



THE UNIVERSITY OF TOKYO



Chemical System Engineering

Chemical Structure Generation Based on Inverse Quantitative Structure-Property Relationship/Quantitative Structure-Activity Relationship

Kimito Funatsu

Chemical System Engineering, The University of Tokyo

Strasbourg Summer School on Chemoinformatics

June 27, 2018

- 1: General introduction
 - Inverse QSPR/QSAR
 - Objective and hypothesis
- 2: Structure generation
- 3: Inverse QSPR/QSAR analysis (from y to \mathbf{x})
- 4: Structure generation based on inverse QSPR/QSAR
- 5: Summary

Molecular design with inverse quantitative structure-property/activity relationship (QSPR/QSAR)

QSPR/QSAR

Data from experiments
(compound, property)



$$y=f(x)$$

Descriptors

Property,
activity..

Chemical
structures

Quantitative structure-property relationship (QSPR)
Quantitative structure-activity relationship (QSAR)

x

Descriptors	Values
MW [g/mol]	180.04
#HBA	3
#NBD	1
#Aromatic Rings	1
TPSA [\AA^2]	63.6

$$y=f(x)$$

y

Property	value
MP [degree]	-5
Viscosity [Pa·s]	8.0
LogP	3
LogS	-1

Inverse QSPR/QSAR

Data from experiments
(compound, property)

$$y=f(x)$$

Descriptors

Property,
activity..

Chemical
structures

Chemical
structures

Structure
generator

X

Descriptors	Values
MW [g/mol]	180.04
#HBA	3
#NBD	1
#Aromatic Rings	1
TPSA [Å ²]	63.6

x: Explanatory
variables
(Descriptors)



Inverse
QSPR/QSAR
analysis



y: Objective variable
(property, activity)



Chemical structures
(chemical graphs)

Obtaining **x** information from **y**

x: Explanatory
variables
(Descriptors)

Inverse
QSPR/QSAR
analysis

y: Objective variable
(property, activity)



Chemical structures
(chemical graphs)

Generating structures based on **x** information

x: Explanatory variables
(Descriptors)



Inverse
QSPR/QSAR
analysis



y: Objective variable
(property, activity)



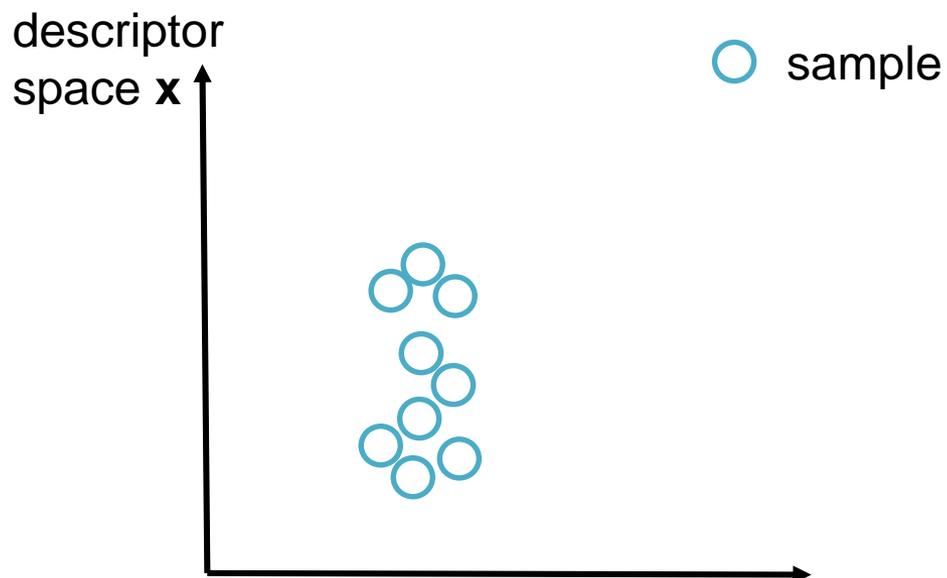
Chemical structures
(chemical graphs)

Obtaining **x** information from **y**

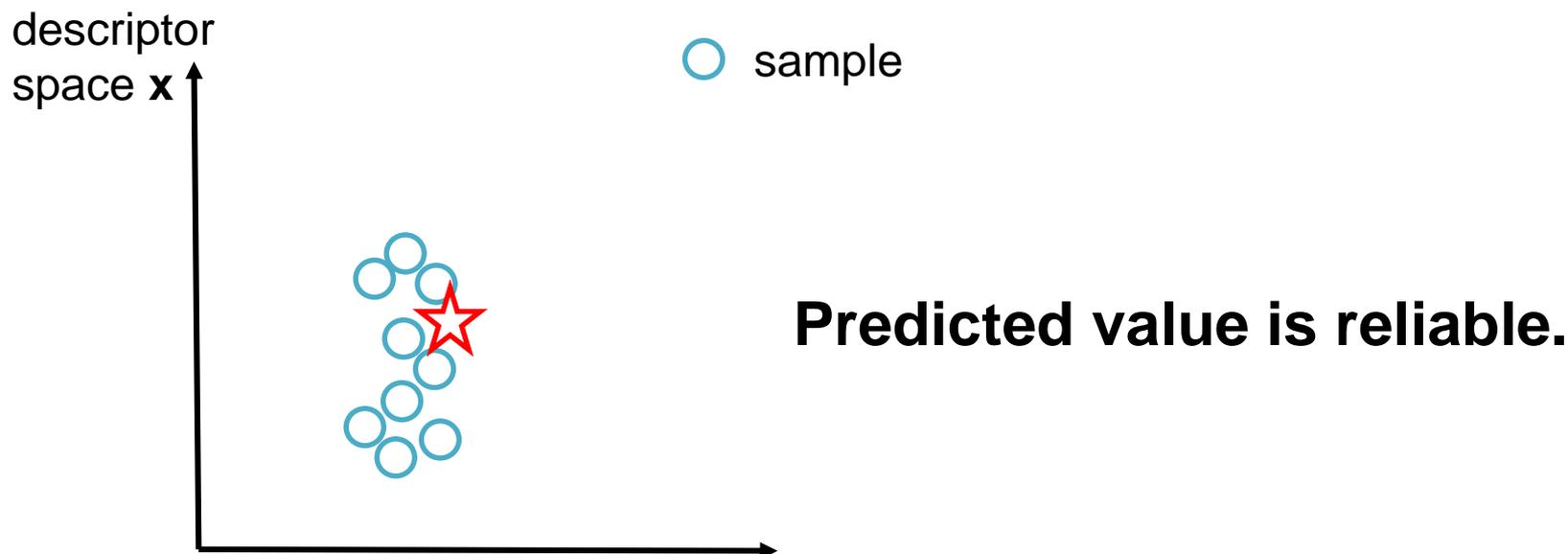
Not considering applicability domain
(AD)

Poor predictability by multiple linear
regression(MLR) model

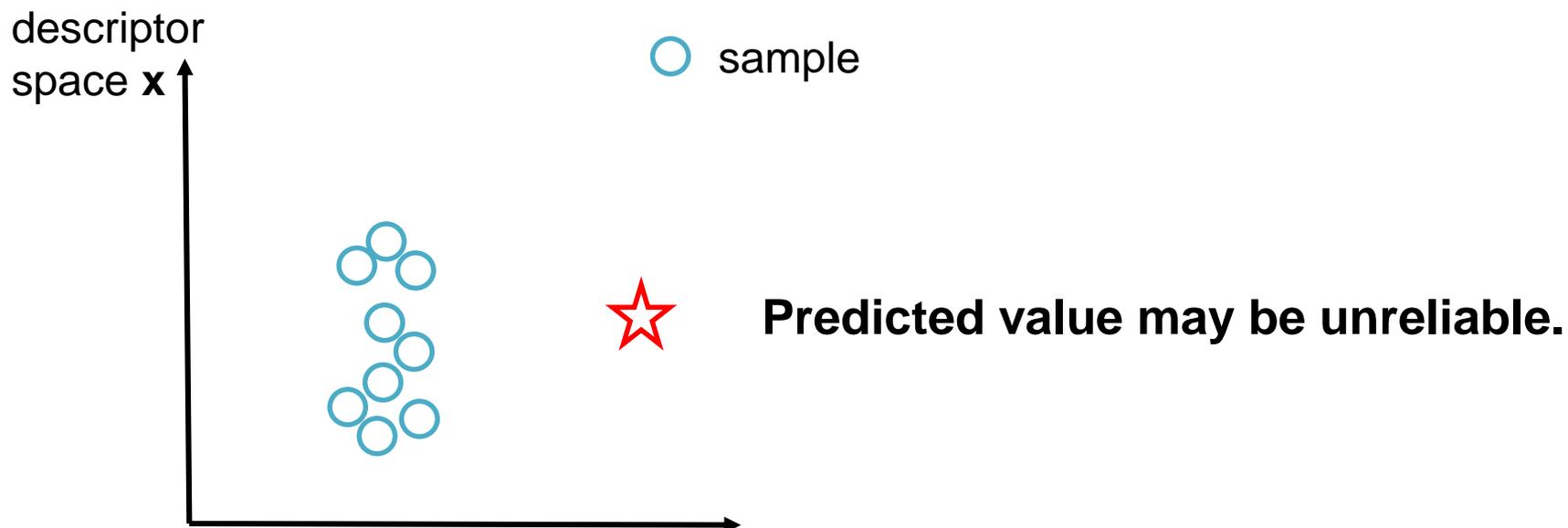
- Only inside AD, predicted values produced by regression models should be reliable.
 - Density-based method
 - Ensemble-based method



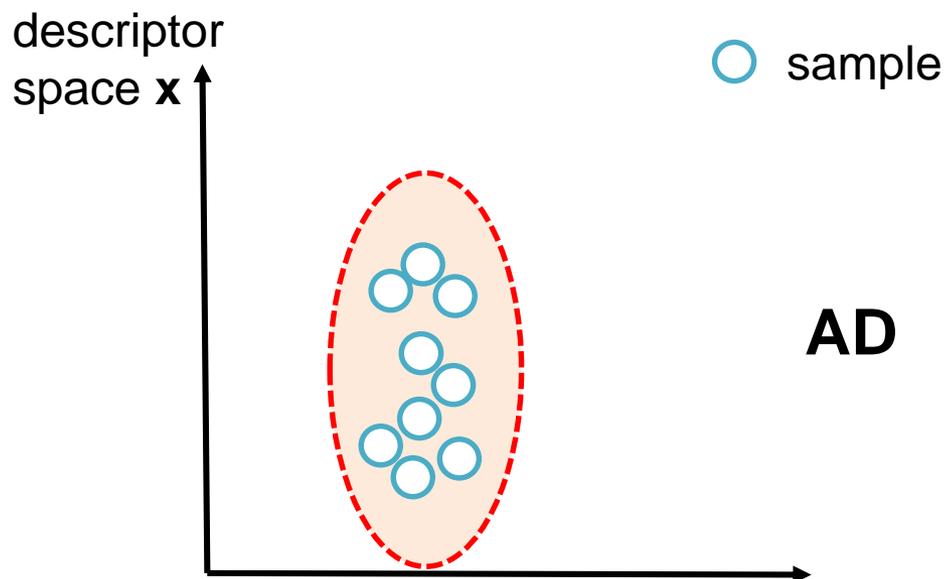
- Only inside AD, predicted values produced by regression models should be reliable.
 - Density-based method
 - Ensemble-based method



- Only inside AD, predicted values produced by regression models should be reliable.
 - Density-based method
 - Ensemble-based method



- Only inside AD, predicted values produced by regression models should be reliable.
 - Density-based method
 - Ensemble-based method

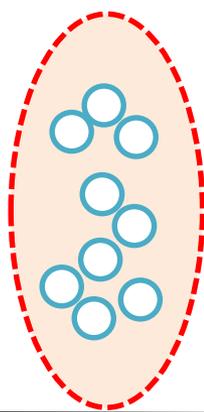


AD In a density-based method, Probability density function $p(\mathbf{x})$ (PDF) is a criterion for AD

- In inverse QSPR/QSAR analysis, AD has not been considered.
 - Not considering training data information
 - Extrapolation is allowed without limitation.

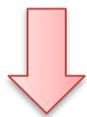
descriptor
space \mathbf{x}

○ sample



AD In a density-based method, Probability density function $p(\mathbf{x})$ (PDF) is a criterion for AD

x: Explanatory
variables
(Descriptors)



Chemical structures
(chemical graphs)

Inverse
QSPR/QSAR
analysis

y: Objective variable
(property, activity)

Chemical structure generation

Treating limited variety (number) of
descriptors

Not considering universal AD

- Specific type of descriptors is employed.
 - Kier indices, X indices, Wiener index.
 - Signatures
- Proper descriptor set varies from projects to projects.

M. I. Skvortsova, I. I. Baskin, *et al.*, *J. Chem. Inf. Comput. Sci.*, 33, 4, 630–634, 1993.

C. J. Churchwell, *et al.*, *J. Mol. Graph. Model.*, 22, 263–273, 2004.

Kirkpatrick, P. and Ellis, C. *Nature*, 432, 823–823, 2004.

- Universal AD is an abstract concept, which is irrelevant with models
 - Determined based only on the training data before constructing any QSPR/QSAR models.

- Simple example: boiling point model

$$\text{bp}(\text{°C}) = -126.19 + 33.42N_c - 6.286T_m$$

N_c : Number of **carbon** atoms

T_m : Number of **terminal carbon** atoms

$$n = 39$$

$$s = 5.86$$

$$r^2 = 0.987$$

- Simple example: boiling point model

$$\text{bp}(\text{°C}) = -126.19 + 33.42N_c - 6.286T_m$$

This equation is valid for **C2-C8** alkanes.



Structures to be generated
should be restricted to alkanes.

Challenges in inverse QSPR/QSAR analysis

Chemical structure generation

Treating limited variety (number) of descriptors

Not considering universal AD

Obtaining x information from y

Not considering AD

Poor predictability by MLR

To develop a practical chemical structure generation system based on inverse QSPR/QSAR by overcoming the challenges.

x: Explanatory variables
(Descriptors)



Inverse
QSPR/QSAR
analysis



y: Objective variable
(property, activity)



Chemical structures
(chemical graphs)

Workflow of QSPR/QSAR analysis

Experimental data
(compound, property/activity)

```
graph TD; A[Experimental data (compound, property/activity)] --> B[x: Molecular descriptors]; B --> C[QSPR/QSAR model]; C --> D[y: predicted property/activity]; E[a specific y value];
```

x: Molecular descriptors

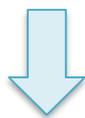
QSPR/QSAR model

y: predicted property/activity

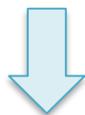
a specific y value

Workflow of QSPR/QSAR analysis

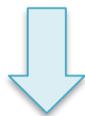
Experimental data
(compound, property/activity)



x: Molecular descriptors



QSPR/QSAR model



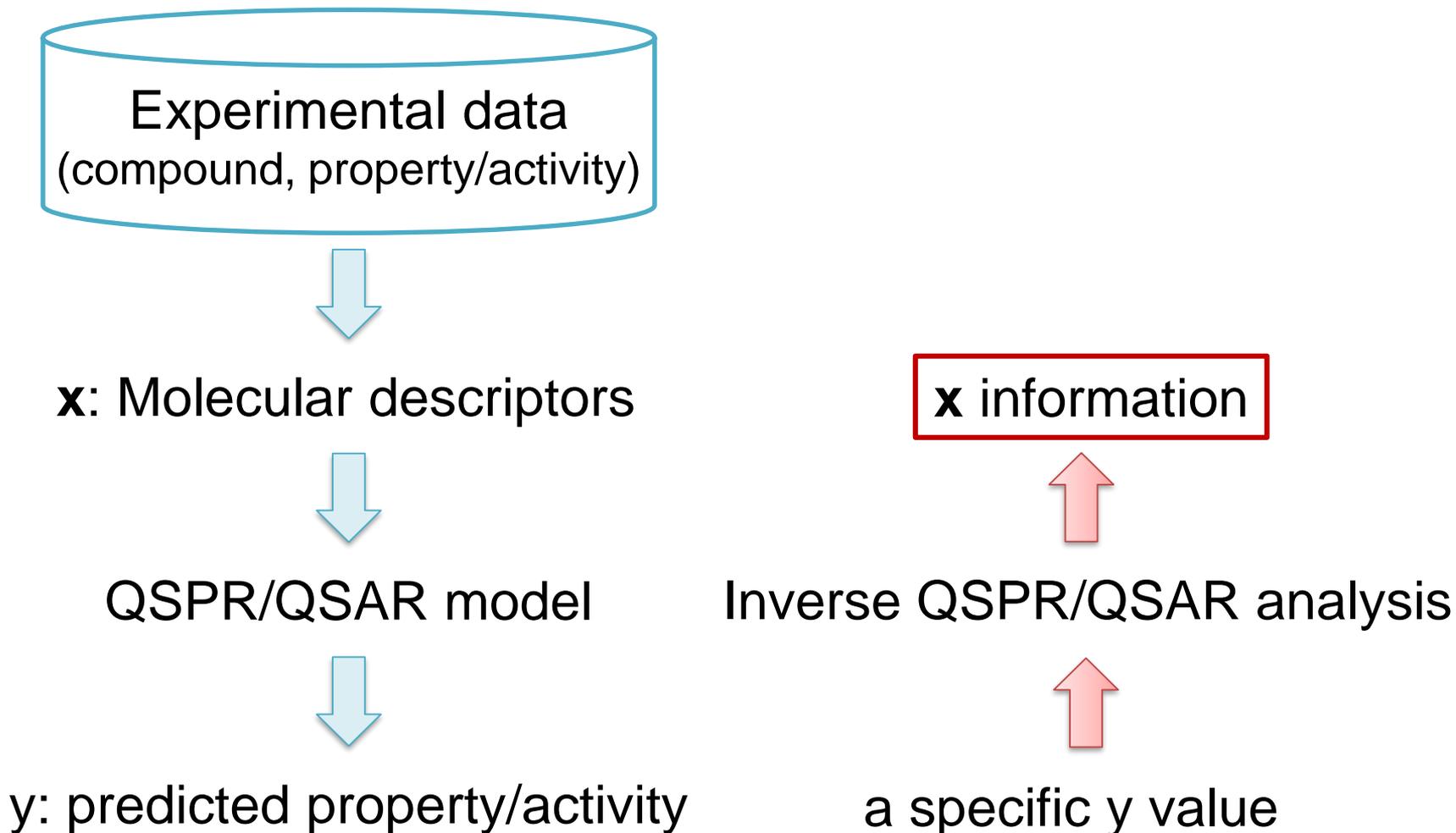
y: predicted property/activity

Inverse QSPR/QSAR analysis



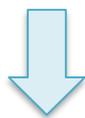
a specific **y** value

Workflow of QSPR/QSAR analysis

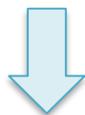


Workflow of QSPR/QSAR analysis

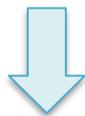
Experimental data
(compound, property/activity)



x: Molecular descriptors



QSPR/QSAR model



y: predicted property/activity

Exhaustive chemical
structures based on
x information



x information



Inverse QSPR/QSAR analysis



a specific **y** value

Chemical structures giving a specific y can be exhaustively generated in inverse QSPR/QSAR by

- ✓ considering local and universal ADs,
- ✓ using efficient structure generator,
- ✓ introducing variety of descriptors.

To develop a practical chemical structure generation system based on inverse QSPR/QSAR

Chemical structure generation

Explaining a methodology for using variety of descriptors

Describing algorithms for treating chemical graphs

Obtaining x information from y

Introducing probability density for treating AD

Explaining a non-linear regression methodology

Goal To develop a structure generator that overcomes challenges as follows:

Challenges

Chemical structure generation

Treating limited variety (number) of descriptors

Not considering universal AD

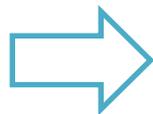
- Using ring systems and atom fragments as building blocks in structure generation.
- Introducing monotonous changing descriptors (MCDs)

Challenges

Chemical structure generation

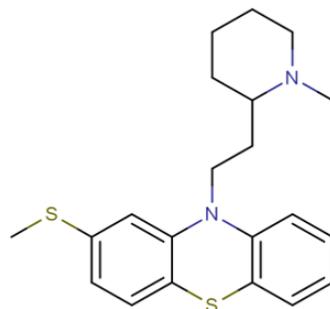
Treating limited variety (number) of descriptors

Not considering universal AD

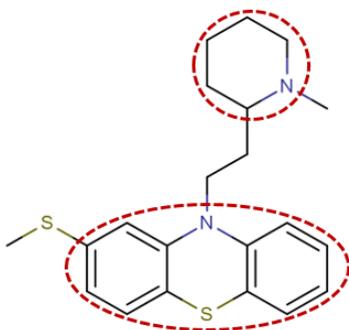


Building blocks (number and kinds)

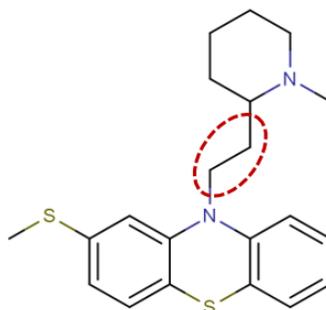
- Ring systems in the training dataset.
- Elements in the training dataset.



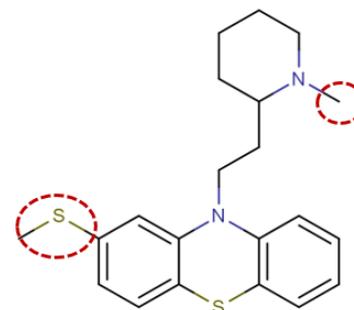
Thioridazine



Ring systems

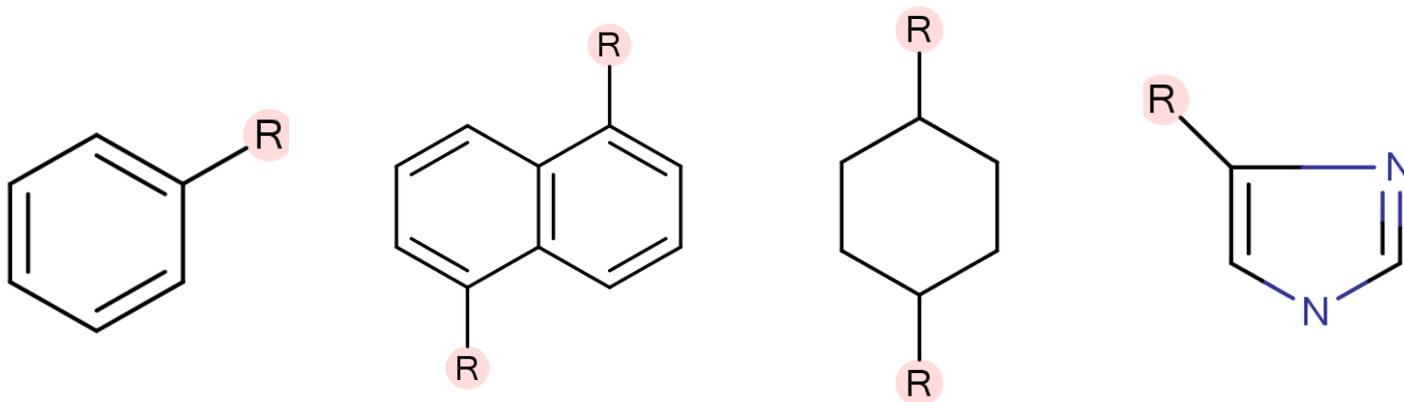


Linker



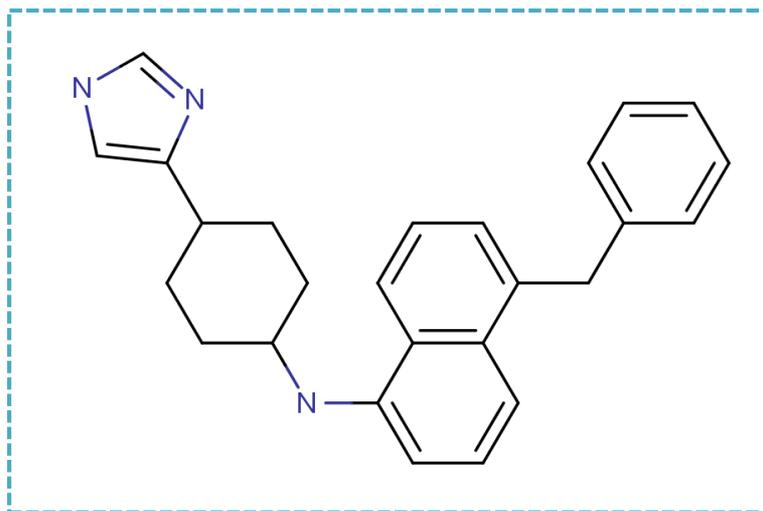
Side chains

Ring systems and atom fragments are combined to form a chemical graph.

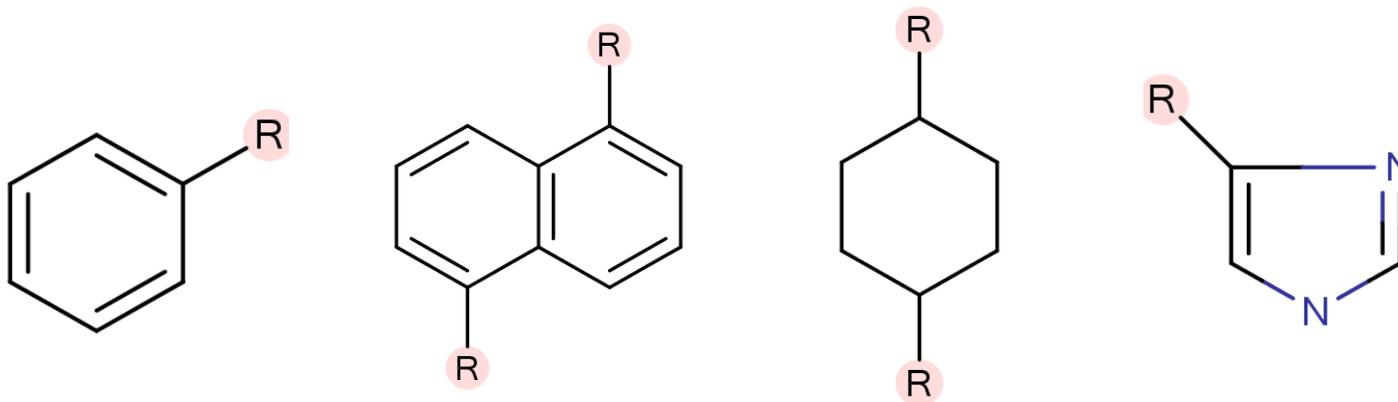


CH₂

NH

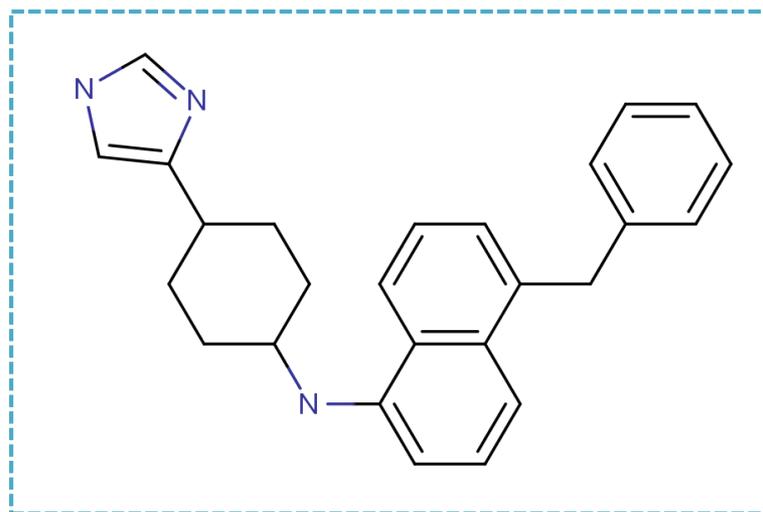


- Generating duplicate structures
 - Combinatorial explosion



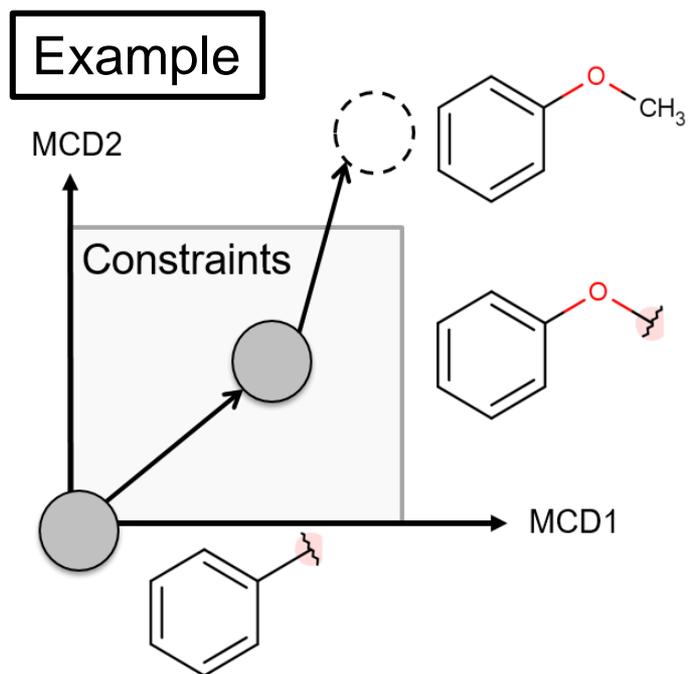
CH₂

NH



- Modifying the canonical construction path method to treat building blocks
 - Assure the uniqueness and exhaustiveness of the generated structures
- Using reduced graphs instead of ring systems
 - For speeding up graph operation during structure generation
- Combining building blocks in a tree-like approach.

- Monotonous changing descriptors (MCDs)
 - MCDs are descriptors whose values change monotonously by adding a building block to a growing structure.
 - molecular weight, topological indices



- MCDs: 409 extracted from DRAGON 5 (790 descriptors)
- Data set: Ligands for alpha 2A adrenergic receptor (GVK)
 - y: pK_i
 - training data: 500, test data: 143
- Regression method: partial least squares regression

	Opt. Compt. ^a	Q_{5fold}^2	$RMSE_{cv}$	R^2	$RMSE_{pred}$	R^2_{pred}
MCD	10	0.832	0.354	0.891	0.385	0.836
DRAGON	11	0.859	0.324	0.918	0.347	0.867

- Propose efficient structure generation algorithms by combining ring systems and atom fragments
- Using all MCDs in DRAGON has molecular description ability compatible to the comprehensive descriptors
 - High predictability of a PLS regression model for alpha 2A adrenoceptor.

Goal To develop a inverse QSPR/QSAR methodology that overcomes challenges as follows:

Challenges

Obtaining \mathbf{x} information from y

Not considering applicability domain

Poor predictability by MLR

- Probability density function (PDF)
 - Gaussian mixture models (GMMs)
- Pseudo nonlinear regression methodology
 - GMMs and cluster-wise multiple linear regression (GMMs/cMLR)

Challenges

Obtaining \mathbf{x} information from y

Not considering applicability domain

Poor predictability by MLR

- GMM: $p(\mathbf{x})$

y value



Posterior density: $p(\mathbf{x}|y)$

- GMMs/cMLR: $p(y|\mathbf{x})$

In order to consider AD

- GMM: $p(\mathbf{x})$

y value



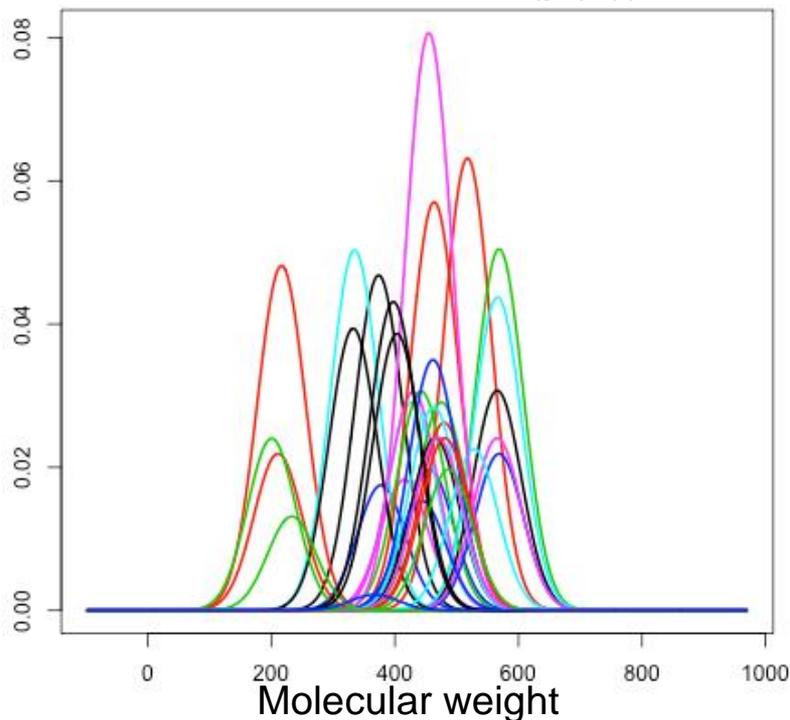
Posterior density: $p(\mathbf{x}|y)$

In order to consider AD

- GMMs/cMLR: $p(y|\mathbf{x})$

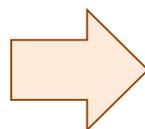
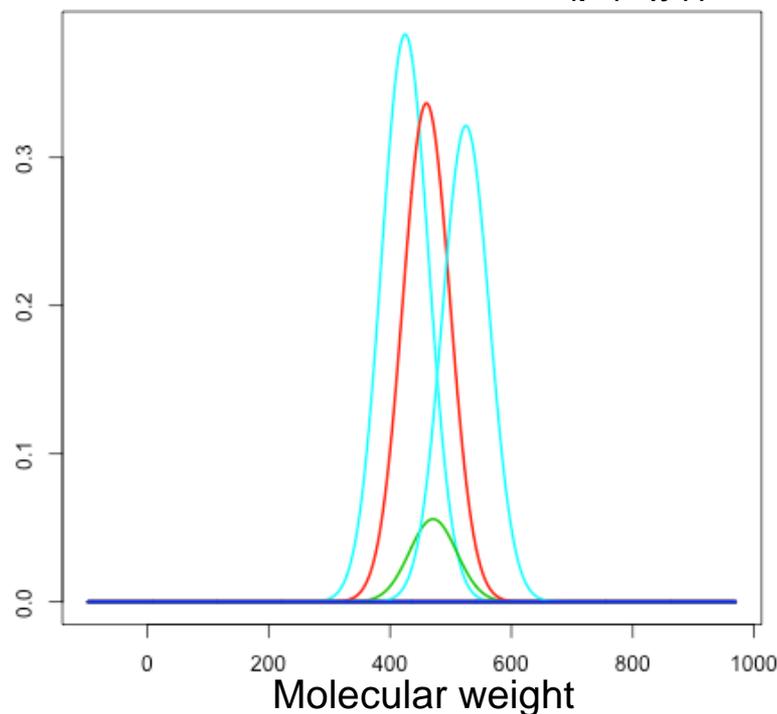
$$p(\mathbf{x}) = \sum_{i=1}^M \pi_i N(\mathbf{x} | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

Prior distribution ($p(\mathbf{x})$)



$$p(\mathbf{x} | y) = \sum_{i=1}^M \omega_i N(\mathbf{x} | \Delta\{\mathbf{A}^T \sigma^{-2} y + \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i\}, \Delta)$$

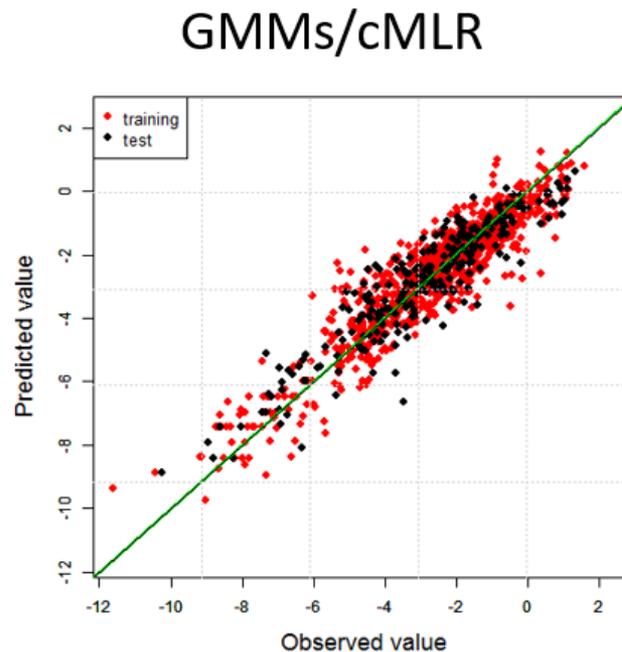
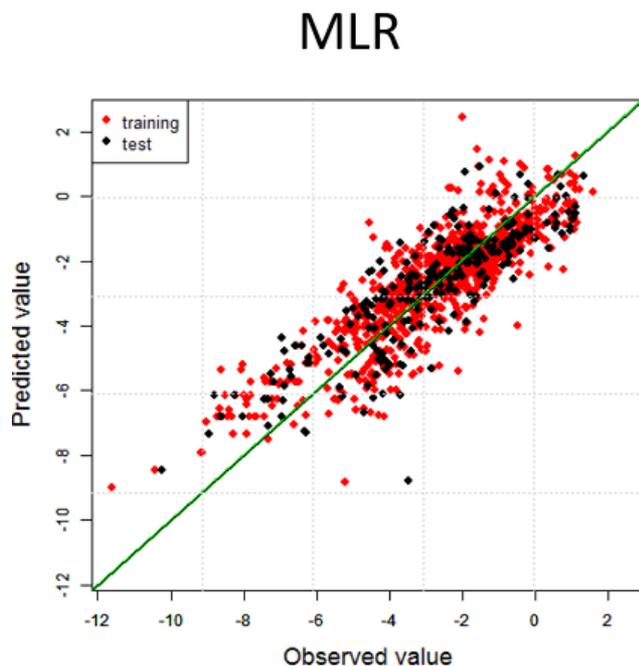
Posterior distribution ($p(\mathbf{x}|y)$)

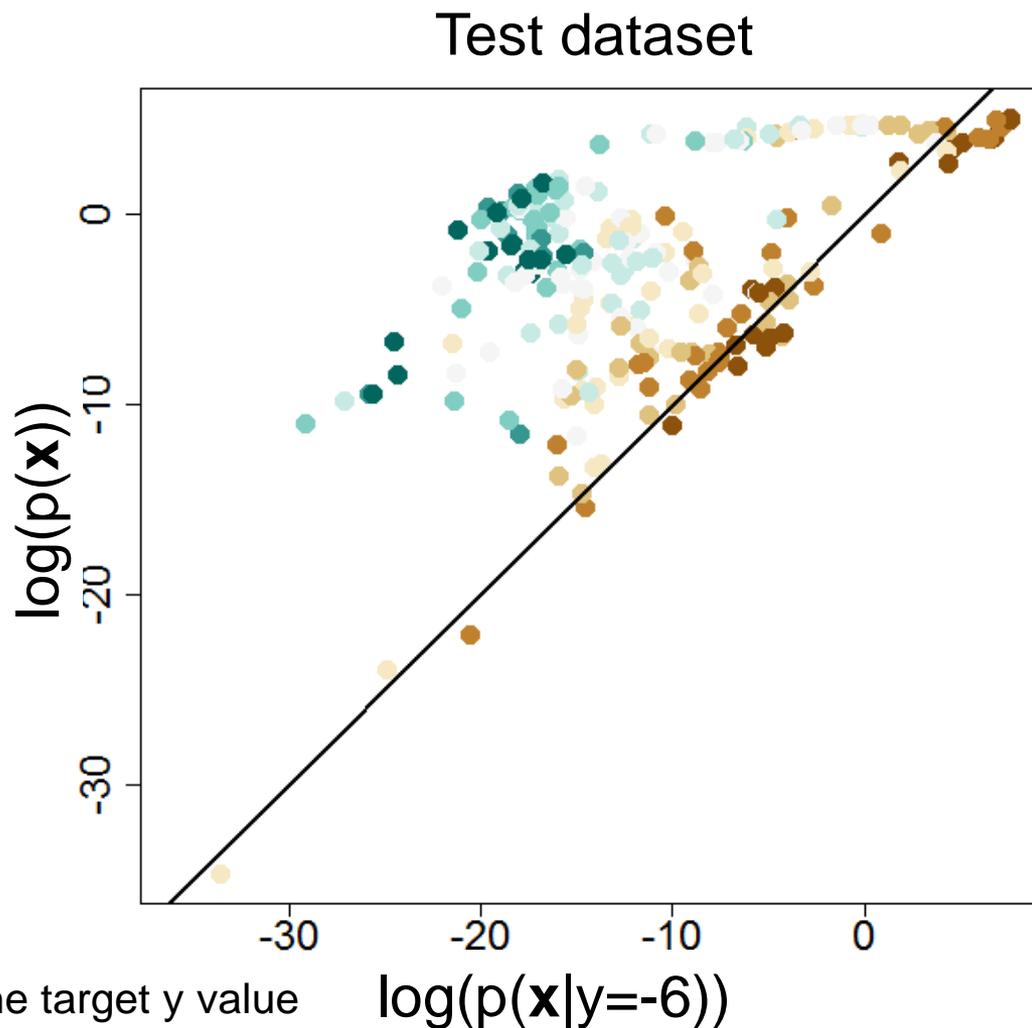


- Aqueous solubility dataset
 - training data: 900
 - test data: 254
 - Objective variable: LogS
 - Descriptors (6):
 - Molecular weight (MW)
 - Hydrogen bond donor (HBD)
 - Hydrogen bond acceptor (HBA)
 - Number of rings (CIC)
 - Topological polar surface area (TPSA)
 - Number of rotatable bonds (nBR)

- 7 Gaussians formed a prior PDF: $p(\mathbf{x})$
- With the Gaussians, GMMs/cMLR model was constructed.

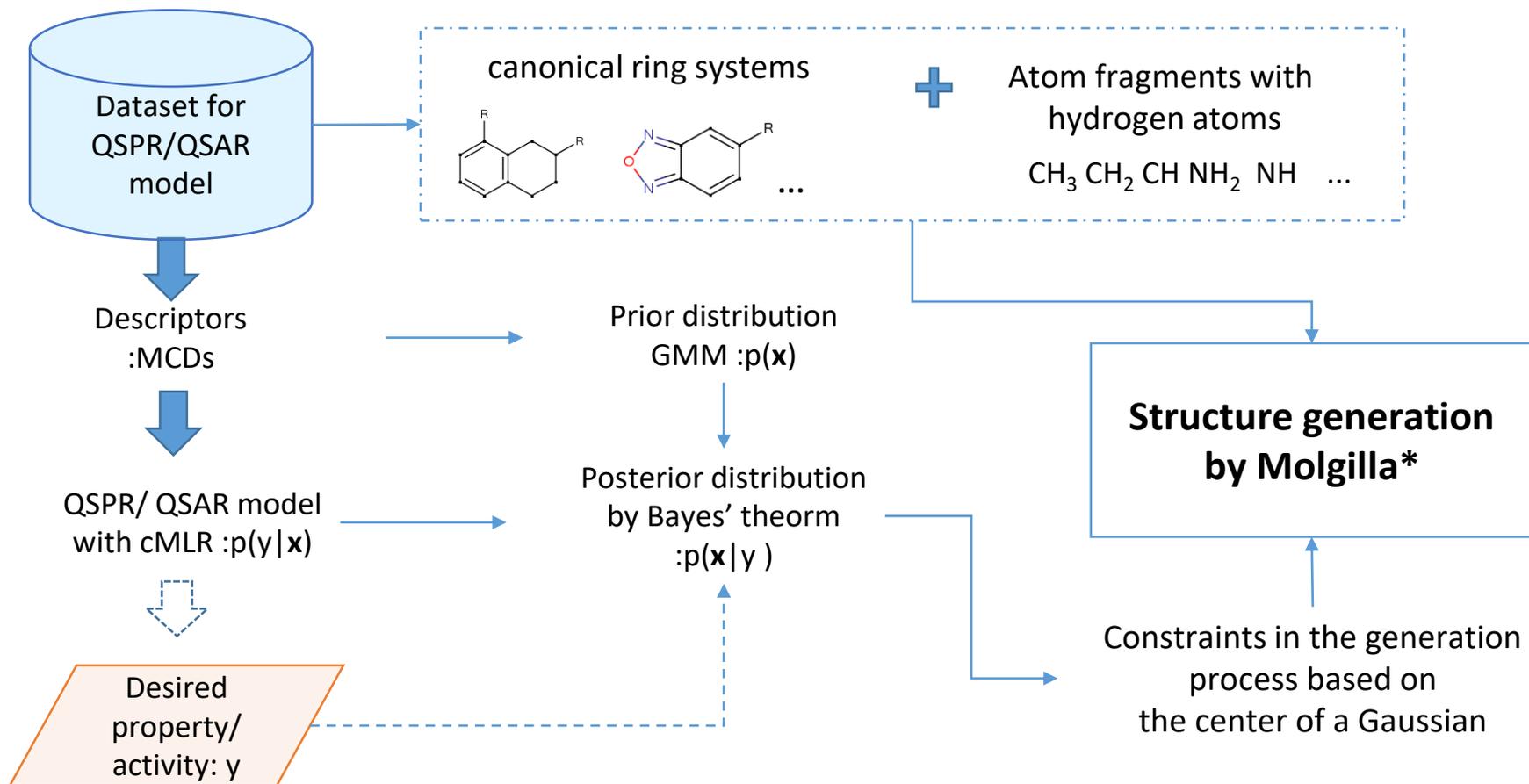
	R^2	RMSE	R_{pred}^2	$\text{RMSE}_{\text{pred}}$
MLR	0.736	1.061	0.722	1.131
GMMs/cMLR	0.853	0.791	0.854	0.820





- Could inherit the prior distribution's feature

- AD was considered by the posterior PDF with GMMs and cluster-wise multiple linear regression(cMLR)
- Posterior PDF gains information from the degree of closeness to a target y value.
- GMMs/cMLR showed better predictability than MLR did.



*Molgilla is a structure generator developed by our lab.

- Target: Thrombin
- Dataset: from ChEBML 20, 1705 samples annotated with pK_i (inhibition constant)
 - confidence score > 7
 - bioactivity type K_i
 - assay type = B
- Elimination of peptide (more than 10 amide bonds or $MW > 1,000$)
- Descriptors: 27 MCDs

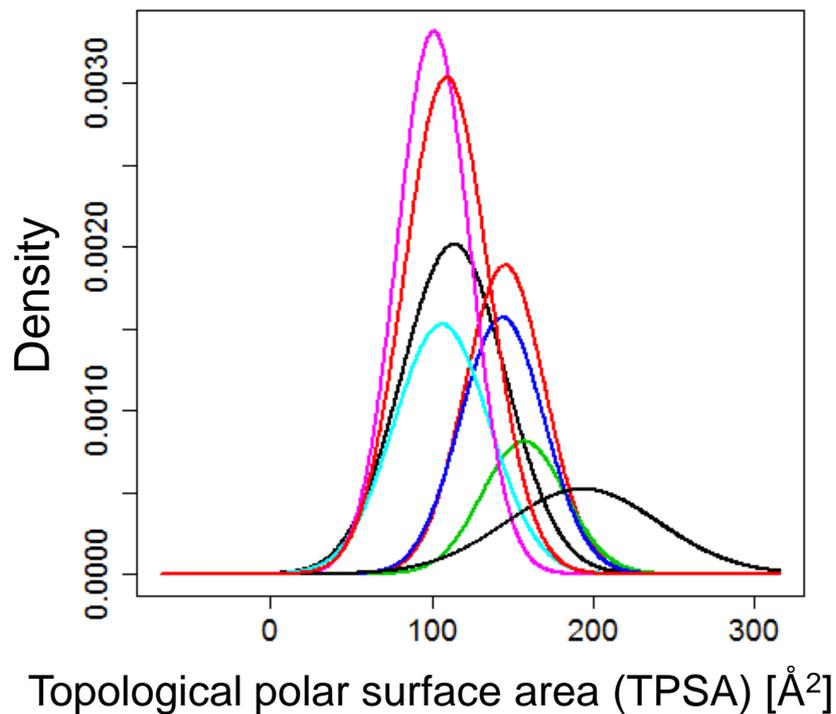
CIC	R05	aR	ZM1V	nBM	nHAcc Lipin	nCH ₂ R ₂	nCHR ₃	nCH ₃ R
nCH ₃ X	nOH	=O	nArNR ₂	nArCO	TPSA	LL	LD	LP
AA	AP	AN	DD	RL	RA	RD	RP	RR

Blue: Sum of topological distance-based descriptors

- inspired by the Chemically Advanced Template Search (CATS)

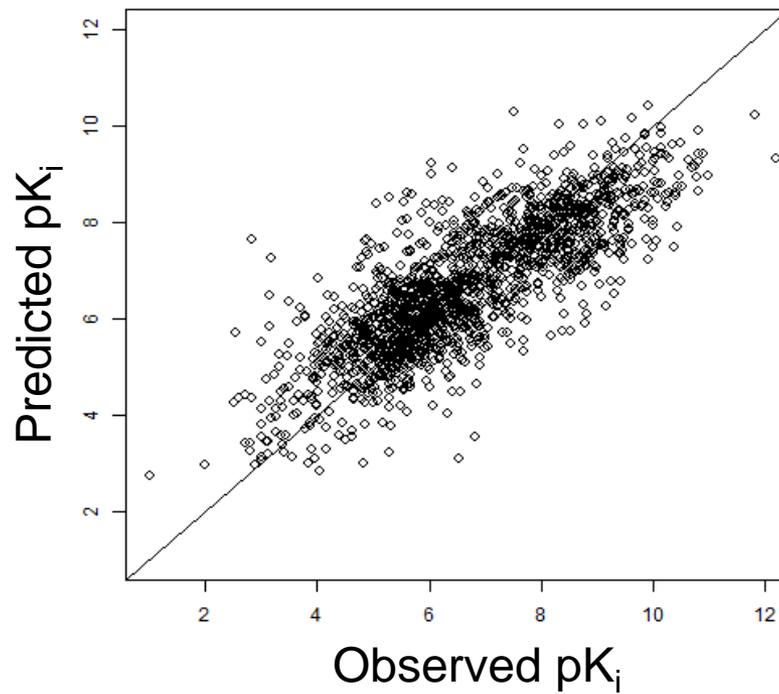
$p(\mathbf{x})$

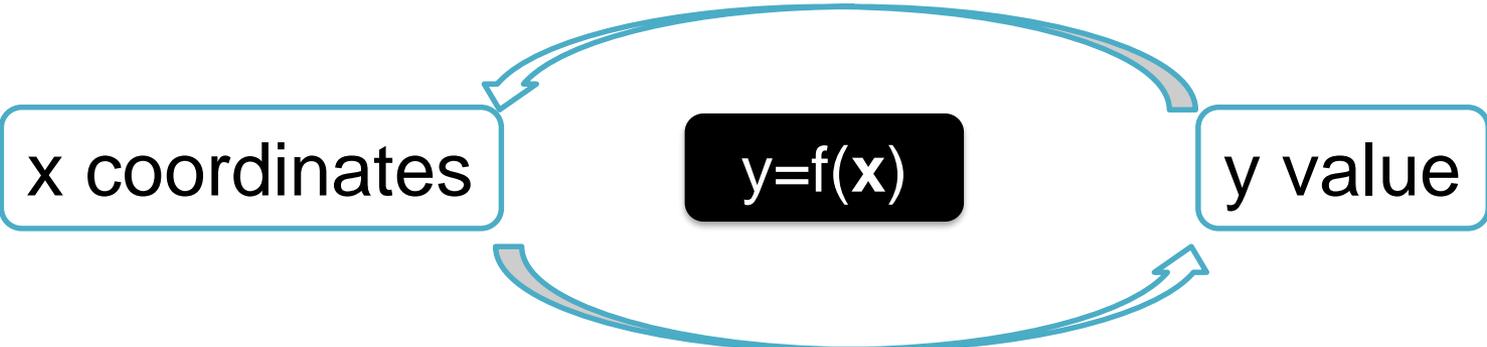
- 8 Gaussians formed $p(\mathbf{x})$



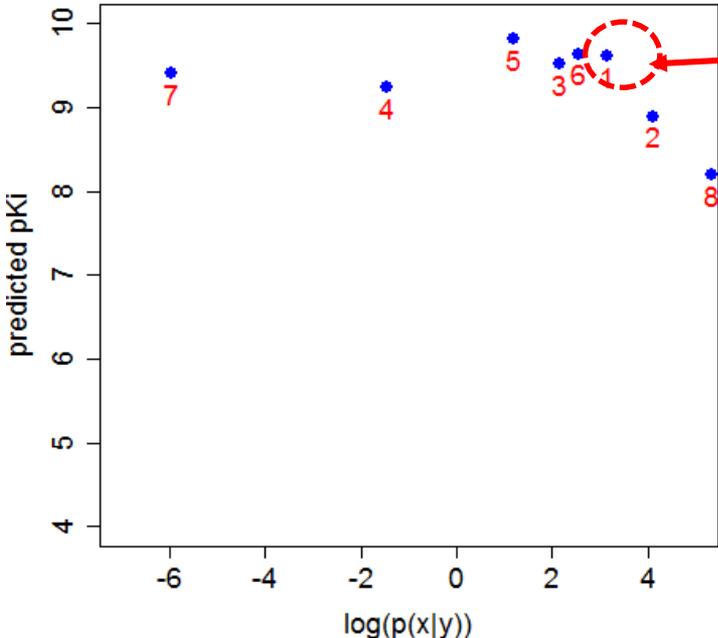
$p(y|\mathbf{x})$

- RMSE: 0.993
- R^2 : 0.656





$$p(\mathbf{x}|y = 11)$$



Target Gaussian for structure generation.

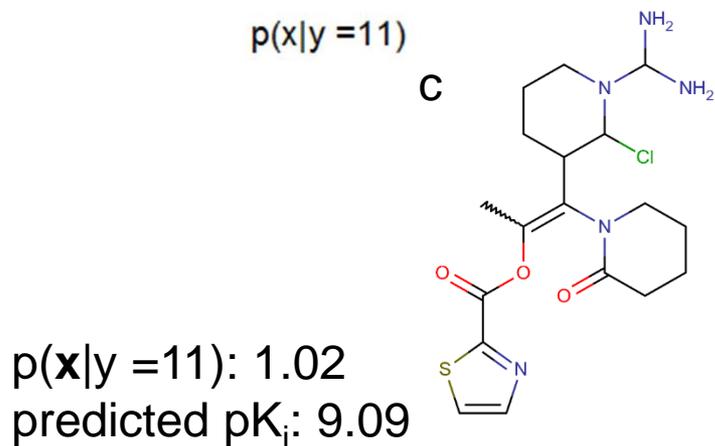
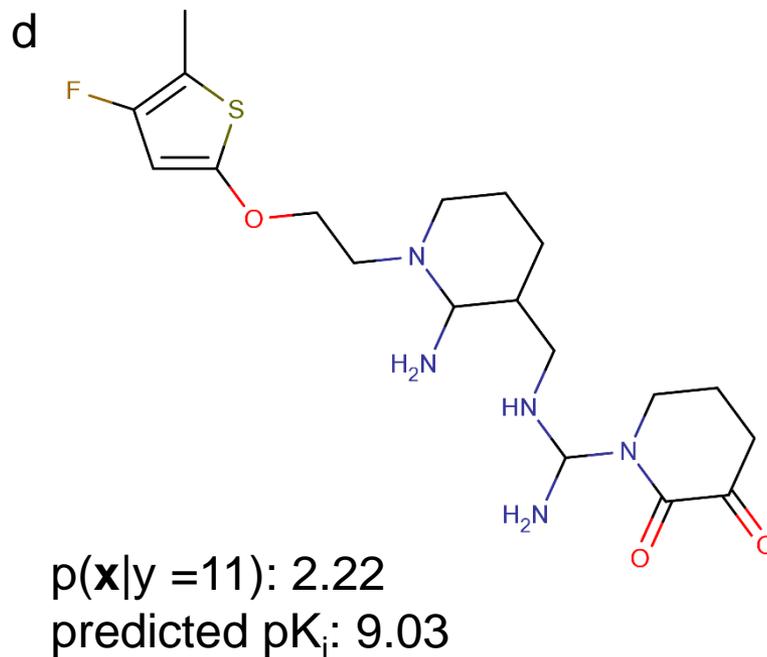
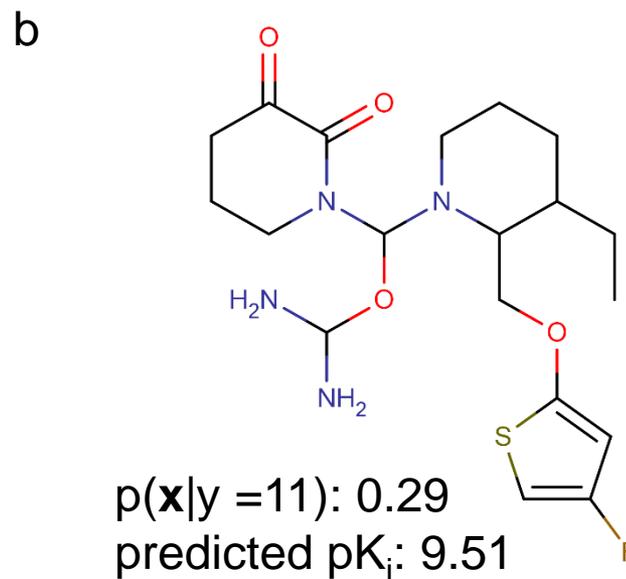
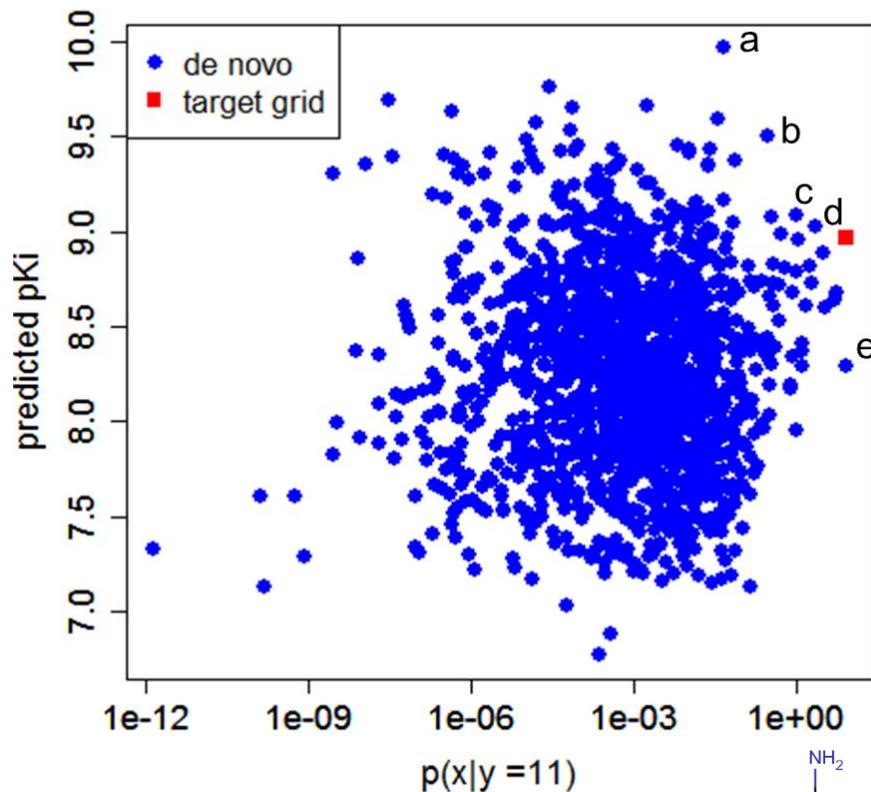
Numbers: the Gaussian centers

- Stochastic generation was conducted in order to generate structures having more building blocks.
- Prohibition rules were applied to avoid generating structures having reactive and unstable substructures
- Ring systems: 289
- Atom fragments: C, N, O, F, Cl, Br, I

Generation Result

Generated structures

1,739



- Structure generation system based on inverse QSPR/QSAR was constructed
 - A Gaussian center of posterior PDF $p(\mathbf{x}|y)$ is selected for structure generation constraints
- As a practical application, small molecules for thrombin inhibitor candidates were generated.

- Efficient chemical graph construction algorithms were introduced.
 - ring systems and atom fragments combination
 - constraints by MCDs are considered during generation
- Inverse QSPR/QSAR analysis methodology was proposed.
 - by introducing PDFs with GMMs/cMLR
 - AD consideration
 - higher predictability than with MLR
- A structure generation system by combining these methodologies was proposed in order to generate chemical structures *de novo*.

- Dr. Miyao, Tomoyuki
- Dr. Kaneko, Hiromasa
- Dr. Escobar, Matheus
- Dr. Tanaka, Kenichi

- Prof. Dr. Schneider, Gisbert
- Dr. Schneider, Petra