Predicting Crystal Packing of Olanzapine Solvates using Random Forest

Rajni M. Bhardwaj,^a Susan Reutzel-Edens,^b Blair F. Johnston, ^a and Alastair J. Florence ^{a*}

^aStrathclyde Institute of Pharmacy and Biomedical Sciences, University of Strathclyde, 161 Cathedral Street, Glasgow G4 0RE, United Kingdom

^b Eli Lilly and Company, Indianapolis, Indiana 46285, United States

* To whom correspondence should be addressed. E-mail: alastair.florence@strath.ac.uk

Table of Contents

Contents

1.	Schematic Workflow of Random Forest	2
2.	Crystal Packing in solvates based on SC_{31} , SC_{32} and SC_{33}	3
3.	Preparation of the Dataset	3
4.	Model Building	6
5.	Important Descriptors Assessment	8
6.	References	9

1. Schematic Workflow of Random Forest



Random Forest ---- Building the Model Forest

Fig. 1. Schematic workflow of building of Random Forests model.

Bootstrap sampling and random selection of input descriptors are used to induce randomness in the input data used to develop the RF model. This ensures that the classification trees grown in the forest are dissimilar and not correlated to each other. Using bootstrap sampling, classification trees are grown using 2/3rd of the dataset and remaining 1/3rd of the dataset [Out Of Bag (OOB) data] is employed to obtain unbiased estimates of correct classification rates (internal estimates of error). Compared to a single classification tree, this algorithm yields better prediction rates and is more robust in dealing with noise in the data set because the forest of trees are grown to the full extent. The generalisation error of a forest of trees classifiers depends on the strength of the individual tree in the forest and the correlation between them.

2. Crystal Packing in solvates based on SC₃₁, SC₃₂ and SC₃₃



Fig. 2. Representative crystal structures of olanzapine solvates based on (a) SC_{31} - olanzapine-1,4dioxane solvate hydrate; (b) SC_{32} -Olanzapine-acetic acid solvate and (c) SC_{33} -Olanzapine-TBMe hydrate; showing the position of solvent molecules within the crystal lattice

3. Preparation of the Dataset

All the solvent molecules were drawn using Chemdraw Ultra (version 11.0) and reliable 3-D conformations generated in Discovery Studio using the Pipeline Pilot interface (Accelrys, 2010). 2-D (185) and internal 3-D (i3-D) (123) molecular descriptors were calculated using MOE.¹ 2-D molecular descriptors are defined to be numerical properties and calculated from the atoms and connection table of the molecule. 3-D molecular descriptors can be classified in two categories: one that depend on internal coordinates only and 2nd that depend on absolute orientation of molecule. A brief explanation of the list of calculated 2- and 3-D molecular descriptors which were used to model the solvent library are given in Table 1 (Source: http://www.chemcomp.com/journal/descr.htm).

Table 1. Molecular descriptors and brief explanation that were calculated for solvent molecules.

Descriptors	Category						
2-D descriptors							
apol, bpol, Fcharge, mr, SMR, Weight,	physical	Physical properties are calculated					
logP (o/w), SlogP, vdw_vol, density,	properties	from the connection table of a					
vdw-area		molecule					
SlogP_VSA0-SlogP_VSA9,	subdivided	The Subdivided Surface Areas are					
SMR_VSA0 - SMR_VSA7	surface areas	descriptors based on an approximate					
		accessible van der Waals surface area					
		calculation for each atom, v_i along					
		with some other atomic property, p_i .					
a_aro, a_count, a_heavy, a_ICM, a_IC,	atom count and	The atom count and bond count					
a_nH, a_nB, a_nC, a_nN, a_nO, a_nF,	bond count	descriptors are functions of the counts					
a_nP, a_nS, a_nCl, a_nBr, a_nl,		of atoms and bonds					
D_{1101N} , D_{1101K} , D_{a1} , D_{c0011} , b double b beauty b rot b rot P							
b single b triple VAdiMa VAdiFa							
chi0 chi0 C chi1 chi1 C chi0y	Kier&Hall	The Kier and Hall kappa molecular					
chi0y C chi1y chi1y C Kier1 -	Connectivity	shape indices compare the molecular					
Kier3. KierA1 - KierA3. KierFlex.	and Kappa	graph with minimal and maximal					
zagreb	Shape Indices	molecular graphs, and are intended to					
	1	capture different aspects of molecular					
		shape.					
balabanJ, diameter, petitjean, radius,	Adjacency and	The adjacency matrix, M, of a					
VDistEq, VDistMa, weinerPath,	Distance Matrix	chemical structure is defined by the					
weinerPol	Descriptors	elements [Mij] where Mij is 1 if					
		atoms i and j are bonded and zero					
		otherwise. The distance matrix, D, of					
		a chemical structure is defined by the					
		of the shortest path from atoms i to i					
		zero is used if atoms i and i are not					
		part of the same connected					
		component.					
a acc, a acid, a base, a don, a hyd,	Pharmacophore	The Pharmacophore Atom Type					
vsa_acc, vsa_acid, vsa_base, vsa_don,	Feature	descriptors consider only the heavy					
vsa_hyd, vsa_other, vsa_pol	Descriptors	atoms of a molecule and assign a type					
		to each atom					
Q_PC+ PEOE_PC+, Q_PC-	Partial Charge	Descriptors that depend on the partial					
PEOE_PC-, Q_RPC+ PEOE_RPC+,	Descriptors	charge of each atom of a chemical					
Q_RPC- PEOE_RPC-, Q_VSA_POS		structure require calculation of those					
PEOE_VSA_POS, Q_VSA_NEG,		partial charges.					
PEOE_VSA_NEG, Q_VSA_PPOS, DEOE_VSA_DDOSQVSA_DNEG							
PEOE VSA PNEG O VSA HVD							
PEOF VSA HVD O VSA POL							
PEOE VSA POL. O VSA FPOS							
PEOE VSA FPOS, O VSA FNEG							
PEOE VSA FNEG, Q VSA FPPOS							
PEOE_VSA_FPPOS, Q_VSA_FPNEG							
$ PEOE_VSA_FPNEG, \overline{Q}_VS\overline{A}_FHYD $							
PEOE_VSA_FHYD, Q_VSA_FPOL							
PEOE_VSA_FPOL, PEOE_VSA+6 -							
PEOE_VSA+0, PEOE_VSA-0 -							
reue_vsa-o 3 D Descriptors							
S-D Descriptors							

E, E_ang, E_ele, E_nb, E_oop, E_sol,	Potential Energy	The energy descriptors use the MOE	
E_stb, E_str, E_strain, E_tor, E_vdw,	Descriptors	potential energy model to calculate	
E_rele, E_rsol, E_rvdw	_	energetic quantities from stored 3D	
		conformations.	
ASA, dens, glob, pmi, pmiX, pmiY,	Surface Area,	Descriptors depend on the structure	
pmiZ, rgyr, std_dim1 -std_dim3, vol,	Volume and	connectivity and conformation	
VSA	Shape		
	Descriptors		
ASA+, ASA-, ASA_H, ASA_P, DASA,	Conformation	Descriptors depend upon the stored	
CASA+, CASA-, DCASA, dipole,	Dependent	partial charges of the molecules and	
diploeX, dipole, dipoleZ, FASA+,	Charge	their conformations.	
FASA-, FCASA+, FCASA-,	Descriptors		
FCASA H, FCASA P	_		

A correlation matrix was prepared using a Pearson correlation coefficient by using a Pipeline Pilot interface. Molecular descriptors which showed zero variance and covariance (threshold of Pearson correlation coefficient >90%) were removed from the dataset. The resultant dataset comprised of 151 calculated molecular descriptors.

4. Model Building

Out of 35 OZPN solvates, 28 are based on a common 2-D SC_{21} and differ only in 3-D originated by different stacking of SC_{21} sheets. A RF model was built to confirm the role of solvent molecules in directing the crystal packings of these 28 OZPN solvates in 3-D.

25 out of 28 solvates were taken for random forest classification as there were no calculated descriptors for 3 solvents (1,4-butanediol, t-butanol and ethylene glycol (EG). OZPN forms two dihydrates i.e. dihydrate B and dihydrate E which fall in class 2 and class 1 respectively. Control over hydrate (dihydrate B vs. dihydrate E) formation can be achieved via appropriate crystallisation conditions. For this classification, hydrate was considered only in one class.

The error plot in Fig. 3 provides an overall OOB error of prediction and prediction accuracy for the classification model and a confusion matrix in Table 2 provides information on the prediction accuracy and OOB error rate associated with each class.



Fig. 3. Error Plot of Random Forest classification error for molecular packing types of olanzapine solvates. The black line represents the error associated with the overall model. Red, blue and green dotted lines represent the error associated with the molecular packing type classes; SC_{31} , SC_{32} and SC_{33} .

Table 2. Confusion matrix generated from Random Forest classification for the three molecular packing types of olanzapine solvates.

Predicted \rightarrow	Class 1	Class 2	Class 3 (SC ₃₃)	Class Error
Observed 🗸	(30 ₃₁)	(30 ₃₂)		(70)
Class 1(SC ₃₁)	13	4	0	23.5
Class 2 (SC ₃₂)	2	2	0	50
Class 3 (SC ₃₃)	0	0	4	0

OZPN solvates based on SC_{31} and SC_{33} have water molecules in the asymmetric unit and so the effect of incorporation of water molecules on the predictive capability of the RF classification model was also investigated. However the same results were obtained after taking account of water as a categorical descriptor (present/not present) and even by taking it as a numerical descriptor whose value was equivalent to the number of water molecules in the asymmetric unit of the respective OZPN solvate.

5. Important Descriptors Assessment

The RF algorithm also assesses the importance of descriptors used in building of the classification model. It is assessed by replacing each descriptor in turn by random noise and the resulting deterioration in the model quality is a measure of descriptor importance. The deterioration in the RF model quality is assessed by mean decrease in accuracy (based on OOB data). Variable dependency plots obtained for all three molecular packing type classes are shown in Fig. 4. Numerical values of descriptors are provided in Table 1 of the main article.



Fig. 4. Variable dependency plots (SMR, apol, VdW_vol and b_count) for molecular packing types, SC_{31} , SC_{32} and SC_{33} .

6. References

 MOE, 2002, Chemical Computing Group, 1010 Sherbrooke St. W, Montreal, Quebec, H3A 2R7.