

# Finding Multi-parameter Rules for Successful Optimization

ACS Spring National Meeting, April 10<sup>th</sup> 2013 Matthew Segall, Iskander Yusof, Edmund Champness

Patent pending

© 2013 Optibrium Ltd Optibrium™, StarDrop™, Auto-Modeller™ and Glowing Molecule™ are trademarks of Optibrium Ltd.

#### **Overview**

- Probabilistic scoring for multi-parameter optimization (MPO)
- Finding multi-parameter rules for drug discovery
- Methods
  - Rule induction
- Illustrative applications
  - 'Drug-like' properties
  - Oral CNS compounds
- User interaction
- Conclusions

# Probabilistic Scoring for MPO





# The Objectives of MPO

Identify chemistries with an optimal balance of properties

- Quickly identify situations when such a balance is not possible
  - -Fail fast, fail cheap
  - -Only when confident



# **Requirements for MPO in Drug Discovery**

- Interpretable
  - Easy to understand compound priority and how to improve compounds' chances of success
- Flexibility
  - Define criteria depending on therapeutic objectives of a project
- Weighting
  - Take into account relative importance of different endpoints to the success of a project
- Uncertainty
  - Take uncertainty into account, avoid missed opportunities

#### Probabilistic Scoring Scoring Profile



<sup>2013</sup> Optibrium Ltd. \* Segall *et al.* Chem. & Biodiv. **6** p. 2144 (2009)

#### StarDrop Prioritisation Probabilistic Scoring

- Property data
  - Experimental or predicted
- Criteria for success
  - Relative importance
- Uncertainties in data
  - Experimental or statistical

- Score (Likelihood of Success)
- Confidence in score



<sup>3 Optibrium Ltd.</sup> \* Segall *et al.* Chem. & Biodiv. **6** p. 2144 (2009)

# The Next Challenge

- How do we choose an appropriate scoring profile?
- Two approaches:
  - Domain/expert knowledge
  - Find the profile *automatically* using existing data
- Can we score compounds automatically without losing the benefits of expert knowledge?
  - Avoid 'black boxes'
  - Maintain interpretability and interactivity

#### Finding Multi-Parameter Rules for Drug Discovery





# **Objectives and Challenges**

- Use historical data to find scoring profiles with which to identify compounds with improved chance of success
  - Any drug discovery objective, e.g. clinical, PK, toxicity...
  - Once developed, profile can be applied prospectively to find new compounds
- Identify most important data with which to distinguish between successful and unsuccessful compounds
  - Any data can be used as input, calculated or experimental
- Explore multi-parametric data
  - Consider properties simultaneously, not individually
  - Avoid 'over counting' of correlated factors
- Rules must be interpretable and modifiable
  - Avoid black boxes
  - Synergy between computer and experts

#### What is a Rule?

 A Rule is a set of property criteria that in combination identify 'good' compounds, e.g.



• For example, Lipinski RoF:

logP<5	MW<500
HBD<5	HBA<10

# What is a Rule?

- A Rule is a box in multi-dimensional property space containing significantly more 'good' than 'bad' compounds
  - Equivalent to a scoring profile



# Methods





- The Patient Rule Induction Method (PRIM) by Friedman and Fisher is an effective way to find rules
- **Top-down peeling:** Start with a box covering all the compounds
- Then repeatedly peel the "worst" sides of the current box



- Bottom-up pasting: "Paste" back regions that we overzealously peeled
- We stop when pasting provides no improvement



- This peeling-and-pasting process gives us a **peeling sequence** of boxes
- We select a single box from the peeling sequence based on its performance over the validation set
- Resulting box corresponds to a rule for selection of successful compounds



• After finding one box, we remove the box's compounds from the dataset and start over



• The result is a series of boxes, each corresponding to an individual rule



#### **Measuring Rule Performance**

- Mean = Average objective value in box
  - Reported as % increase over objective value for full set
- Support = Proportion of data set 'covered' by box
  - Reported as % coverage
- Specificity vs. Sensitivity trade-off

# Variable Importance

- PRIM does not tell us the relevance of each property criterion to a given rule's predictions
- Cannot answer questions like:
  - Should we trade off solubility for potency?
  - Would it be valuable to generate data for a particular property?

# Variable Importance

•  $\alpha_i$  = false-negative rate of property criterion *i* 



• Importance =  $1 - \alpha_i$ 

#### **Illustrative Results**





#### Example: Drug-Like Properties QED\*

- Quantitative Estimate of Drug-Likeness (QED)\*
- Combine values for 8 properties

MW	logP	HBD	НВА
PSA	ROTB	AROM	ALERT

- For each individual property desirability function fitted to distribution for 771 oral drugs
- QED calculated as geometric mean of individual desirabilities



#### Example: Drug-Like Properties RDL\*

- Relative Drug Likelihood (RDL)\*
- Compare characteristics of 771 oral drugs with 1000 randomly selected compounds from ChEMBL database
  - What property values increase likelihood of compound being an oral drug?
- Used same 8 properties as QED
- RDL calculated as geometric mean of individual likelihoods



© 2013 Optibrium Ltd. \*Yusof and Segall, Drug Discov. Today DOI: 10.1016/j.drudis.2013.02.008

#### Example: Drug-Like Properties Rule Induction

- Rule induction applied to data set of 771 oral drugs and 1000 randomly selected compounds from ChEMBL
  - Random split 70:30 training:validation sets
- Used same 8 properties as QED and RDL as inputs
- Minimum coverage values compared
  - 20%, 30%, 40%, 50%

#### Example: Drug-Like Properties Rule Induction

• Minimum coverage 20% - 2 Rules

Profile		Desired Value	Importance
▲ Rule1			
MW	≤	444.855	
AROM	≤	1.01	
ALERTS	≤	1.01	

Set	Mean Improvement (%)	Coverage (%)
Train	60	22
Val	57	24

Profile		Desired Value	Importance
🔺 Rule 2			
ROTB	≤	4.04	
ALOGP	≤	2.727	

Set	Mean Improvement (%)	Coverage (%)
Train	51	23
Val	46	23

#### Example: Drug-Like Properties Rule Induction

• Minimum coverage 30%

Profile		Desired Value	Importance
AROM	≤	2.02	
MW	≤	334.775	

Set	Mean Improvement (%)	Coverage (%)
Train	51	36
Val	45	37

#### • Minimum coverage 40%

Profile		Desired Value	Importance
AROM	≤	2.02	
MW	≤	444.855	
ALERTS	≤	1.01	

Set	Mean Improvement (%)	Coverage (%)
Train	44	44
Val	42	46

#### • Minimum coverage 50%

Profile		Desired Value	Importance
AROM	≤	2.02	
MW	≤	432.745	

Set	Mean Improvement (%)	Coverage (%)
Train	35	57
Val	35	58

#### Example: Drug-Like Properties Results

 Applied to independent test set of 247 oral drugs and 1000 compounds randomly selected from ChEMBL



#### Example: Oral CNS Rule Induction

- Data set of 1191 drugs labelled as orally administered and CNS active or not
  - By approved route of administration and therapeutic indication (noisy)
- Divided into training (667), validation (286) and test (238) sets
- Calculated ADME properties from StarDrop<sup>™</sup> used as input:

logP	Solubility (logS)	Human Intestinal Absorption category (HIA)
Blood-brain barrier penetration (BBB log)	Plasma protein binding category	P-gp substrate category

• Minimum coverage 20%

#### Example: Oral CNS Results



# **User Interaction**





#### **Interactively Explore Profile Building**



# Conclusion

- MPO is a powerful approach to select and design compounds with a high chance of success
- Rule Induction helps to guide the development of scoring profiles to select compounds for a drug discovery objective
  - Apply to any objective
  - Use experimental or calculated data
  - Not black box synergy between computer and expert
- Identify most important data to guide selection of successful compounds
  - Optimise screening strategy and prioritise experimental resources
- For more information:
  - <u>matt.segall@optibrium.com</u>
  - <u>www.optibrium.com</u>



#### Acknowledgements

- Tatsu Hashimoto MIT
- Optibrium team, including:
  - Iskander Yusof
  - Ed Champness
  - Chris Leeding
  - James Chisholm
  - Hector Garcia Martinez