Automatic Model Generation

Joelle Gola, Olga Obrezanova

19 June 2007



Copyright © 2006 Galapagos NV



- Background
- Process
- Techniques
- Illustrative examples
- Next steps





Application of ADME Predictions







Quality Volume Descriptors Set split Techniques Validation Model selection

Uncertainty Chemical space Interpretation





Objectives

- Simplify and speed up the model building process
- Understand and make use of the model output
 - Compound in, prediction out....
- Flexible approach
 - Expert vs. non expert
- Facilitate integration of the model with the decision making process
 - Chemical space analysis
 - Probabilistic scoring











Checking input data







- ~330 2D SMARTS based descriptors
- Input own descriptors (additional columns)
- User defined 2D SMARTS based descriptors







- Less than 4% occurrences
- Standard deviation: 0.0005
- Pair-wise correlated descriptors:
 - \succ Correlation >=95%







- Require 3 sets:
 - > Training: Building models
 - > Validation: Selecting best model
 - Fest: Checking best model
- Clustering based method:
 - Structural fingerprint
 - Tanimoto level: 0-1
 - \succ % compounds in the training set
- 3 user defined sets





Classification models:
 Decision Tree

- Continuous models:
 - Partial Least Squares
 - Gaussian Processes
 - Radial Basis Functions
- Automatically run all the appropriate models
- User selects preferred techniques



- Compare statistical results on validation set:
 - Set of rules to automatically select the best model
 - User choice







- Uncertainty in prediction
- Glowing Molecule



Techniques: Classification Models

- Recursive partitioning approach to building classification models
- Based on C4.5 approach developed by Ross Quinlan

Rules:

- Generate up to 20 different models by automatically varying method settings:
 - Simple decision trees with various stopping conditions
 - Pruned decision trees
 - Rule sets built from decision trees





Techniques: Continuous Models

- Suite of modeling techniques include:
 - Partial Least Squares
 - Gaussian Processes
 - Radial Basis Functions





Techniques: Gaussian Processes

• Powerful machine learning technique based on a Bayesian statistical approach

> Bayes Theorem:

$$P(f|Y,X) \propto P(Y|f,X)P(f)$$

Posterior distribution

Prior distribution

Gaussian process defines the distribution over functions



Techniques: Gaussian Processes

- Automatic determination of the model parameters:
 - Inherent ability to select relevant descriptors
- Implemented 5 techniques*:

Techniques	Descriptor Selection	
GP-Fixed	-	
GP-2DSearch	_	Increasing computational demand
GP-FVS	\checkmark	
GP-RFVS	\checkmark	
GP-OPT	\checkmark	

* Obrezanova et al., J.Chem.Inf.Model., accepted



Techniques: Radial Basis Functions

For given Y and x, RBF is used to find function f(x)
 f(x) = Y + noise

• f(x) chosen as:

$$f(x) = \sum_{i=1}^{N} a_i \|x - x^{(i)}\|$$

• RBF requires f(x) to pass through all training points:

$$y_j = \sum_{i=1}^N a_i \| x^{(j)} - x^{(i)} \|, \quad j = 1...N$$

weights a_i can be found from this linear system of equations



Techniques: Radial Basis Functions

- Good method for small or large data sets
 - > But sensitive to noise created by excessive descriptors

Solution:

- Coupled with a genetic algorithm, GA-RBF
- Automatically select RBF or GA-RBF:

# Compounds per descriptor	GA-RBF	RBF
< 5	√	-
≥ 5	-	√
Bringing balance to optimisation		BioFocus DP

Techniques: Performance Measures

Classification models

 $Accuracy = \frac{Correctly \ predicted \ compounds}{Total \ number \ of \ compounds}$

 $Kappa = \frac{(Observed agreement - Chance agreement)}{(Total - Observed agreement)}$

Continuous models

$$R^{2} = 1 - \frac{\sum_{i}^{i} (y_{i}^{pred} - y_{i}^{obs})^{2}}{\sum_{i}^{i} (y_{i}^{obs} - \overline{y}^{obs})^{2}} \qquad RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^{N} (y_{i}^{pred} - y_{i}^{obs})^{2}}$$



Techniques: Choosing Best Model

 Selection of the best model is done by its performance on the validation set.

Rules:

- Classification model
 - > The higher Kappa statistic (0<Kappa< 1)
 - If identical Kappa statistics then choose model with best statistical results on training set
- Continuous model
 - Smaller RMSE





Techniques: Uncertainties in Prediction



Descriptor 1

 The position of compounds relative to the chemical space of the model is reflected in the reported standard deviation (SD)



Techniques: Interpretation

• Molecule in, prediction out...

"Why did *this* model predict *that* value for *my* molecule?" "What can I do to my molecule to improve the prediction?"

- Models already encode some of these answers!
 - Reveal non-linear as well as linear relationships between property being modelled and the selected descriptors
 - Proprietary algorithm: Glowing Molecule





Illustrative Examples

- Example 1: BBB classification model
- Example 2: Continuous hERG model
- Example 3: Building a local model





• BBB± data set from a CNS library published by Zhao et al. Originally prepared by Adenot and Lahana

- > J. Chem. Inf. Model, 2007,47,170
- 1593 compounds

Chemically Diverse Set







	Training (1273 cpds-112 Descriptors)			Vali	dation	Test	
				(158 cpds)		(160 cpds)	
Models	# Descriptors	Карра	Accuracy	Карра	Accuracy	Карра	Accuracy
DT1	7	0.78	0.94	0.88	0.96		
DT2	5	0.77	0.94	0.88	0.96		
DT4	6	0.77	0.94	0.88	0.96		
DT7	1	0.63	0.91	0.70	0.90		
DT13	13	0.87	0.96	0.91	0.97		
DT14	15	0.89	0.97	0.90	0.96		
DT15	15	0.86	0.96	0.95	0.98		
DT16	20	0.84	0.96	0.91	0.97		
DT17	13	0.84	0.96	0.90	0.96		
DT19	28	0.88	0.97	0.92	0.97		
DT20	26	0.86	0.96	0.91	0.97		



	Training (1273 cpds-112 Descriptors)			Vali	dation	Т	est
				(158 cpds)		(160 cpds)	
Models	# Descriptors	Карра	Accuracy	Карра	Accuracy	Карра	Accuracy
DT1	7	0.78	0.94	0.88	0.96		
DT2	5	0.77	0.94	0.88	0.96		
DT4	6	0.77	0.94	0.88	0.96		
DT7	1	0.63	0.91	0.70	0.90		
DT13	13	0.87	0.96	0.91	0.97		
DT14	15	0.89	0.97	0.90	0.96		
DT15	15	0.86	0.96	0.95	0.98		
DT16	20	0.84	0.96	0.91	0.97		
DT17	13	0.84	0.96	0.90	0.96		
DT19	28	0.88	0.97	0.92	0.97		
DT20	26	0.86	0.96	0.91	0.97		



	Training (1273 cpds-112 Descriptors)			Vali	dation	Test	
				(158 cpds)		(160 cpds)	
Models	# Descriptors	Карра	Accuracy	Карра	Accuracy	Карра	Accuracy
DT1	7	0.78	0.94	0.88	0.96		
DT2	5	0.77	0.94	0.88	0.96		
DT4	6	0.77	0.94	0.88	0.96		
DT7	1	0.63	0.91	0.70	0.90		
DT13	13	0.87	0.96	0.91	0.97		
DT14	15	0.89	0.97	0.90	0.96		
DT15	15	0.86	0.96	0.95	0.98	0.95	0.98
DT16	20	0.84	0.96	0.91	0.97		
DT17	13	0.84	0.96	0.90	0.96		
DT19	28	0.88	0.97	0.92	0.97		
DT20	26	0.86	0.96	0.91	0.97		



Validat	ion set	Pred	icted
Vandation Sec		BBB -	BBB +
Obs.	BBB -	38	3
	BBB +	0	117

Test set		Predicted			
		BBB -	BBB +		
Obs.	BBB -	38	1		
	BBB +	2	119		



Rule 1:

If PSA > 108.8 and q137 > 26 and calclogP <= 3.5 then

BBB - confidence = 0.98

Rule 2:

If thioEther = 0 and PSA = < 128.2 and ed70 = 0 and Negative Charge = 0 then

BBB+ confidence = 0.97



Example 2: Continuous hERG Model

• 177 pIC_{50} values from the literature

> Patch clamp measurements in mammalian cells – mainly HEK293

	Training			Validation		Test	
	(124 cpds-14	13 desc	riptors)	(27	cpds)	(26 cpds)	
Models	# Descriptors	R ²	RMSE	R ²	RMSE	R ²	RMSE
PLS	143	0.69	0.770	0.61	0.953		
GPFixed	143	0.85	0.538	0.66	0.883		
GP2DSearch	143	0.83	0.576	0.65	0.892		
GPFVS	41	0.78	0.661	0.68	0.864	0.71	0.822
GPRFVS	19	0.81	0.611	0.65	0.897		
GPOPT	42	0.84	0.558	0.63	0.925		
GA-RBF	23	1	0	0.64	0.910		



Example 2: Continuous hERG Model



Example 2: Continuous hERG Model

- Important chemical features for hERG binding:
 - ≻Lipophilicity
 - ➢Negative charge
 - Positively charged nitrogen at pH 7.4
 - Aromaticity index
- Redesign using Glowing Molecule







Example 3: Building a Local Model

- Austin et al. measured logD values at pH 7.4 for 78 compounds
 - > J Med Chem. 2003 Jul 17;46(15):3210-20
- QSAR model highlighted the importance of lipophilicity and ionization in controlling beta(2) duration
- Wants to predict logD7.4 values for new compounds





Example 3: Global Model

• Current global model was built on 1044 compounds

logD values measured at pH 7.4 were extracted from StARLITe[™]



Training Set: GPRFVS Model, 52 descriptors $R^2 = 0.91 RMSE = 0.57$

Validation set: $R^2 = 0.82$ RMSE = 0.75 Test set: $R^2 = 0.86$ RMSE = 0.68



Example 3: Applying Global Model

Results for the 78 Austin compounds:



Observed logD7.4

Bringing balance to optimisation



Example 3: Building a Local Model

Build a model on the 78 Austin compounds.

	Training		Validation		Test		
	(54 cpds – 7	2 descri	iptors)	(12 cpds)		(12 cpds)	
Models	# Descriptors	R ²	RMSE	R ²	RMSE	R ²	RMSE
PLS	72	0.44	0.54	0.37	0.57		
GPFixed	72	0.70	0.39	0.32	0.60		
GP2DSearch	72	0.86	0.26	0.39	0.57		
GPFVS	22	0.82	0.31	0.55	0.49	0.62	0.45
GPRFVS	6	0.81	0.32	0.46	0.54		
GPOPT	18	0.87	0.26	0.47	0.53		
GA-RBF	7	1	0	-0.13	0.78		





Example 3: Building a Local Model



Validation set: R² = 0.55 **RMSE = 0.49**

Test set: R² = 0.62 RMSE = 0.45



» Next Steps

- New modeling techniques:
 - Gaussian Processes with Nested Sampling
 - > Apply Gaussian Processes to categorical problems
- New descriptors:
 - Automatically designed descriptors to take into account of features not present in current set of descriptors
- Set split:
 - Look at other techniques





Acknowledgments

- Matthew Segall
- Edmund Champness
- Olga Obrezanova
- Christopher Leeding
- Andre Kramer



admensa-support@glpg.com



