

# QSAR in Virtual Screening and Lead Optimisation

Dr. Peter Gedeck

Novartis Institutes for BioMedical Research



### Who are we?

Novartis at a glance

- World's third largest pharmaceutical company by sales
- One of 20 largest companies by market capitalisation
- Ranked among most respected companies worldwide
- Unique portfolio to meet changing healthcare needs:
  - Leading innovative pharmaceuticals
  - High-quality, low-cost generics (Sandoz)
  - Preventative vaccines
  - Eye care and Ophthalmic (CibaVision)
  - Consumer health products









# QSAR – definition



QSAR: Quantitatively relating biological activity to structure.

**NOVARTIS** 

 However, I will use the term QSAR very general and not only talk about activity, but any possible endpoint

4 | QSAR in Lead Optimisation and Lead Finding | Peter Gedeck | 12 June 2008



#### **Models**

Example: Modeling log*BB* (Blood-brain barrier)

Applications: QSAR in lead optimisation and lead finding



# QSAR – definition

#### Rational QSAR

• Use structural properties related to Free Energy Relationships

#### Pragmatic QSAR

- Use similarity, in effect predicting compounds according to similarity to those in the training set.
- Relevance of descriptors to Free Energy is unknown.



# **Rational QSAR**

- Use structural properties related to Free Energy Relationships
- Non-covalent Drug-Macromolecular interactions are equilibrium processes (as is partitioning)
  - Dissociation constants are determined by the difference in free energy of bound to unbound states.

 $\Delta G = -RT \log K_i$ 

 $K_i$  is proportional to  $IC_{50}$ , therefore  $\Delta G$  is proportional to  $-\log IC_{50}$ 

**NOVARTIS** 

 $\log IC_{50} = c \Delta G$ 

• Use –log  $IC_{50}$  or p $IC_{50}$  to model  $\Delta G$ 



# **Pragmatic QSAR**

- Neighbourhood principle:
  - Molecules with similar structural properties are likely to have similar biological activities
- Don't worry about theory, just work out combination of properties that give good activity
- Variety of descriptors derived for a structure.
  - Examples: Topological indices, pharmacophores, IR spectra, 2D fragments, structural fingerprints etc.
- However beware: Correlation without causation!



# **Correlation and Causation**

# A new parameter for sex education

SIR—There is concern in West Germany over the falling birth rate. The accompanying graph<sup>1,2</sup> might suggest a solution that every child knows makes sense.



Universität Düsseldorf, Moorenstrasse 5, D-4000 Düsseldorf, FRG

- Fachserie Gebiet und Bevölkerung (Statistisches Bundesamt, Kohlhammer, Stuttgart, 1984).
- Bauer, S. & Thielcke, G. Die Vogelwarte 31, 183-191 (1982).

Source: Sies H, Nature 332 (1988) 495. April issue!

9 | QSAR in Lead Optimisation and Lead Finding | Peter Gedeck | 12 June 2008

- The correlation of storks and babies on the left is 0.99
- Even if a descriptor is highly correlated with activity, it may not be its cause
- This is particularly true for pragmatic QSAR models





#### Models

#### Example: Modeling logBB (Blood-brain barrier)

Applications: QSAR in lead optimisation and lead finding

**NOVARTIS** 

10 | QSAR in Lead Optimisation and Lead Finding | Peter Gedeck | 12 June 2008

# Data quality

#### What does this graph show?



- Results for two different in-vitro assays for the same target (drift in assay conditions)
- Be realistic: QSAR models can only be as good as your data
   11 | QSAR in Lead Optimisation and Lead Finding | Peter Gedeck | 12 June 2008

#### Data quality

- Even within an assay we see variability of the data
  - Comparison of results from repeated measurements
  - Sample concentration, performance of cells



Be realistic: Garbage in – garbage out



# Data quality [Example logBB]

- Blood-brain distribution in mice
- Data distribution
  - Most results scattered around centre (logBB~0)
- Data quality
  - Variation of repeated measurements large for some compounds
- Remove compounds with large measurement error









#### Models

#### Example: Modeling logBB (Blood-brain barrier)

Applications: QSAR in lead optimisation and lead finding

**NOVARTIS** 

14 | QSAR in Lead Optimisation and Lead Finding | Peter Gedeck | 12 June 2008

# Descriptors

Classification based on representation of molecule

- 1D molecular formula
- 2D molecular connectivity / topology
- 3D molecular geometry / stereochemistry
- 4D/5D/... conformational ensembles











#### 16 | QSAR in Lead Optimisation and Lead Finding | Peter Gedeck | 12 June 2008

#### Descriptors

- Descriptors characterize various molecular properties by numerical values.
  - Hydrophobic, electronic, steric / size / shape, hydrogen bonding
- More then 2000 descriptors published
  - R. Todeschini, V. Consonni, Handbook of Molecular Descriptors, Wiley, 2000
- There are many approaches to select descriptors (forward selection, genetic algorithms, ...)





# Descriptors [logBB example]

- 1D/2D-Descriptors calculated using MOE [187 descriptors]
- Removed columns which don't change [160 descriptors]
- Range normalization of descriptors to interval [0,1]



Heatmap showing descriptor values. Rows and columns are ordered using hierarchical clustering



#### Models

#### Example: Modeling logBB (Blood-brain barrier)

Applications: QSAR in lead optimisation and lead finding

**NOVARTIS** 

18 | QSAR in Lead Optimisation and Lead Finding | Peter Gedeck | 12 June 2008

#### Models

- Model relates activity to the descriptors
  - Regression
  - Classification
- The ideal model
  - avoids overfitting (=perfect fit of training data but no generalization)
  - provides a robust estimate of the reliability and applicability of the model
  - is interpretable
  - works with all datasets
- An ideal model doesn't exist



# Regression

 Model a continuous dependent variable (activity) as a function of independent variables (descriptors)

$$y = f(X)$$

- Often used methods
  - Multiple Linear Regression
  - Partial Least Squares regression:
    - Like MLR but prevents overfitting
  - Artificial Neural Networks:
    - "Black Box" generates complex rules relating the activity to the descriptors.
  - Decision Trees, Random Forests:
    - Collection of logical rules that allow prediction of activity
  - Support vector machines, gaussian processes:
    - Novel machine learning approaches

# Regression [Decision trees and Random Forests]

- Decision tree: A sequence of rules links properties with activity
- Ensemble of trees / random forest: Combine prediction of multiple trees





# Regression [logBB example]

- Performance of random forest model on training data
  - r<sup>2</sup> = 0.92 [Excellent ... really?]
- Performance of model on test data
  - r<sup>2</sup><sub>pred</sub> = 0.53 [not so good anymore ... what went wrong?]



# Model validation

#### Internal cross-validation r<sup>2</sup><sub>CV</sub>/q<sup>2</sup> value

- Leave one out / leave many out
- Train several models and predict left out compounds
- Want q<sup>2</sup> to be >0.4

#### External Validation Training set/Test set r<sup>2</sup><sub>pred</sub> value

- Take a 20% subset of compounds as independent Test set.
- Other 80% the Training Set. Derive a model with Training set and predict the Test set.

#### Y-Scrambling

 Mix up the values for the biological activity (Y, dependent variable) attempt to build a model. If possible to model the scrambled data then we have a problem.

U NOVARTIS

23 | QSAR in Lead Optimisation and Lead Finding | Peter Gedeck | 12 June 2008

#### Model validation

- Regression statistics
  - Multivariate predictive r<sup>2</sup><sub>pred</sub>

$$r_{pred}^{2} = 1 - \frac{\sum_{i \in test} (y_{i}^{pred} - y_{i}^{act})^{2}}{\sum_{i \in test} (y_{i}^{act} - \overline{y}^{act})^{2}}$$

Correlation actual versus predicted r<sup>2</sup><sub>corr</sub>

$$r_{corr}^{2} = \frac{\left(\sum_{i \in test} (y_{i}^{pred} - \overline{y}^{pred})(y_{i}^{act} - \overline{y}^{act})\right)^{2}}{\sum_{i \in test} (y_{i}^{pred} - \overline{y}^{pred})^{2} \sum_{i \in test} (y_{i}^{act} - \overline{y}^{act})^{2}}$$



# Model validation [logBB example]

Training performance: r2=0.92 rmse=0.203 correlation=0.97



Internal validation: q2=0.52 rmse=0.491 correlation=0.72



External validation: pred-r2=0.53 rmse=0.477 correlation=0.73



#### Repeated training/test splits: Boxplot shows variation of statistics



#### Y-scrambling:

r<sup>2</sup> still very good, but q<sup>2</sup> values below 0



# Classification

 Model a discrete dependent variable (class) as a function of independent variables (descriptors)

$$c = f(X)$$

- Many methods exist, e.g.:
  - Linear Discriminant analysis
  - Decision Trees
  - Random Forests
  - Convert regression model into classifier
- Validation
  - Internal cross-validation
  - External validation using test sets
  - Scrambling





### Classification [logBB example]

 Prediction from Random Forest Regression converted into Classification

$$c = \begin{cases} blood & \text{if } f(X) < 0\\ brain & \text{if } f(X) > 0 \end{cases}$$





# Classification [logBB example]

- Confusion matrix
  - Table
  - Graph
- Usual statistics
  - Accuracy: 73% (percent correctly predicted)
  - Recall (percent correctly predicted by predicted class):
     65% brain, 80% blood
  - Precision (percent correctly predicted by actual class) 76% brain, 71% blood

#### Actual Predicted Blood Brain Blood 0.41 0.1 Brain 0.17 0.32





#### Models

Example: Modeling log*BB* (Blood-brain barrier)

Applications: QSAR in lead optimisation and lead finding



# Lead optimisation: a typical flow chart



# Lead optimisation: a typical flow chart



#### HQSAR Descriptors

Break structure into fragments and count occurrences

Combine all counts for all possible fragments into a vector of numbers = hologram

Reduce length of vector by folding = reduced hologram









#### PLS:

Multiple linear regression model suitable for cases with fewer observations than descriptors.

HQSAR-model: Prediction = f(reduced hologram)



- HQSAR model (left: medium/low quality)
  - Potent compounds are underestimated
- Note assay variability as shown on the right





- Method is fast and can easily be applied to larger virtual libraries
- Using Novartis in-house data (about 500 assays), we showed that using HQSAR
  - for 35% of assays we get  $r_{pred}^2$ >0.5 and for
  - for 55% of assays we get  $r_{pred}^2$ >0.4

#### Lead optimisation – the more rational approach

- Protein-ligand interactions do not care about the specific molecular structure
- The receptor only recognises the electron cloud, rather than the bonding structure, allowing us to propose new chemotypes
- But how can QSAR capture this?





- Comparative Molecular Field Analysis
- What is a molecular field?







- Align all molecules
- Calculate molecular field at discrete points in space around each molecule individually





- Large number of highly correlated descriptors
- Number of descriptors much larger than number of compounds
- Determine regression with activity using Partial Least Squares (PLS) method



NOVARTIS





Steric Probe: Green steric bulk favoured Yellow steric bulk disfavoured



Electrostatic Probe: Red electronegative field favoured Blue electropositive field favoured



40 | QSAR in Lead Optimisation and Lead Finding | Peter Gedeck | 12 June 2008

- CoMFA model identifies points in space that most significantly determine activity
  - It will only learn what controls levels of activity
  - If you don't test a possible interaction, it will not highlight it as positive or negative (e.g. hinge binder in kinases)

() NOVARTIS

- Analysis of CoMFA models allows to suggest regions for modifications
  - For example where changing the steric bulk, would increase or decrease binding
- Structural alignment of compounds is crucial
  - It can impact the interpretation of model





#### Models

Example: Modeling log*BB* (Blood-brain barrier)

Applications: QSAR in lead optimisation and lead finding



# QSAR in Lead finding

#### Lead optimisation

- Few data points (100s)
- High quality assay
- Biased towards actives

#### Lead finding (HTS)

- Huge number of data points (1.000.000s)
- Low quality assay
- Majority of compounds is inactive (0.1-1% hit rate)



# The problem in lead finding

20 % false positives



90% false positives

**U** NOVARTIS

44 | QSAR in Lead Optimisation and Lead Finding | Peter Gedeck | 12 June 2008

# The problem in lead finding



45 | QSAR in Lead Optimisation and Lead Finding | Peter Gedeck | 12 June 2008

# The problem in lead finding

- Can we reduce the number of false positives and rescue false negatives?
- A very pragmatic approach:
  - Build a classification model to predict compounds to be active or inactive
  - or inactive
    Glick M, et al. Enrichment of extremely noisy HTS data using a naïve Bayes classifier. J Biomol Screen 2004; 9; 32 DOI: 10.1177/1087057103260590



#### Extended connectivity descriptors

- Extended connectivity descriptors (ECFP)
  - aka circular substructure fingerprints
  - Implementation in PipelinePilot (SciTegic / Accelrys)
- Fragment based descriptor scheme
  - Atoms and their environment



# Naïve Bayesian Classifier

Simple classifier based on Bayes' theorem

$$p(C|F_1, \dots, F_n) = \frac{1}{Z}p(C)\prod_{i=1}^n p(F_i|C)$$

- Probability that compound with features F<sub>i</sub> belongs to class C is determined by the product of probabilities to find Feature F<sub>i</sub> when it is in class C
- The features used in this study are the occurrence of the extended connectivity fragments



#### Results

- Classification models successfully improved hit rate
- Amazing robustness of models against noise levels



Similar publications on the use of Naïve Bayes classifier to improve virtual screening / high-throughput docking results

**NOVARTIS** 



#### Models

Example: Modeling log*BBB* (Blood-brain barrier) Applications: QSAR in lead optimisation and lead finding





# Summary

- Data quality
  - It's important to understand your data and their quality
- Descriptors
  - Correlation means not necessarily causation
- Methods
  - No methods works for all cases
  - Validation is an important aspect of model building
- Applications
  - Activity and property prediction (best when rational QSAR)
  - Assistance in compound and library design
  - Understanding of changes required to gain good properties
  - Lead finding
- Importantly we must not be surprised to fail

### Acknowledgments

- Richard Lewis
- Peter Hunt
- Katrin Spiegel
- Lilya Sviridenko
- Michael Carter
- Meir Glick
- John Davies

