# QSAR: an introduction

Wim Dehaen, 30.1.2024 7ADD, Olomouc

# Short introduction: Who Am I

- Researcher in D. Svozil's group (UCT Prague) Applied cheminformatics
- Researcher in P. Perlikova's group (UCT Prague) Computational support for medicinal chemistry, mainly SBDD
- My interests:
  - Cheminformatics
  - Medicinal chemistry and CADD
  - Organic chemistry
  - Digital signal processing (of audio)

0...



**Department of Informatics and Chemistry** 



**Department of Organic Chemistry** UCT PRAGUE

## Overview

- What is QSAR?
- The history of QSAR
- Structure activity relationships
- Data and featurization in QSAR
- QSAR methods
- QSAR validation
- Applicability domain and interpretability
- QSAR applications

# What is QSAR

#### • QSAR is:

 $\circ$  Quantitative Structure Activity Relationships

- QSOR (odorants, olfactory chemicals)
- QSPR (nonbiological properties, e.g. max. absorption and emission wavelength)
- Medicinal chemistry + Physical organic chemistry + Statistics (my opinion)
- "An application of data analysis methods and statistics to developing models that could accurately predict biological activies or properties of compounds based on their structures" - A. Tropsha

Predicted biological activity = Function (structural features) + error

# What is QSAR

#### • Overview:

- $\circ$  Data
  - Data source (molecular structure + end points)
  - Regularization of data
  - Featurization
  - Regularization of features
- $\circ$  Model-building
  - Appropriate choice of model
  - Appropriate data splits
- $\circ$  Evaluation
  - Metrics
  - Interpretation
  - Applicability Domain
- $\,\circ\,$  Application
  - Virtual screening
  - Hit optimization
  - ADME(T) filters



# What is QSAR

- Machine learning:
  - $\odot$  QSAR is a form of "traditional" supervised machine learning
    - Non-deep:
      - Small data regime
      - Interpretability
  - $\circ$  "Deep QSAR"
    - Seems mostly a semantic distinction between deep QSAR and deep learning for molecular property prediction
- Virtual screening:
  - $\odot$  QSAR models can be applied for virtual screening:
    - Property filters (solubility, aggregation, toxicity)
    - Activity
  - $\odot$  But Applicability Domain issues!

- Legendre, Gauss (1805):
  - $\odot$  Method of least squares used to determine orbits of spatial bodies based on astronomical observations
- Crum-Brown and Fraser (1865):

 $\odot$  Relation between structure and physiological action

• Hammett equation (1935):

 Empirically derived substituent constants derived from substituted benzoic acid hydrolysis rate



### 159 years ago!

Although we cannot obtain a rational explanation of the connection between the chemical and physiological characters of a substance until we know more of the *modus operandi* of poisons, it might be supposed that a careful examination and comparison of known facts would lead to the discovery of some empirical law or laws by means of which we could deduce the action from the chemical constitution. Unfortunately, however,

Hantsch equation

$$\log \frac{1}{C} = k_1 D_1 + k_2 D_2 + \dots + k_i D_i$$

with k<sub>i</sub> and D<sub>i</sub> the ith constant and descriptor respectively

• Does not have to be linear, for example Hantsch famous plant hormone work had this form:

$$\log \frac{1}{C} = -k_1 (\log P)^2 + k_2 (\log P) + k_3 \sigma + k_4$$

- Free-Wilson formalism
- Biological activity (log scale) is determined by the sum of the activity of the reference compound and substituent contributions
- With µ the biological activity of the unsubstituted analog and sum(a<sub>i</sub>) the sum of substituent contributions.

$$\log \frac{1}{C} = \sum a_i + \mu$$

#### • CoMFA

Comparative Molecular Field Analysis

 $\odot$  "The first 3D QSAR"

• Shape, electrostatic, H-bonding, ... 3D features

- placing aligned conformations in grid and probing with e.g. lipophilic probe
- Resulting model is somewhat like a pharmacophore in that it captures spatial distributions of features

#### • OECD and QSAR use in regulatory context

 $\odot$  Hazard assessment of chemicals

OECD QSAR toolbox

 $\circ$  e.g. Genotoxicity

 $\circ$  e.g. Biodegradation

 $\circ$  e.g. Skin sensitization

 $\circ$  Intended use: filling gaps in (eco)toxicity data

# **OECD QSAR principles**

- A defined endpoint
- An unambiguous algorithm
- A defined domain of applicability
- Appropriate measures of goodness-of-fit, robustness, and predictivity
- A mechanistic interpretation, if possible

## Structure activity relationship

- Very important concept in medicinal chemistry
- Intuitive and ad hoc (non quantitative!)
- "An understanding of the SAR for a set of molecules allows one to rationally explore chemical space, which in the absence of "sign posts" is essentially infinite" R.Guha, In Silico Models for Drug Discovery, Methods in Molecular Biology, vol. 993

# Structure activity relationship

- Who recognizes this man?
- Voted "3rd Greatest Belgian"
  - 1. Father Damian (Catholic missionary)
  - 2. Eddy Merckx (Cyclist)
  - 3. ???



## Paul Janssen of "Janssen" fame

- Brought 80 drugs to market

   Fentanyl
   Haloperidol
- Emblem of the trial-and-error analog exploration SAR approach taken by medicinal chemists

Ο

Fentanyl is the result of SAR exploration and optimization of meperidine /

## Data source

- Publicly available:
  - $\circ$  ChEMBL
  - $\circ$  PubChem
  - $\odot$  Community benchmarks such as SAMPL
  - $\odot\,\text{Ad}$  hoc data sets in the literature
- Commercial databases:
  - $\circ$  Reaxys
  - $\circ\, \text{Scifinder}$ 
    - "Please do not use our product in this way" Scifinder Rep when i asked them if I can use the result of Scifinder Queries in open source cheminformatics workflows

# Data cleanup

- Endpoints (e.g. Assay data)
  - $\circ$  Experimental noise
  - $\odot$  Annotation inconsistencies and errors
  - Fundamental incompability between different measurements (e.g. different assay conditions)
  - $\odot$  Some publications more trustable than others
    - (e.g. HTS hits are more noisy and have more chance to have assay artefacts)

#### Combining IC50 or Ki Values From Different Sources is a Source of Significant Noise

10 January 2024, Version 1

Working Paper

Gregory A. Landrum 💿, Sereina Riniker 💿

Show author details ~

# Data cleanup

- Molecules:
  - Molecular identity and normalization (tautomers, protonatino, kekulization, charges, salts, ...)
  - $\odot$  Solvent, pH, T dependent
  - $\odot$  Parsability in software (some SMILES will load in obabel but not RDKit abd vice versa)
  - Undesired properties (e.g. atoms that forcefield can't deal with)
  - $_{\odot}$  In some case, chemical entities can also be mixtures, reactions, specific conformations, ...

### Some examples



Source: So you think you understand tautomerism?



Formally tautomers, but won't interconvert:

$$H_3C \longrightarrow CH_3 \implies H_2C \longrightarrow CH_2$$

## Featurization of molecules

• Descriptors:

 $\odot$  Graph invariants for a chemical graph

- Take into account chemical graph is colored, weighted, non-directed and has extra data for each vertex like charge, hybridization status, ...
- $\odot$  Should be invariant to:
  - Vertex order of the graph
  - But also rotation and translation (if vertices have 3d coordinates included)
- "There are three keys to the success of any QSAR model building exercise: descriptors, descriptors, and descriptors" - Alexander Tropsha

# Descriptors

- Many, many possibilities
- A fingerprint is essentially a list of descriptors resulting from a common process
  - $_{\odot}$  Usually binary vector, sometimes integer, occasionally float
- Binary descriptors:
  - $\circ$  Presence (1) or non-presence (0) of one or several structural motifs
  - $\,\circ\,$  Topological pharmacophore
  - Bag of fragments (e.g. a bit in Morgan fingerprint)
  - $\circ$  One specific fragment (a bit in a structural key such as Klekota-Roth FP)
- Integer descriptors:
  - $\circ\,$  Count of one or several structural motifs:
  - $\odot$  Same possibilities as above, but count occurence
  - $_{\odot}$  Some Lipinski-like descriptors: HBD, HBA, HAC

# Descriptors

- Real descriptors:
  - $\circ$  Continuous range
  - $\odot$  Topological indices
    - Zagreb, Kier hall, Randic, Kappa, Wiener, etc
  - Physicochemical descriptors (MW, TPSA, ...)
  - o cLogP (this is in a sense calculated by a QSPR model itself)
  - $\odot\,\text{QM}$  calculated properties
  - Empirical (e.g. boiling point, NMR shift, partition coefficient, ...)

# Descriptors

#### • 1D:

 $\circ$  Molecular weight, formula, ...

 "Whenever the underlying biological observables may depend on transport as well as receptor fit, such as passive membrane penetration, "1D" descriptors, such as log P, polar surface area, and pKa, should be considered. " - Richard D. Cramer

#### • 2D:

 Based on "2D structure" of molecule (actually a dimensionless graph), topological indices, substructure presence, ...

#### • 3D:

 Conformation dependent, surface area, volume, shape, QM-calculated descriptors, ...

#### • 4D:

 $\odot$  Dynamics, flexibility of bonds, MD-PLIF, ...

# Fingerprint example: ECFP

- Extended connectivity fingerprint aka Morgan(like) fingerprint
- Molecular fingerprint based on topological neighborhoods around atoms at a given topological radius threshold
- Commonly used as input in QSAR, highly performant and generic
- Implemented in major cheminformatics packages









## Feature normalization

- Descriptors are often correlated
- Mixing binary, integer and real data
- Some features have no relation with the end point at all
- The fewer features, the faster and more interpretable the model is

### Feature normalization

• Scaling descriptors:

e.g. using z scores
 Scale between [0,1] by linearly scaling [min(D),max(D)]

### • Collinearity:

 $\circ R^2$  between descriptor should be under a given threshold (e.g. 0.8)

### Feature selection

- Pruning after model building based on feature importances o If method permits this, e.g. RF
- Building models with one or more descriptors excluded
- Principal component analysis

# Training a model

- In general, QSAR tends to use "standard machine learning" algorithms and simpler methods like logistic regression

   As opposed to more complex deep learning based approaches
- Supervised learning (has an endpoint guiding it)
  - $\circ$  Classification:
    - The endpoint is categorical: two or more labels. E.g. Active/Inactive
  - Regression:
    - The endpoint is continuous. E.g. pIC50 between 3.0-9.0
- Unsupervised learning:
  - Clustering:
    - The data set gets divided into two or more clusters
    - Is this really QSAR?

# Training a model

• Linear approaches

PArtial Least Squares, Multiple Linear Regression, Free-Wilson

- Non-linear approaches
  - $\odot$  "All the rest"
  - $\circ \, \text{See next slide}$

# Typically used approaches

### • Regression:

- $\circ$  Decision tree
- $\circ$  Support Vector Machine
- o Neural Network
- $\circ$  Random Forest
- $\circ$  k Nearest Neighbors
- $\odot$  Multiple linear regression
- $\odot$  Partial least squares
- $\circ$  Gaussian process

# Typically used approaches

### • Classification:

- $_{\odot}$  Decision tree
- $\circ$  Support Vector Machine
- o Neural Network
- $\circ$  Random Forest
- $\circ$  k Nearest Neighbors
- $\circ$  Logistic regression
- $\circ\,\text{Naive}$  Bayes

# Typically used approaches

- Clustering
- Not really QSAR



### **Example: Random forest**

• Decision tree:



# **Example: Random forest**

Random Forest



ntro\_rt.html

system-

# Data splitting

• Why? Over- and underfitting

Classification



Source: https://www.mathworks.com/discovery/overfitting.html

# Data Splitting



# Data splitting

Random split

 $_{\odot}$  This can make tasks too easy: similar scaffolds in test and train

Scaffold split

Shows if model is able to hop between scaffolds (to some extent)

• Time split

 $_{\odot}$  To mimic actual discovery process

## Metrics

• Binary classification metrics:

Often calculated from confusion matrix (see next slide)

 $\odot$  Source of slide: wikipedia

		Predicted condition		Sources: [11][12][13][14][15][16][17][18][19] view · talk · edit	
	Total population = P + N	Predicted Positive (PP)	Predicted Negative (PN)	Informedness, bookmaker informedness (BM) = TPR + TNR - 1	Prevalence threshold (PT) = $\frac{\sqrt{TPR \times FPR} - FPR}{TPR - FPR}$
condition	Positive (P)	True positive (TP), hit	False negative (FN), type II error, miss, underestimation	True positive rate (TPR), recall, sensitivity (SEN), probability of detection, hit rate, power $= \frac{TP}{P} = 1 - FNR$	False negative rate (FNR), miss rate $= \frac{FN}{P} = 1 - TPR$
Actual	Negative (N)	False positive (FP), type I error, false alarm, overestimation	True negative (TN), correct rejection	False positive rate (FPR), probability of false alarm, fall-out $= \frac{FP}{N} = 1 - TNR$	True negative rate (TNR), specificity (SPC), selectivity $=\frac{TN}{N} = 1 - FPR$
	$\frac{Prevalence}{P + N}$	Positive predictive value (PPV), precision = $\frac{TP}{PP} = 1 - FDR$	False omission rate (FOR) = $\frac{FN}{PN} = 1 - NPV$	Positive likelihood ratio (LR+) = $\frac{TPR}{FPR}$	Negative likelihood ratio (LR–) = $\frac{FNR}{TNR}$
	$\begin{array}{l} \text{Accuracy} \\ \text{(ACC)} \\ = \frac{\text{TP} + \text{TN}}{\text{P} + \text{N}} \end{array}$	False discovery rate (FDR) = $\frac{FP}{PP} = 1 - PPV$	Negative predictive value (NPV) = $\frac{TN}{PN}$ = 1 - FOR	Markedness (MK), deltaP (Δp) = PPV + NPV – 1	Diagnostic odds ratio (DOR) = $\frac{LR+}{LR-}$
	Balanced accuracy (BA) = $\frac{TPR + TNR}{2}$	$F_{1} \text{ score}$ $= \frac{2 \text{ PPV} \times \text{TPR}}{\text{PPV} + \text{TPR}} = \frac{2 \text{ TP}}{2 \text{ TP} + \text{FP} + \text{FN}}$	Fowlkes-Mallows index (FM) = √PPV×TPR	Matthews correlation coefficient (MCC) =√TPR×TNR×PPV×NPV -√FNR×FPR×FOR×FDR	Threat score (TS), critical success index (CSI), Jaccard index $= \frac{TP}{TP + FN + FP}$

## **Metrics**

. . .

- Regression metrics:
- Q<sup>2</sup> Squared leave-one-out cross-validation correlation coefficient
- R<sup>2</sup> Coefficient of determination
- MAE Mean absolute error

$$Q_{abs}^2 = 1 - \sum_{\gamma} \left(Y_{exp} - Y_{LOO}\right)^2 / \sum_{\gamma} \left(Y_{exp} - \langle Y \rangle_{exp}
ight)^2$$

$$R_{abs}^2 = 1 - \sum_{\gamma} \left( Y_{exp} - Y_{pred} \right)^2 / \sum_{\gamma} \left( Y_{exp} - \langle Y \rangle_{exp} \right)^2$$

$$\textit{MAE} = \sum_{\textit{Y}} \left|\textit{Y} - \textit{Y}_{\mathsf{pred}} \right| / n$$

Source: Best Practices for QSAR Model Development, Validation, and Exploitation

# Interpretability

- Simple models like Hantsch are sometimes naturally interpretable
- Some models are "black boxes" but can be interpreted
- Some models offer feature importances directly
- Some tricks exist for probing explanations (model agnostic):

   Shapley scores (See the Tutorial) "For the given prediction, how much has each feature contributed compared to the average prediction?"
   Counterfactuals: "what changes will result in an alternate outcome?"

# **Applicability Domain**

Several ways of assessing it:

 Leverage (thresholds on williams plot)
 Conformal prediction



- In general, a standard QSAR model has a SMALL domain of applicability
- Remember: "In general classical statistics is far too optimistic when validating a QSAR, because its underlying assumptions about data distributions are contradicted by the extraordinarily heterogeneous nature of chemical structures and mechanisms of biological response. Restricting the structural scope of a QSAR should help, but the distribution of "local" structural variations, within a series undergoing lead optimization, is also unlikely to be uniform." Richard A Kramer

# **Applicability Domain**



X

Representations on Landscape Topology and the Formation of Activity Cliffs.

# **QSAR** applications

• Virtual screening



Source: J. Chem. Inf. Model. 2014, 54, 2, 634-647

# **QSAR** applications

- Multiparameter optimization: e.g. ADME(T)+antitarget+potency where one or several of the parameters are calculated by discrete QSAR models
- E.g. a reward function in RL can often include "SA" or "QED"

# **QSAR** applications

- ADME(T) filters
- Aggregation detection: SCAM detective
- Tox21 challenge



# **QSAR** tutorial

- Open source software:
  - RDKit (cheminformatics), scikit-learn (machine learning), python data science stack (scipy, numpy, pandas), shap (model explainability), Jupyter notebooks (interactive python environment)
- Some python knowledge is useful, but notebooks are constructed so they can be re-used and adjusted without too much tinkering

### Thanks for the attention

- Thank you to my funders; CZ-OPENSCREEN and Czech Science Foundation
- And see you @ the tutorial!