

Increasing the performance, trustworthiness and practical value of machine learning models: a case study predicting hydrogen bond network dimensionalities from molecular diagrams

Supplementary Information

A. List of molecular descriptors calculated using RDkit package. Marked in bold are the two descriptors that were discarded during the data pre-processing stage, because their values were constant across all dataset examples.

MolWt	SMR_VSA10	Chi2v
ExactMolWt	SMR_VSA2	Chi3n
HeavyAtomMolWt	SMR_VSA3	Chi3v
NumHAcceptors	SMR_VSA4	Chi4n
NumHDonors	SMR_VSA5	Chi4v
NumRotatableBonds	SMR_VSA6	HallKierAlpha
FractionCSP3	SMR_VSA7	Kappa1
NHOHCount	SMR_VSA8	Kappa2
NOCount	SMR_VSA9	Kappa3
NumHeteroatoms	SlogP_VSA1	BertzCT
HeavyAtomCount	SlogP_VSA10	TPSA
NumAliphaticCarbocycles	SlogP_VSA11	Ipc
NumAromaticCarbocycles	SlogP_VSA12	FpDensityMorgan1
NumSaturatedCarbocycles	SlogP_VSA2	FpDensityMorgan2
NumAliphaticHeterocycles	SlogP_VSA3	FpDensityMorgan3
NumAromaticHeterocycles	SlogP_VSA4	
NumSaturatedHeterocycles	SlogP_VSA5	
RingCount	SlogP_VSA6	
NumAliphaticRings	SlogP_VSA7	
NumAromaticRings	SlogP_VSA8	
NumSaturatedRings	SlogP_VSA9	
NumValenceElectrons	EState_VSA1	
NumRadicalElectrons	EState_VSA2	
qed	EState_VSA3	
MolLogP	EState_VSA4	
MolMR	EState_VSA5	
BalabanJ	EState_VSA6	
MaxAbsPartialCharge	EState_VSA7	
MinAbsPartialCharge	EState_VSA8	
MaxPartialCharge	EState_VSA9	
MinPartialCharge	EState_VSA10	
MaxAbsEStateIndex	EState_VSA11	
MinAbsEStateIndex	VSA_EState1	
MaxEStateIndex	VSA_EState2	
MinEStateIndex	VSA_EState3	
PEOE_VSA1	VSA_EState4	
PEOE_VSA10	VSA_EState5	
PEOE_VSA11	VSA_EState6	
PEOE_VSA12	VSA_EState7	
PEOE_VSA13	VSA_EState8	
PEOE_VSA14	VSA_EState9	
PEOE_VSA2	VSA_EState10	
PEOE_VSA3	LabuteASA	
PEOE_VSA4	Chi0	
PEOE_VSA5	Chi0n	
PEOE_VSA6	Chi0v	
PEOE_VSA7	Chi1	
PEOE_VSA8	Chi1n	
PEOE_VSA9	Chi1v	
SMR_VSA1	Chi2n	

B. Statistical methods considered for model generation, along with optimised parameters and mean 5-fold cross validation accuracy.

Method	Optimal Parameters	Accuracy
Random Forest	max_depth: 30, n_estimators: 100	0.57
SVM Linear kernel	C: 50	0.57
SVM RBF kernel	C: 50, gamma: 0.001	0.58
SVM Sigmoid kernel	gamma: 0.001	0.53
Logistic Regression (one vs rest)	C: 200	0.54
Logistic Regression multinomial	C: 100	0.56
K Nearest Neighbours	n_neighbors: 50	0.51
Gaussian Naive Bayes	var_smoothing: 0.01	0.43
Gaussian Process (one vs rest)	random_state: 0	0.47
Gradient Boosting Classifier	learning_rate: 0.1, max_depth: 10, n_estimators: 100	0.57