Catalyst Design Using Machine Learning

2020.10.10 Literature Seminar B4 Junichi Taguchi

Contents

1. Introduction

2. Data-Driven Molecular Design in Asymmetric Catalysis (Bull. Chem. Soc. Jpn. 2019, 92, 1701)

3. Prediction of higher-selectivity catalysts by computer-driven workflow and machine learning

(Science 2019, 363, eaau5631)

Contents

1. Introduction

2. Data-Driven Molecular Design in Asymmetric Catalysis (Bull. Chem. Soc. Jpn. 2019, 92, 1701)

3. Prediction of higher-selectivity catalysts by computer-driven workflow and machine learning

(Science 2019, 363, eaau5631)

What is Machine Learning (ML)

- Tool to find rules and patterns in the data
- Identification and prediction are the main uses
- It is likely to be more accurate than traditional methods



www.homesciencetools.com > flasks < このページを訳す

Chemistry Flasks | Shop for All Your Science Flask Needs on ...

Results 1 - 38 of 38 - Chemistry flasks in all shapes and sizes: 50-2000ml, volumetric, distilling, Erlenmeyer, filtering & more. Sets or individual flasks from Pyrex and HST.

Flask Stand, polypropylene · Erlenmeyer Flask, 250 ml · Erlenmeyer Flask, 100 ml

よく一緒に購入されている商品



総額: ¥1.432 ポイントの合計: 14 pt (1%) 3点ともカートに入れる



What is Machine Learning (ML)



See also 160916_LS_Yuki_NAKAGAWA and 180411_LS_Yusuke_Imamura

Catalyst Design Using Molecular Field Analysis (MFA)

MFA: A regression analysis between <u>an objective variable</u> and molecular fields calculated from 3D-molecular structures

an objective variable = the enantiomeric ratios of products for the purpose of molecular design in asymmetric catalysis

Procedure

1

2

3

Construction of an in silico library of catalysts

Calculation of molecular fields

Collection of the experimental data

Application of machine learning methods to generate models

step 1: Construction of an in silico library of catalysts



in silico library

step 2: Calculation of molecular fields

regression analysis

$$y = \beta_0 + \beta_1 x_1 + \cdot \cdot + \beta_n x_n$$

y : the logarithms of product enantiomeric ratios (k_{rel}) (i.e. $\Delta\Delta G^{\ddagger} = -RTlogk_{rel}$)

x_n: parameters for chemical properties of compounds (i.e. "descriptor(記述子)")

β: regression coefficient (Constants determined by performing a regression analysis)

step 2: Calculation of molecular fields



step 4: Application of machine learning methods to generate models



Key Points in This Lecture

<u>Differences in molecular field analysis methods</u>



Differences in the amount of the training data

Bull. Chem. Soc. Jpn. 2019, 92, 1701

Science **2019**, 363, eaau5631

using experimental approach

a large amount of reactions

a <u>small</u> amount of reactions using theoretical approach

1) Yamaguchi, S.; Sodeoka, M. Bull. Chem. Soc. Jpn. 2019, 92, 1701

Dr. Mikiko Sodeoka and Prof. Scott E. Denmark



- 1981 B.Pharm. @ Chiba University
- 1983 M.Pharm. @ Chiba University (Prof. T. Hino)
- 1989 Ph.D Pharmacy @ Chiba University (Prof. M. Nakagawa)
- 1992- Research Associate @ University of Tokyo (Prof. M. Shibasaki)
- 1999- Associate Professor @ University of Tokyo
- 2000- Professor @ Tohoku University
- 2006- Chief Scientist, Synthetic Organic Chemistry Laboratory, RIKEN

Research Interests :

transition metal chemistry/organic synthetic chemistry/chemical biology

- 1975 S.B. @ Massachusetts Institute of Technology
- 1980 D.Sc.Tech @ ETH Zurich (Prof. Albert Eschenmoser)
- 1980- Assistant Professor @ University of Illinois
- 1986- Associate Professor @ University of Illinois
- 1987- Professor @ University of Illinois

Research Interests :

synthetic organic chemistry/invention of new organic reactions/ total synthesis of natural products

1) https://chemistry.illinois.edu/sdenmark
2) https://www.ist.go.jp/erato/sodeoka/english/profile/index.html

Contents

1. Introduction

2. Data-Driven Molecular Design in Asymmetric Catalysis (Bull. Chem. Soc. Jpn. 2019, 92, 1701)

3. Prediction of higher-selectivity catalysts by computer-driven workflow and machine learning

(Science 363, eaau5631)

Model Reaction System (step 1)









R = Me, Et, Bn, *i*Pr, Cy, *t*Bu (6 substrates)



Yamaguchi, S.; Sodeoka, M. *Bull. Chem. Soc. Jpn* **2019**, *92*, 1701
Hamashima, Y.; Yagi, K.; Takano, H.;Tamás, L.; Sodeoka, M. *J. Am. Chem. Soc.* **2002**, *124*, 14530







Optimization of Pd-enolate structures (B3LYP/LANL2DZ(Pd) and 6-31G(d) level)

Alignment of the set of intermediates

Atoms except for the β -ketoester and equatorial Ar-groups on the ligands were removed.



Application of ML Methods to Generate Models (step 4)



Application of ML Methods to Generate Models (step 4)

Si-face

Re-face



Molecular Design to Improve Enantioselectivity



Molecular Design to Improve Enantioselectivity



Better than any of the data in the training set

Short Summary

Molecular field analysis methods:

Intermediate structures in an asymmetric induction step were employed for MFA



• The amount of the training data:

<u>A little amount of reactions (24 reactions)</u> were used as a training data

Contents

1. Introduction

2. Data-Driven Molecular Design in Asymmetric Catalysis (Bull. Chem. Soc. Jpn. 2019, 92, 1701)

3. Prediction of higher-selectivity catalysts by computer-driven workflow and machine learning

(Science 2019, 363, eaau5631)

The catalyst scaffold (step 1)

target: The BINOL (1,1'-bi-2-naphthol)-derived family of chiral phosphoric acids



- synthetic accessibility
- ease of diversification by installation of an array of substituents
- the acidity of the phosphoryl group can be unsaturated
- the backbone can be unsaturated or saturated
- can be used for a vast number of synthetically useful reactions

806 chiral phosphoric acid catalysts in the silico library



1) Yamaguchi, S.; Sodeoka, M. *Bull. Chem. Soc. Jpn.* **2019**, *92*, 1701 2) Zahrt, A. E.; Henle, J. J.; Rose, B. T.; Wang, Y.; Darrow, W. T.; Denmark, S. E. *Science* **2019**, *363*, eaau5631

Average Steric Occupancy (ASO) (step 2)



The grid points that included the van der Waals radii of any atoms were counted as 1, or were otherwise counted as 0.

Average Steric Occupancy (ASO) (step 2)



ASO and electronic descriptors for reactants and products (+catalyst) \rightarrow individual reaction profiles including substrate properties.

Universal Training Set (UTS) (step 3)

Conventional methods

• A training set was chosen randomly from experimental data.

 \cdot There is a risk that the model would not be sufficiently accurate, depending on the variability of the results used for the training set.



Universal Training Set (UTS) (step 3)



24 catalysts



Model Reaction (step 3)



1) Henle, J. J.; Zahrt, A. E.; Rose, B. T.; Darrow, W. T.; Wang, Y.; Denmark, S. E. *J. Am. Chem. Soc.* **2020**, *142*, 11578 2) Ingle, G. K.; Mormino, M. G.; Wojtas, L.; Antilla, J. C. Org. Lett. **2011**, *13*, 4822

Test Catalysts with Averages for All Substrate Combinations (step 3)



Modeling Study

Could this tool be used to predict the results of either new substrate combinations or new catalysts?



Modeling Study

Could this tool be used to predict the results of either new substrate combinations or new catalysts?

test set

a) new substrate



b) new catalyst



c) new substrate & new catalyst





Modeling Study

Could this tool be used to predict the results of either new substrate combinations or new catalysts?



• Mean Absolute Deviation (MAD) : The average of the absolute value of the difference between the predicted and actual values

$$MAD = \frac{1}{n} \sum_{i=1}^{n} |x_i - \overline{x}|$$

- a) MAD: 0.161 kcal/mol
- b) MAD: 0.211 kcal/mol
- c) MAD: 0.236 kcal/mol

→ The model made good predictions

 $\label{eq:shared} \begin{array}{l} \Delta\Delta G^{\ddagger} = 3.0 \ kcal/mol \rightarrow about \ 99 \ \% \ ee \\ 0.5 \ kcal/mol \rightarrow about \ 40 \ \% \ ee \end{array}$

Summary

The main differences between two methods



1) Yamaguchi, S.; Sodeoka, M. *Bull. Chem. Soc. Jpn.* **2019**, 92, 1701 2) Zahrt, A. E.; Henle, J. J.; Rose, B. T.; Wang, Y.; Darrow, W. T.; Denmark, S. E. *Science* **2019**, 363, eaau5631

Appendix

ML and Organic Chemistry



Prediction of organic reaction outcome

See also 160916_LS_Yuki_NAKAGAWA and 180411_LS_Yusuke_Imamura

1) Szymkuć, S.; Gajewska, E.; Klucznik, T.; Molga, K.; Dittwald, P.; Startek, M.; Bajczyk, M.; Grzybowski, B. A. Angew. Chem. Int. Ed. 2016, 55, 5904.

2) Wei, J. N.; Duvenaud, D.; Aspuru-Guzic, A. ACS Cent. Sci. 2016, 2, 725.

Computer-assisted synthetic planning

Conventional Process of Catalyst Design



Catalyst/Substrate Optimization



1) Jang, K. J.; Hutson, G. E.; Johnson, R. C.; McCusker, E. O.; Cheong, P. H. -Y.; Scheidt, K.A. *J. Am. Chem. Soc.* **2014**, *136*, 76.

Proposed Reaction Mechanism



Regularized Regression Model

An analysis method that reduces the estimate by adding constraints to the ordinary least squares method

Features:

- 1) Calculation of estimation
- 2) Variable selection

Examples:

1) Ridge regression

- 2) Lasso regression
- 3) Elastic Net regression





A Example of Least Squares Method

Ridge Regression and Lasso Regression

Ridge regression

It is used to improve the accuracy of the model.

Lasso regression

The number of variables included in the model is limited, making it easier to interpret.



RSS: Residual Sum of Squares

Elastic Net has the strong points of both the ridge regression and the lasso regression.

Visualized important structural information with the structures



Data Used for the MFA



Catalyst/Substrate	Me	Et	Bn	<i>i</i> Pr	Су	<i>t</i> Bu	Bzh
1Pd	36 (0.44)	28 (0.34)	40 (0.49)	38 (0.47)	46 (0.59)	52 (0.68)	77 (1.19)
2Pd	60 (0.81)	63 (0.86)	69 (0.99)	77 (1.19)	77 (1.19)	81 (1.32)	94 (2.0)
3Pd	38 (0.47)	30 (0.36)	40 (0.49)	48 (0.61)	55 (0.72)	61 (0.83)	83 (1.39)
4Pd	29 (0.34)	19 (0.22)	33 (0.40)	29 (0.35)	37 (0.46)	37 (0.45)	75 (1.13)
5Pd	60 (0.8)	63 (0.86)	70 (1.0)	78 (1.23)	79 (1.24)	81 (1.32)	97 (2.44)

Arrhenius equation



 $\Delta\Delta G^{\ddagger}$: a difference in free energy between competing transition states leading to different enantiomers

 $k_{rel} = \frac{100 + \% ee}{100 - \% ee}$, R: gas constant, T: temperature

NFSI conformation of the transition states





TS-1Pd-*i*Pr_S

The lowest energy conformers are showed in this page. Distances between oxygen atoms on NFSI and some hydrogen atoms on the Pd-enolate complexes in the structures The units of atomic distances are Å. NFSI conformations are different.

Distortion Energy



Average Steric Occupancy (ASO) (step 2)



1) Zahrt, A. E.; Henle, J. J.; Rose, B. T.; Wang, Y.; Darrow, W. T.; Denmark, S. E. Science 2019, 363, eaau5631

Electrostatic Potential (ESP) (step 2)

Conventional methods (Hammett parameters)

 $\boldsymbol{\cdot}$ represent the through-bond electronic perturbation a substituent has on a system

 However, the 3,3'- substituents in the *in silico* library are <u>too diverse to</u> <u>be represented with experimentally derived Hammett parameters</u>

	substituent	$\sigma_{\mathbf{m}}$	σ_{p}
1.	BF ₂	0.32	0.48
2.	Br	0.39	0.23
3.	GeBr ₃	0.66	0.73
4.	SiBr ₃	0.48	0.57
5.	Cl	0.37	0.23
6.	HgCl	0.33	0.35
7.	SÕ ₂ Cl	1.20	1.11
8.	SCÍ	0.44	0.48
9.	ICl ₂	1.10	1.11

 \rightarrow a new calculable parameter had to be developed that reflects the perturbation of the substituent on a charged particle: ESP

1) Zahrt, A. E.; Henle, J. J.; Rose, B. T.; Wang, Y.; Darrow, W. T.; Denmark, S. E. *Science* **2019**, *363*, eaau5631 2) Henle, J. J.; Zahrt, A. E.; Rose, B. T.; Darrow, W. T.; Wang, Y.; Denmark, S. E. *J. Am. Chem. Soc.* **2020**, *142*, 11578 3) Hansch, C.; Leo, A.; Taft, R. W. *Chem. Rev.* **1991**, *91*, 165

3,3'- substituents

Electrostatic Potential (ESP) (step 2)



Example MIF calculated for 4-nitrobenzyltrimethylammonium cation.

1) Zahrt, A. E.; Henle, J. J.; Rose, B. T.; Wang, Y.; Darrow, W. T.; Denmark, S. E. *Science* **2019**, *363*, eaau5631 2) Henle, J. J.; Zahrt, A. E.; Rose, B. T.; Darrow, W. T.; Wang, Y.; Denmark, S. E. *J. Am. Chem. Soc.* **2020**, *1*42, 11578

Electrostatic Potential (ESP) (step 2)



Evaluation of ESP_{MAX} descriptor by correlating relative ESP_{MAX} with Hammett parameters.

1) Zahrt, A. E.; Henle, J. J.; Rose, B. T.; Wang, Y.; Darrow, W. T.; Denmark, S. E. *Science* **2019**, *363*, eaau5631 2) Henle, J. J.; Zahrt, A. E.; Rose, B. T.; Darrow, W. T.; Wang, Y.; Denmark, S. E. *J. Am. Chem. Soc.* **2020**, *1*42, 11578

Universal Training Set (UTS) (step 3)



Principal Component Analysis (PCA)

Principal component analysis (PCA) is a technique for reducing the dimensionality of such datasets, increasing interpretability but at the same time minimizing information loss.



Kennard-Stone Algorithm (step 3)

This algorithm is one of the famous ways to select a sample of training data in an evenly distributed manner







Calculate Euclidean distances between the not yet selected samples and all the samples selected so far. The minima of those distances become the representative distance for each sample.

0.6

Machine Learning Method (step 4)



Entry	Model	Q ² (k=5)	R ²	MAD (Train, $\Delta\Delta G$)	MAD (SubTest, ∆∆G)	MAD (CatTest, ΔΔG)	MAD (Sub/CatTest, ∆∆G)
1	RF	0.731	0.991	0.048	0.211	0.211	0.294
2	SVR_POLY2	0.748	0.953	0.096	0.161	0.211	0.238
3	LASSOLARS	0.572	0.906	0.160	0.165	0.206	0.224

Parametric method (step 4)



These parameters are determined by training sets. The training sets are never used after training.

Kernel method (step 4)



Predict the value of y when $x=x_0$

It is thought to take a value close to y of the data (\bigcirc) in the vicinity of x_0

On the other hand, there is no relationship between the data far from x_0 and the value of y when $x=x_0$.



Support Vector Machine (SVM) (step 4)

It is thought that SVM is has the best pattern classification in the common machine learning methods.



This model tends to make a wrong prediction

"margin" is a distance between a support vector and a boundary line

Examples of Good Prediction



Another Modeling Study

Could this tool be used to identify new reactions that are more selective than any reaction in the training data?

In this model study...



Note that the UTS was not used in this model study.

Another Modeling Study

Could this tool be used to identify new reactions that are more selective than any reaction in the training data?



Development of a Computer Workflow for Catalyst Optimization



1) Henle, J. J.; Zahrt, A. E.; Rose, B. T.; Darrow, W. T.; Wang, Y.; Denmark, S. E. J. Am. Chem. Soc. 2020, 142, 11578