Predicting Reaction Performance Using Machine Learning





Contents

- **1. Introduction**
- 2. Prediction in Intermolecular Reaction Screening (main paper)
- **3. Prediction in Usual Reaction Screening**

Machine Learning (ML)

(前略)...先日、あの羽生善治永世七冠に直接お会いする機会があり、<mark>AI</mark>に関する非常に興味深い話を 伺いました。羽生永世七冠は、ご自身がAI技術に対して非常に大きな興味を持っており、将棋を通じてAI 技術も研究しているそうです。<mark>AI</mark>がその膨大なデータ処理能力を通じて、人間が長年かけて築き上げて きた将棋の世界の境界をときに飛び越え、これまでとは違う将棋の姿を見せてくれるというのです。永世 七冠の羽生さんの、新しい可能性と向かい合い、面白いと思って挑戦し続ける姿勢は、ぜひ、参考にして もらいたいと思います。…(後略) 平成29年度東京大学学位記授与式 総長告辞

「AIが仕事を奪う」 脅威論への素朴な疑問

B! 🛕 PUSH通知

4月13日 (金) 9時0分 ITmedia NEWS



「人丁知能(AI)が原因で失業する」と信じて いる人は大勢います。では私たちの周りに、人工 知能に仕事を奪われた人は居るでしょうか? 少 なくとも、今の所は私の周りにはいません。もし かして私たちは、居るはずのない幽霊にただおび えているのではないでしょうか。

その上うか図診が語られ始めたのけ 第2次↓

Why ML is hot area now? Improvement in calculation ✓ Availability of big data Information sharing

⑤ 写真を拡大

ます。 オズボ

EMPL(4/18(水) 7:00配信



JOBS TO COMPUTERISATION? 」という論文 47%が機械に代替されるリスクがある」という主

然とするのも当然です。

に花開きました。たった20年で労働者の半数が歩シマンテックが、同社のEDR製品「Symantec Advanced Threat Protection (ATP) | において、機械学習を適用して標的型攻撃の検出などを行う 新機能の追加を発表した。

エ知能 シマンテック、EDR製品で標的型攻撃を検知する機械学習活用

1) https://www.u-tokyo.ac.jp/gen01/b_message29_09_j.html

2) http://www.itmedia.co.jp/news/articles/1804/13/news018.html 3) http://ascii.jp/elem/000/001/665/1665363/

What is ML?

Machine learning is one of research topic in Artificial Intelligence (AI).

Machine (computer) learns pattern from large amount data, and then make a model for classification or regression (回帰).

cf. Al...Intelligence demonstrated by machine



Application in Daily Life



Recommendation in Amazon





Siri

Self driving technology

Application in Chemistry

Biomedical science

Virtual screening of libraries of drug-like molecules for biological function.

- See also 160723_LS_Shunichiroh_KATOH
- Organic chemistry

Assist with synthetic planning via retrosynthetic pathway.

- See also 160916_LS_Yuki_NAKAGAWA

Predict the product(s) of chemical reactions given a set of reactants and conditions. (Wei, J. N.; Duvenaud, D.; Aspuru-Guzik, A. *ACS Cent. Sci.* 2016, *2*, 725.)



Predict the products of textbook questions.

Challenges... It is difficult to use for non-specialists Multidimensionality of chemical structure

Predicting Reactions

• Predict from substrate scope



Chemical theory to understand the rule is needed.

1) Nielsen, M. L.; Ugaz, C. R.; Li, W.; Doyle, A. G. J. Am. Chem. Soc. 2015, 137, 9571

Predicting Reactions

• Linear regression model (線形回帰)

Site selective C-H amination (Bess, E. N.; Du Bois, J. et al., J. Am. Chem. Soc. 2014, 136, 5783)



Prof. Abigail G. Doyle



1998-2002 A.B. and A.M. @Harvard Univ. (Prof. Eric Jacobsen) Research topic: iron-catalyzed epoxidations of alkenes

2002-2003 Pre-D. @Stanford Univ. (Prof. Justin Du Bois) Research topic: gold catalysts for the hydration of unactivated alkenes

2003-2008 Ph.D. @Harvard Univ. (Prof. Eric Jacobsen) Ph.D. Thesis: Engaging Alkyl Halides and Oxocarbenium Ions in Asymmetric Catalysis.

2008-2013 Assist. Prof. @Princeton Univ. Department of Chemistry 2013-2017 Associate Prof. 2017- Professor



Ni-catalyzed Cross Coupling





See also 171118_LS_Daiki_Kamakura ⁹

Contents

1. Introduction

2. Prediction in Intermolecular Reaction Screening (main paper)

3. Prediction in Usual Reaction Screening

Predicting reaction performance in C–N cross-coupling using machine learning

Derek T. Ahneman¹, Jesús G. Estrada¹, Shishi Lin², Spencer D. Dreher^{2,*}, Abigail G. Doyle^{1,*} + See all authors and affiliations

Science 13 Apr 2018: Vol. 360, Issue 6385, pp. 186-190 DOI: 10.1126/science.aar5169

Buchwald-Hartwig Reaction ¹⁾



Useful reaction to synthesize pharmaceuticals. But incopatible with heterocycles that contain heteroatom-heteroatom bonds (ex. isoxazole).

Pharmaceuticals containing an isoxazole core



1) a) Paul, F.; Patt, J.; Hartwig, J. F. J. Am. Chem. Soc. 1994. 116, 5969.
b) Guram, A. S.; Buchwald, S. L. J. Am. Chem. Soc. 1994. 116, 7901.

Intermolecular Reaction Screening ¹⁾



16×**24**×**3**×**4** = 4608 reactions

1) Collins, K. D.; Glorius, F. Nat. Chem. 2013, 5, 597.

Ultra-High-Throughput Experimentation



Descriptor

Descriptor: Variable that descripts chemical property of the substrate.



• Additive Descriptors (*n* = 19)

(E_{HOMO}, E_{LUMO}, Dipole Moment, Electronegativity, Hardness, Molecular Volume, Molecular Weight, Ovality, Surface Area, *C3 NMR Shift, *C3 Electrostatic Charge, *C4 NMR Shift, *C4 Electrostatic Charge, *C5 NMR Shift, *C5 Electrostatic Charge, *N1 Electrostatic Charge, *O1 Electrostatic Charge, V1 Frequency, V1 Intensity)

- Aryl Halide Descriptors (*n* = 27)
- Base Descriptors (n = 10)
- Ligand Descriptors (n = 64)

Prediction Procedure



Training set and test set must be differrent to make reliable algorithm.



Prediction Results



Random forest model was best fitting.

Random Forest Model



d1, d2, d3, d4: randomly chose descriptors

See also 160723_LS_Shunichiroh_KATOH

Problems in ML

"Activity cliffs"

...modest changes in chemical structure can lead to notable changes in reaction outcome (Cruz-Monteagudo, M. *et al.*, *Drug Discov. Today* 2014, *19*, 1069.)



 \rightarrow training data need to spread across the chemical space of interest.

Out-of-Sample Prediction



Out-of-Sample Prediction

Train: 15 additives \rightarrow Test: other 8 additives



Mechanistic Analysis

It is difficult to interpret the outcome of ML. Which is the important descriptor???



1) T. Hastie, R. Tibshirani, J. Friedman, "The Elements of Statistical Learning: Data Mining, Inference, and Prediction" *Springer*, **2009**.

Variable Importance Plot



Authors said that "the descriptors suggest that the propensity of the additive to act as an electrophile influences reaction outcomes."

So author hypothesized oxidative addition of Pd to isoxazole was main side reaction, and it was confirmed by experiment.



However, it is difficult to connect *C3 shift and oxidative addition.

Problem in this method is that if multiple correlated variables are present, each will individually appear to have less importance.

Short Summary



- Yield was predicted with high accuracy.
- All descriptor was calculated only by computer.
- Reaction performance can be predicted without chemical theory.
- Interpretation of valuable importance plot.
- In this report, aryl halides and isoxazoles were used as screening conditions. These compound has large similarity in chemical structures.

Contents

- **1. Introduction**
- 2. Prediction in Intermolecular Reaction Screening (main paper)

3. Prediction in Usual Reaction Screening



Deoxyfluorination with Sulfonyl Fluorides: Navigating Reaction Space with Machine Learning

Matthew K. Nielsen, Derek T. Ahneman, Orestes Riera, and Abigail G. Doyle*®

Dehydroxyfluorination

Nielsen, M. L.; Ugaz, C. R.; Li, W.; Doyle, A. G. J. Am. Chem. Soc. 2015, 137, 9571



The system is useful for acyclic secondary alcohol.

This protocol gave poor yield for activated alcohols, because the intermediate were attacked by DBU. If the reaction didn't proceed well with "best condition", alternative condition must be investigated.

-> Reaction screening was conducted to understand the reactivity of sulfonyl fluorides.

Reaction Screening



5 sulfonyl fluorides:



sterically bulky 32×5×4 = 640 reactions

Alcohols



Alcohols





External Validation Set



Problem in Previous Method



20

Increase in Mean Squared Error (% yield)²

30

Alcohol Surface Area -Alcohol *C1 Electrostatic Charge -Base Dipole Moment -Alcohol *C1 NMR Shift -Base *N1 Electrostatic Charge -Sulfonyl Fluoride E HOMO -

10

-> too many descriptor makes the performance worse.

30

Modified Method

Descriptors relevant to the nucleophilic substitution was selected, and new categorical descriptors was added.

- Calculated descriptors
- alcohol - Alcohol - *C1 electrostatic charge/*C1 exposed area/electronegativity
- Base *N1 exposed area
- Sulfonyl fluoride *S1 electrostatic charge/*F1 electrostatic charge/*O1 electrostatic charge
- Categorical descriptors

- Alcohol – primary/secondary/tertiary/cyclic/4-membered ring/5-membered ring/6-membered ring/7-membered ring/benzylic/allylic/homobenzylic/homoallylic/ α -carbonyl/ β carbonyl/hemiacetal/amino alcohol



sulfonyl fluoride base

Summary

Machine Learning was applyed to predict reaction performance 1. Intermolecular reaction screening



Reaction perfomance can be predicted well by automated caluculation.

2. Usual substrate scope



Reaction perfomance can be predicted well by defining appropriate descriptor.

Ligands





Descriptors

• Additive Descriptors (*n* = 19)

E_{HOMO}, E_{LUMO}, Dipole Moment, Electronegativity, Hardness, Molecular Volume, Molecular Weight, Ovality, Surface Area, *C3 NMR Shift, *C3 Electrostatic Charge, *C4 NMR Shift, *C4 Electrostatic Charge, *C5 NMR Shift, *C5 Electrostatic Charge, *N1 Electrostatic Charge, *O1 Electrostatic Charge, V1 Frequency, V1 Intensity

• Aryl Halide Descriptors (n = 27)

E_{HOMO}, E_{LUMO}, Dipole Moment, Electronegativity, Hardness, Molecular Volume, Molecular Weight, Ovality, Surface Area, *C1 NMR Shift, *C1 Electrostatic Charge, *C2 NMR Shift, *C2 Electrostatic Charge, *C3 NMR Shift, *C3 Electrostatic Charge, *C4 NMR Shift, *C4 Electrostatic Charge, *H2 NMR Shift, *H2 Electrostatic Charge, *H3 NMR Shift, *H3 Electrostatic Charge, V1 Frequency, V1 Intensity, V2 Frequency, V2 Intensity, V3 Frequency, V3 Intensity

• Base Descriptors (*n* = 10)

E_{HOMO}, E_{LUMO}, Dipole Moment, Electronegativity, Hardness, Molecular Volume, Molecular Weight, Ovality, Surface Area, *N1 Electrostatic Charge

• Ligand Descriptors (*n* = 64)

Dipole Moment, *C1 NMR Shift, *C1 Electrostatic Charge, *C2 NMR Shift, *C2 Electrostatic Charge, *C3 NMR Shift, *C3 Electrostatic Charge, *C4 NMR Shift, *C4 Electrostatic Charge, *C5 NMR Shift, *C5 Electrostatic Charge, *C6 NMR Shift, *C6 Electrostatic Charge, *C7 NMR Shift, *C7 Electrostatic Charge, *C8 NMR Shift, *C8 Electrostatic Charge, *C9 NMR Shift, *C9 Electrostatic Charge, *C10 NMR Shift, *C10 Electrostatic Charge, *C11 NMR Shift, *C11 Electrostatic Charge, *C12 NMR Shift, *C12 Electrostatic Charge, *C13 NMR Shift, *C13 Electrostatic Charge, *C14 NMR Shift, *C14 Electrostatic Charge, *C15 NMR Shift, *C15 Electrostatic Charge, *C16 NMR Shift, *C16 Electrostatic Charge, *C17 NMR Shift, *C17 Electrostatic Charge, *H11 NMR Shift, *H11 Electrostatic Charge, *H3 NMR Shift, *H3 Electrostatic Charge, *H4 NMR Shift, *H4 Electrostatic Charge, *H9 NMR Shift, *H9 Electrostatic Charge, *P1 Electrostatic Charge, V1 Frequency, V1 Intensity, V2 Frequency, V2 Intensity, V3 Frequency, V3 Intensity, V4 Frequency, V4 Intensity, V5 Frequency, V5 Intensity, V6 Frequency, V6 Intensity, V7 Frequency, V7 Intensity, V84 Frequency, V8 Intensity, V9 Frequency, V9 Intensity, V10 Frequency, V10 Intensity

Isoxazole Descriptors

No.	*C3_NMR _shift	*C3_elect rostatic_ charge	*C4_NMR _shift	*C4_elect rostatic_ charge	*C5_NMR _shift	*C5_ele rostatio charge	ect *N1_elec c_ rostatic e charge	t *O1_elect _ rostatic_ charge	E_HOMC	E_LUMO
15	5 139.05	0.009	116.54	0.009	154.93	0.4	-0.2	5 -0.158	-0.248	5 -0.0459
23	3 163.09	0.721	95.32	-0.581	155.72	0.2	41 -0.38	3 -0.106	-0.260	5 -0.0587
No.	V1_frequ ncy	ie V1_intei ty	nsi dipole_ men	_mo Electi t gativ	rone vity harc	Iness ⁿ	nolecular n _volume	nolecular _weight	Ovality	surface_ar ea
1	5 892.39	95 10.2	204 3.184	221	0.15	0.1	121.96	119.123	1.152	137.07
2	697.30	62 5.0	3.44	567	0.16	0.1	164.92	171.152	1.343	195.35





Reaction of Isoxazole and Transition Metal

• Direct C-H arylation of Isoxazoles

(Fall, Y.; Reynaud, C.; Doucet, H.; Santelli, M. Eur. J. Org. Chem. 2009, 4041.)



(Shigenobu, M.; Takenaka, K.; Sasai, H. Angew. Chem. Int. Ed. 2015, 54, 9572.)



• Oxidative addition of trinsition metal to isoxazole N-O bond (Yu, S.; Tang, G.; Li, Y.; Zhou, X.; Lan, Yu.; Li, X., *Angew. Chem. Int. Ed.* **2016**, *55*, 8696.)



Reaction of isoxazole with $Pd(PPh_3)_4$



Important descriptor was extracted automatically by ML.

Combined Graph (Fluorination)

