0.047	SloaP VSA5	* 5 -	
		0.279 ± 0.100	0.333 ± 3.33		0.093 ± 0	0.035	$0.134 \pm$	0.007	VSA EState1	1 -	
	КГ	0.090 ± 0.039	0.339 ± 0.0	JD4	0.721 ± 0		$0.070 \pm$	0.042	 LabuteASA	Α-	
					$\frac{1}{2}$ + $\frac{10}{2}$	OW + IC			MaxAbsPartialCharge	e -	
	GBDT	0.863 ± 0.024	0.767 ± 0.0		0.809 ± 0	0.048	$0.858 \pm$	0.026	PEOE_VSA6	5 -	
		0.599 ± 0.036	$0.377 \pm 0.020 \qquad 0.379 \pm 0.020 = 0.020 \pm 0.02$		0.579 ± 0.00	0.183	0.474 ± 0.040		qec	- L	
ure Addition [.]		0.239 ± 0.110	0.333 ± 5.55		0.080 ± 0	0.039	0.108 ± 0.000	0.069	_	0.00	
BCF, logKOW	$\begin{array}{c c c c c c c c c c c c c c c c c c c $										
tures: 200	IVIOIECU	iar ingerprints a	are less compe		nan pnysi				s dalasel		
	Model	Accuracy		Pre	cision	F1 (V	/eighted)		С	omin	
		0.611 ± 0.026	Paneid: ECFP 0 472 ± 0.052			0 601					
	SVC	0.044 ± 0.030 0.654 + 0.032	0.472 ± 0.052 0.470 + 0.042	0.041	± 0.100 + 0.138	0.001	± 0.040				
		0.034 ± 0.032 0.602 + 0.041	0.470 ± 0.042 0.527 + 0.046	0.525	+ 0.130	0.020	3 ± 0.033				
RF	RF	0.644 ± 0.039	0.442 ± 0.044	0.681	± 0.128	0.583	3 ± 0.053				
	MLP	0.622 ± 0.036	0.457 ± 0.078	0.516	± 0.103	0.580	$) \pm 0.057$	GBDT:	gradient boo	sted	
	PanelE: MACCS features							decisio	n tree; SVC:		
00-0-00	GBDT	0.632 ± 0.034	0.469 ± 0.039	0.652	± 0.102	0.603	3 ± 0.040	suppor	t vector class	sifier;	
LR	SVC	0.630 ± 0.043	0.465 ± 0.049	0.571	± 0.109	0.606	5 ± 0.043	LR: log	istic regressi	on; F	
nisms		0.526 ± 0.042	0.533 ± 0.068	0.470	± 0.045	0.541	± 0.039	random	n forest; MLP	': Mu	
rating Chamicala		0.648 ± 0.038	0.463 ± 0.048	0.629	± 0.135	0.613	3 ± 0.047	layer p	erceptron.		
ng Chemicals	Conclusions										
	 Accurate prediction of bioconcentration mechanisms is important for chemical risk assessme 										
dation and	 Our GBDT model achieved a high prediction performance of 0.895 ± 0.022 for physicochemic 										
ig 70%	gene	rated features p	lus logKow an	d logB	CF.				-		
	Perfc	ormance varied	with feature rep	presen	tations, h	ighligh	ting the in	nportanc	e of selecting	g fea	

Results and Discussion

10.863 for RDKit-generated features plus class impact on difficult samples.

ent. ical features and 0.863 ± 0.024 for RDKit-

ture sets for each dataset and model.



es combined with key physicochemical assist users and enable the development of olution for bioconcentration mechanism





ng soon!



April 2025, unpublished results, please do not copy