

# 修士論文

## 有機合成に最適な酵素候補提示のための 特徴量エンジニアリングによる EC 番号予測

EC Number Prediction Using Feature Engineering  
to Present Optimal Enzyme Candidates  
in Organic Synthesis

富山県立大学大学院 工学研究科 電子・情報工学専攻

2255018 武藤 克弥

指導教員 奥原 浩之 教授

提出年月: 令和6年(2024年)2月



# 目次

図一覧	ii
表一覧	iii
記号一覧	iv
第1章 はじめに	1
§ 1.1 本研究の背景	1
§ 1.2 本研究の目的	2
§ 1.3 本論文の概要	3
第2章 有機合成と酵素	4
§ 2.1 有機合成分野の概要	4
§ 2.2 酵素と EC 番号	6
§ 2.3 酵素の探索	8
第3章 機械学習による EC 番号予測	11
§ 3.1 計算機上における化学反応の表現法	11
§ 3.2 EC 番号予測手法	13
§ 3.3 機械学習と特徴量エンジニアリング	14
第4章 提案手法	19
§ 4.1 特徴ベクトルと特徴エンジニアリングによる EC 番号予測	19
§ 4.2 EC 番号予測モデルの構築と予測	21
§ 4.3 提案手法の実装と流れ	23
第5章 実験結果並びに考察	26
§ 5.1 数値実験の概要	26
§ 5.2 実験結果と考察	27
第6章 おわりに	32
謝辞	33
参考文献	34

# 図一覧

2.1	HCHO の化学反応ネットワーク [10]	5
2.2	$\alpha$ -キモトリプシンの立体構造 [14]	6
2.3	反応の進行と必要なエネルギー [15]	6
2.4	モルヌピラビルの合成	7
2.5	EC 番号による酵素の分類	8
2.6	EC3.1.1.2 の 代表的な反応式	8
3.1	KEGG COMPOUND で取得できる構造式と MOL ファイルの例	12
3.2	rdkit を用いた化合物の情報	13
3.3	PubChemPy でグルコースの情報を取得した結果	13
3.4	SMOTE によるオーバーサンプリングのイメージ	16
3.5	決定木におけるクラス分類の様子	17
4.1	特徴ベクトルのイメージ	20
4.2	SMILE 形式のデータセット (一部抜粋および改変) []	21
4.3	モデル構築と予測の流れ	22
4.4	除外リスト	24
5.1	実験 2 の流れ	27
5.2	予備実験 1 の結果	27
5.3	EC 1.X.X に対する分類	29
5.4	EC 2.X.X, 3.X.X に対する分類	29
5.5	EC 4.X.X, 5.X.X に対する分類	30
5.6	EC 6.X.X, 1 桁目に対する分類	30
5.7	EC 番号 1 桁～3 桁目までの分類結果 1	30
5.8	EC 番号 1 桁～3 桁目までの分類結果 2	31

## 表一覧

4.1	実際の特徴ベクトル . . . . .	20
4.2	学習データの内訳と SMOTE 増加数 . . . . .	23
4.3	SMILE 反応式 (左辺) . . . . .	24
5.1	予備実験 2 で選ばれた記述子リスト . . . . .	28

# 記号一覧

以下に本論文において用いられる用語と記号の対応表を示す.

用語	記号		
正例データを正しく正と予測した数	$TP$	負例データを誤って正と予測した数	$FP$
負例データを正しく負と予測した数	$TN$	負例データ誤って正と予測した数	$FN$
正に分類されたデータのうち, 実際に正だったデータの割合	$Precision$	正に属するデータを正しく正に分類されたデータの割合	$Recall$
全体の予測精度	$Accuracy$	Precision と Recall の調和平均	$F1$
近傍データ数	$K$	着目するデータ点	$x_o$
選択されたデータ点	$x_k$	ランダムな生成値	$r$
次元数	$p$	決定木上位ノード内のデータ	$D_p$
決定木上位ノードのデータ数	$N_p$	上位ノードに対する下位の左(右)ノード	$D_{left}, D_{right}$
上位ノードに対する下位の左(右)ノード	$N_{left}, N_{right}$	$D_p$ を含むノードにおける情報利得	$IG(D_P, f)$
ノード $t$ におけるクラス $i$ のデータの割合	$p(i t)$	ノード $t$ 内のクラス数	$c$
ノード $t$ におけるジニ不純度	$I_G(t)$		

## はじめに

### § 1.1 本研究の背景

有機合成化学分野では、化学反応の予測や設計に酵素を生体触媒として利用する機会が増加している。生体触媒となる酵素は、特定の反応物 (基質) に対して触媒反応を示し、化学触媒に比べてしばしば化学反応をより効率的に進め、かつ環境に優しい条件で利用できる。実際、目的の化合物を生成するために従来では 10 ステップの合成を行っていたものを、生体触媒を取り入れることで 3 ステップまで短縮した事例がある [1]。そのような背景から、基質から目的の生成物を生み出す化学反応に対して、最適な酵素を探索することが重要視されている。

酵素には 4 桁の Enzyme Commission number (EC 番号) が割り振られており、どの反応を触媒するか、どの結合や基質に作用するかによって酵素が分類されている [2]。EC 番号は代謝経路の解析で得られる酵素の性質の特定や触媒する化学反応 (酵素反応) の探索に用いられる。代謝経路の解析では遺伝子解析によって得られたタンパク質配列から性質を特定し、酵素反応の探索では、化合物の合成を通して反応効率のよい組み合わせを探索する必要があるが、これらは実験によって正確が解析が行われる。近年、実験コストや時間を短縮する目的で、EC 番号を機械学習を用いて予測する研究が行われている。EC 番号を予測することで対象となる酵素の大まかな性質や触媒する反応物 (基質) の種類などを知ることができるため、実験による解析や探索の効率化が期待できる。EC 番号が登録されている酵素には、酵素反応やタンパク質配列といった情報がデータベース上で付与されており、これらの情報から既存のデータに対して正しい EC 番号を予測する機械学習モデルの開発が行われてきた。既存のデータに対する予測精度を高め、将来的に未知の酵素に対して自動的な EC 番号の割り当てが行われることが目指されている。

EC 番号予測が活発になる以前から情報技術を用いた解析技術が数多く存在しており、化学合成経路の解析や未知酵素と既存酵素のタンパク質の類似性探索を行うデータベースなどの技術が土台となって様々な手法が開発されてきた。EC 番号の予測手法として、タンパク質配列に着目した遺伝子ベースの手法や化合物の構造的特徴や物理・化学的な特性値の変化に着目した化学的手法がある。タンパク質配列では配列の類似性や自然言語処理によって与えられた配列がどの EC 番号に属するのかを予測する [3]。化学的手法では化合物が持っている特徴的な分子の有無に関するバイナリ値 (フィンガープリント) や、化合物に対する物理・化学的指標を複数用いた特徴ベクトルを用い、EC 番号が割り当てられた酵素反応を特徴ベクトルとして学習器に学習させることで EC 番号を予測する [4] [5]。また、EC 番号の予測範囲も 1,2 桁など局所的だったものが拡大していき、代謝経路の解析に重き

を置く遺伝子ベースの手法では4桁全てを予測できるなモデルが開発されるようになってきている。

## § 1.2 本研究の目的

本研究では、有機合成を行うことを主目的として、基質から生成物に変化する際の物理・化学的な特性値の変化に着目し、機械学習を用いて特定の化学反応に最適な酵素の候補をEC番号として予測する。物理・化学特性値はPythonライブラリのRDKitで得られる特徴量(記述子)を用いる。これらの記述子は分子量や親油性、電荷などの125種の物理・化学的特性値と特徴的な構造(分子断片)を持っているかを判定する85種のバイナリ値で構成されている。基質から生成物に変化に着目した従来手法として、物理・化学的特性値を用いた手法[5]と差分フィンガープリントを用いたものがあるが[4]、本研究はこの2つの手法の一部を組み合わせたものとなる。化合物の構造情報だけでなく、物理、化学的な指標を加えることで、化学合成時の特徴をより詳細に捉えることができると考えられる。

本研究の下準備として、初めに酵素を用いた化学反応(酵素反応)を収録した4つのデータベースからなるEC番号が割り当てられている既存の酵素反応のデータセットを提案手法に合わせて加工する。このデータセットは左辺と右辺からなる化学反応式中の化合物の化学構造を文字列で表現したもので構成されている。そのため、RDKitの210種類の記述子を用いた、基質から生成物に変化する際の特性値変化量からなる210次元の特徴ベクトルのデータセットを作成した。次に記述子間の相関係数が1の記述子の片方を削除し、欠損値が含まれる記述子を削除することで200次元の特徴ベクトルとした。提案するEC番号予測モデルは従来と同様にEC番号の1桁目から3桁目を重点的に予測するモデルとする、既存のEC番号が割り当てられている酵素反応に対して、正しいEC番号に予測ができるかの多クラス分類を行う。

初めに、データセットの訓練データに対して1桁目と2,3桁目に対する分類の検証を行い、200種類の中から分類精度に貢献する記述子を1桁目の分類、およびEC1からEC6クラスの2,3桁目分類において、20種類ずつ選出し、共通のものや分類性能に関わる特徴量約30種を選択する。分類器としてはRandom Forest (RF)を用い、記述子はForward法で1つずつ選択する。また、2,3桁分類による記述子選択ではクラス間のデータ分布に偏りが生じる懸念があることから、検証用訓練データにSMOTEを適用しオーバサンプリングを行う。次に、選択した記述子と訓練データを用いて1,2,3桁目を同時に予測する分類モデルを作成する。RFのハイパーパラメータである最大深さと、決定木数に対して、グリッドサーチでパラメータ調整を行ったのち、テストデータに対する分類精度を評価した。この際、それぞれのクラス分類における特徴選択と、グリッドサーチでは5分割交差検証を用いた。SMOTEを適用した検証用訓練データで作成したモデルで検証用テストデータを予測し、スコア平均値が最も高い記述子組み合わせ・パラメータを選択した。

数値実験では3つの実験を行った。1つ目の実験では、従来研究の類似手法で行われたEC3クラスの2,3桁目分類に対して、SMOTEを適用し、F1-macroスコアで比較することでSMOTEの有効性を確認する。2つ目の実験では、EC番号1桁目とEC1からEC6の2,3桁目クラスの記述子選択後に作成される分類モデルに対して、モデル単体での予測精度を評価する。各記述子選択で選択された20種を用いてSMOTEとグリッドサーチを適用



し, F1-macro 値で評価した. 最後の実験では2つめの実験で得られた記述子の中から共通のモノや分類性能に大きく影響する記述子約 30 種を選択し, EC 番号の 1 桁目から 3 桁目の組み合わせのクラスに対して, 多クラス分類を行う.

## § 1.3 本論文の概要

本論文は次のように構成される.

**第 1 章** 本研究の背景と目的について説明した. 背景では, 有機合成における酵素の利用について, 生体触媒を用いることのメリットとその課題について述べた. 目的では, 目的の生成物を得る際に用いる最適な酵素を予測するための, EC 番号を予測するシステムの概要について述べた.

**第 2 章** 有機合成で行われている内容や情報分野との関係, 酵素の性質と EC 番号による分類体系, および酵素探索を行うための EC 番号予測の重要性について述べる.

**第 3 章** 化学反応の機械学習を行うにあたっての計算機上での化学反応の表現法や機械学習による EC 番号予測の従来手法を述べる. また, 本研究で用いる機械学習や特徴エンジニアリング関する手法について述べる.

**第 4 章** 提案手法についての説明, および提案システムの実装と手順について説明する.

**第 5 章** 提案手法による数値実験の概要と実験結果, 考察を述べる.

**第 6 章** まとめと今後の課題について述べる.



# 有機合成と酵素

## § 2.1 有機合成分野の概要

有機合成では人工的に有機化合物を作り出すことを目的としている。古くから病の治療として天然の有機化合物が用いられており、薬として有効な成分のみを取り出すことが近年行われてきた。1805年、F.W.A.Serürner がアヘンから強い麻酔作用を持つ morphine を取り出すに成功したことを皮切りに、薬効成分が次々に抽出されるようになっていき、その発展とともに有機化学が発展していった [6]。また、1828年には、Wöhler が有機化合物として初となる尿素の合成に成功し、その後 Liebig を筆頭として、有機化合物の扱い方、構造式での構造理解が明確化されていった [7]。

現在まで、比較的簡単に入手できる化合物から、天然に存在する薬の成分などを人工的に生成する全合成によって、様々なものが合成されてきた。計算機が発達するようになると、実験結果で得られた情報がデータベースに蓄積されていき、化学の現象を計算機上で上手く表現することで、高速なデータ処理が可能となった。そして、情報技術によって、データベースからデータを分類・予測する分野としてケモインフォマティクスは現在まで発展してきた。

ケモインフォマティクスの研究分野として以下のものが挙げられている [8]。

1. ケモインフォマティクス情報検索，データベース，グラフ理論，反応設計など
2. マテリアルズ・インフォマティクス構造物性相関など
3. バイオインフォマティクス
4. 計算機科学
5. 理論・計算科学 (量子科学，分子軌道法，分子科学)
6. コンビナトリアルケミストリー
7. 通信・システム (コンピュータネットワーク，並列化，専用機，コンピュータグラフィックスなど)
8. ラボラトリーオートメーション
9. 関連する化学教育・学習システム

化学分野では情報学に適用できる問題が多く存在する。例えば、化合物の構造に着目したとき、原子の部分の頂点、結合部分を辺とみなすことでグラフ理論の問題になる。合成時の反応経路の設計においては、どの化合物から出発し、いかにステップ数やコストなどを抑え、かつ効率的に目的物を生成していくかという最適化問題に帰着できる。例として、化

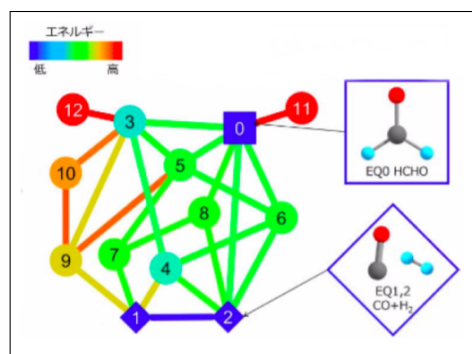


図 2.1: HCHO の化学反応ネットワーク [10]

学反応ネットワーク中の最適な反応経路に関する研究がある [10]. ここでは, 化学反応における安定平衡構造を頂点, 遷移状態構造を辺とした化学反応ネットワークにおける最短経路候補について検討している. 莫大なパターンの経路を調べる代わりに, 反応が経由するそれぞれの状態における最大・最小エネルギー差が小さい経路に絞り, 合成経路設計, 反応予測, 逆合成解析の 3 つの場合において, K 最短経路問題などを適用して, 計算機実験による探索の性能評価を行っている. 図 2.1 は原子 H, C, H, O で構成される化合物の化学反応ネットワークを表している.

化合物を合成する過程においては, 化学合成経路設計と, 化学反応予測の 2 つのアプローチがある [11]. 化学合成経路設計では, 目的とする最終的な生成物を設定し, 何の物質から出発して, どのような合成経路をたどって合成していくか, という逆合成的な手順を用いた手法である. 化学反応予測は, 出発物質を決め, 目的の生成物を得るための反応は実際に起こるのか, 副産物は何が生成されるのかといった手順をたどる. 効率的に合成を進めるため, これらの手法を計算機上で行うためのシステム開発が行われてきた. 化学合成経路を設計するシステムは, 1970 年頃から多数開発されてきた一方で, 化学反応の予測を行うシステムはあまり開発されてこなかった. それは, 化学反応は様々な要因が複雑に絡み合うことで発生するため, 反応の予測が困難であるためである. しかし, 近年では計算機上での化学反応特性の表現方法や機械学習の発展によって, 予測のハードルが下がりつつある. 化学反応時に関わってくる要因を計算機側で上手く表現し, 一部の要因に重点を置いて予測を行うことで, 困難性を解消している.

研究例として, 化学反応時の電子移動に関する, 極性反応とラジカル反応について予測したものがある [12]. ここでは, 1110 個の極性反応, 103 個のラジカル反応からデータベースを作成し, 機械学習の分類を行っている. 分子の反応部位や構造情報, 原子の複数の性質などを, 量子科学計算によって数値化および特徴ベクトルとして表現し, 10 分割交差検証によって, 分類精度を評価している.

このように, 化学反応や化合物における特定の特徴を数値化し, 計算機上で扱いやすく, かつ機械学習に組み込みやすい形式を開発にすることによって, 精度の良い反応予測を可能にしている.



図 2.2:  $\alpha$ -キモトリプシンの立体構造 [14]

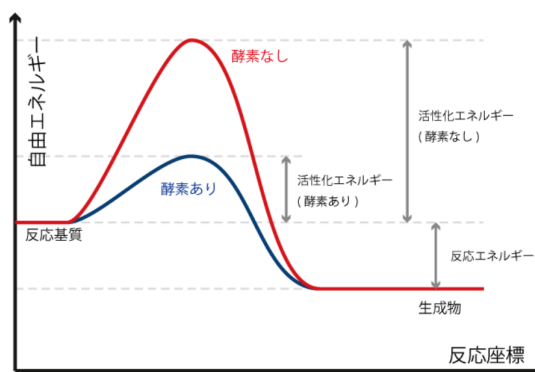


図 2.3: 反応の進行と必要なエネルギー [15]

## § 2.2 酵素と EC 番号

酵素は生体内に必要な化学反応を触媒するタンパク質で、生物が生きていくためには必要不可欠なものである。酵素には基質特異性という、特定の反応物 (基質) のみに触媒反応を示す特性を持っている。これは、Emil.H.Fischer が唱えた「鍵と鍵穴説」と呼ばれる、基質を鍵、酵素を鍵穴とする考え方が用いられている。基質が酵素に結合することで反応が始まり、基質が生成物へと変化すると結合が外れる。このとき、酵素自体は変化することなく元の状態に戻るため、触媒として繰り返し利用できる。酵素の構造イメージを図 2.2 に示す。

より多くの基質と結びついて作用する用途の広い触媒とするために、基質特異性を広げるタンパク質工学と呼ばれる分野がある。ここでは、アミノ酸配列の一部を置き換えることで酵素の性質を改変したり、ランダムに変異させた変異体ライブラリを作成し、スクリーニングによって所望の触媒機能をもったものを選択するといったことが行われている [13]。

酵素を用いることのメリットとして以下のことが挙げられる。

### 反応速度の増加

酵素は生体触媒として、基質の化学反応をより早く、安定して進めることができる。化学反応において、反応が進むにつれてエネルギーが増加し、遷移状態をピークに減少していく。この反応開始から遷移状態になるために、必要なエネルギーを活性化エネルギーと呼び、大きいほど反応が進みにくくなる。しかし、酵素を用いることで必要な活性化エネルギーが低下し、反応を速く進めることができる。酵素を用いた場合と用いなかった場合のエネルギー遷移の様子を図 2.3 に示す。

### 環境への影響の低減

通常、化学触媒は高温や高圧といった条件下で使用する事が適している場合が多い。一方、生体触媒は常温、常圧で使用する事ができ、これらの条件下での反応であれば、高温・高圧にするためにかかるエネルギーの削減につながる。

### 高選択性

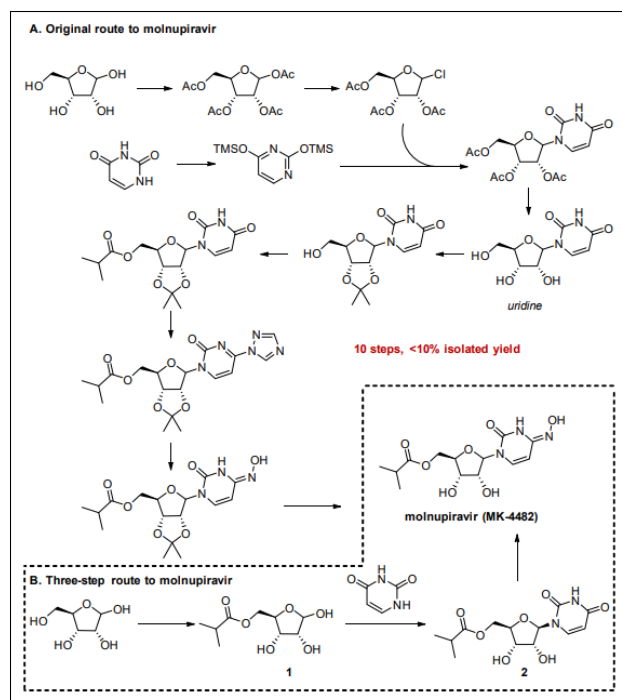


図 2.4: モルヌピラビルの合成

選択性とは、反応が起こりうる化合物の部位が複数ある中で、どれだけ特定の部位のみに反応するかの度合いを表す。選択性が高いほど、特定の構造を持つ部位のみを選んで、反応させることができるため、効率的な合成につながる。

上記の理由から、生体触媒を使って、医薬品を生成する事例が増えている。例として新型コロナウイルスの治療薬として、治験が進められているモルヌピラビル (MK-4482) の合成がある [1]。ここでは従来 10 ステップで行っていたものを、生体触媒を取り入れることによって 3 ステップまで短縮している。図 2.4 にモルヌピラビルの従来における合成方法と、提案された方法の比較を示す。合成ステップを短縮することは、使用する試薬などのコストが減り、結果として環境にも優しい。

## 酵素番号 (Enzyme Commission numbers : EC 番号)

酵素は EC 番号という、4 組の数字の組み合わせからなる番号で管理されており、酵素の性質ごとに分類されている。EC ○.○.○.○ というように番号が振られ、1 桁目の数字 (class) はどの反応を触媒するかによって、1(酸化還元酵素),2(転移酵素),3(加水分解酵素),4(離脱酵素),5(異性化酵素),6(合成酵素),7(輸送酵素) の 7 つに分類されている [16] [17]。2 桁目の数字では、どの結合に作用するか、3 桁目の数字ではどの基質 (化合物) に反応するかや、必要とする補酵素情報というように分類され、4 桁目の数字で 1 から 3 桁目までの組み合わせ番号 (EC ○.○.○) に属する酵素の名前 (登録順) を表している。図 2.5 に EC 番号分類のイメージを示す。EC3(加水分解酵素) を例に見ると、エステル結合に作用する EC 3.1, グリコシド結合に作用する EC 3.2,・・・と分類されている [2]。さらに、EC 3.1 の下層に注目すると、カルボン酸エステルに作用する EC 3.1.1, チオエステルに作用する EC 3.1.2,・・・と分か

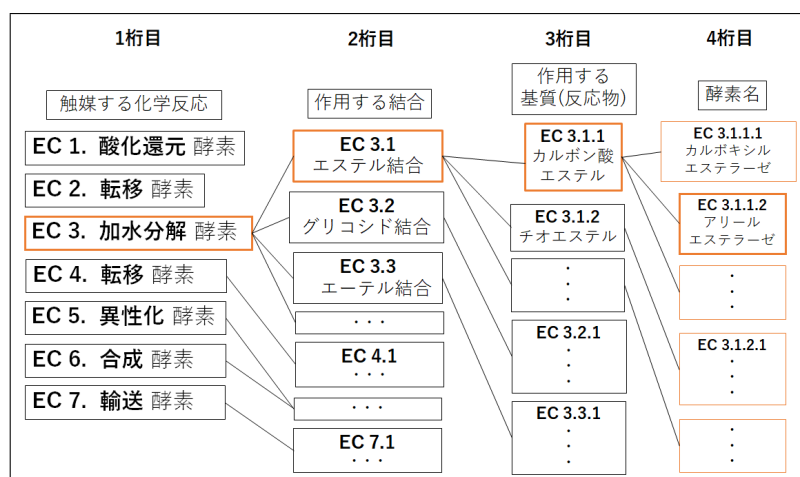


図 2.5: EC 番号による酵素の分類

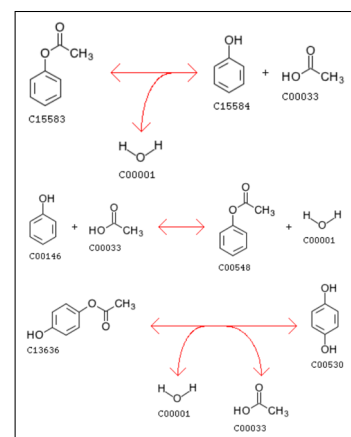


図 2.6: EC3.1.1.2 の代表的な反応式

れ, EC 3.1.1 からは EC 3.1.1.1 のカルボキシルエステラーゼ, EC 3.1.1.2 のアリールエステラーゼ, ... と分類されている. EC 番号を扱うデータベースでは, 4桁の EC 番号に分類されている酵素の名前と触媒する反応式が記載されている. 例えば, Kyoto Encyclopedia of Genes and Genomes (KEGG) [18] では, 図 2.6 のように自然界で起こる酵素を用いた代表的な化学反応 (酵素反応) が 1 つまたは複数登録されている. このように, 既存の酵素反応データを機械学習に用いることで, EC 番号が分かっている酵素反応を入力した際, 正しい EC 番号を予測できるかの推定を行うことができる.

## § 2.3 酵素の探索

ここでは, 代謝経路の解析と有機合成における酵素探索の概要に触れ, 機械学習を用いて EC 番号を予測する背景を述べる.

### 代謝経路解析における EC 番号予測手法の利用)

酵素が触媒することで起こる生体内の反応が別の反応を次々に引き起こす流れを代謝経路と呼ぶ. 代謝経路内で起こる酵素反応を実験的に解析し, 性質を特定 (アノテーション) していくことで, 生体機能の解明につながる. また, 酵素がどの化学反応を触媒するのかを解析することで, 材料や新薬の開発につながる. 遺伝子配列の解析技術の発展によって, タンパク質配列を特定する速度が向上している一方で, 実験等によるアノテーションが追いついておらず, 機能が未解明の配列が蓄積されている. そのため, 機械的な機能予測や自動的なアノテーション法を開発することが重視されている [3]. 従来行われてきた割り当て方法として, アノテーションを必要とするタンパク質配列を BLAST などで検索する方法がある.

例えば, 類似する配列は類似の性質を有するという類似性の概念によって, 類似度の高い既存のタンパク質配列の情報を元に性質が特定されていた. しかし, 検索結果に表示される配列が類似度の低いもののみの場合, 信頼性のあるアノテーションが行えないという問題がある [9]. 機械学習による予測はそのような問題を解決するために用いられ, 類似性



検索の改善や類似性を超えた自然の法則をデータによって推定することで、より高性能な割り当てを可能にすることが期待されている。具体的には、EC 番号によって分類されている酵素のタンパク質配列を数値化し、ベクトルを学習器に入力することで、EC 番号の分類体系に沿った分類モデルを作成できる。

## 有機合成における酵素探索

2.1 節で述べたように有機合成では逆合成的なアプローチで化学合成を行うことがある。最終的に得たい化合物を設定し、切断が容易な部位を探し、2つの化合物に分解する。同様に分解した化合物からさらに切断部位を設定する手順を繰り返し、低コストで入手しやすい化合物が現れるまで分解する。このように合成経路と用いる化合物を決定した後、化合物を合成していく。スムーズに目的の化合物を得るため、様々な試薬や溶媒、触媒、温度、気圧を設定し、最も収率の良い組み合わせを実験によって選択する。このとき、反応を効率良く進める生体触媒として酵素が用いられることがある。一般的に酵素反応は自然界で発生する反応であり、特定の基質のみに作用するため、非天然の化合物には反応しない。そのため、非天然物から目的化合物を作る場合、酵素の働きを活性化させたり、作用する基質の種類を増やすような試薬・溶媒の探索が行われる。具体例として、図 2.4 下部の 1 と示された、ribose と *o*-isobutyrl oxime から 5-isobutyryl ribose を作成する過程を挙げる [1]。初めに 8 種類の酵素製品を用いて最も効率的に 5-isobutyryl ribose が得られる酵素を選択する実験が行われている。ここでは、EC 3.1.1.3 である Novozym 435 が選択された。次に Novozym 435 を用いて 5-isobutyryl ribose の収率が最も高くなる、溶媒や Novozym 435 の濃度が選択され、最終的な合成経路として決定される。

将来的に EC 番号を予測するモデルを有機合成の場で用いることができれば、合成に用いるべき酵素の候補を同じ EC 番号内の酵素製品に絞り込むことができる。それによって、酵素候補を選出する過程を短縮し、次の過程である、酵素で目的物を効率よく得られる実験環境の探索にスムーズに進むことができる。また、酵素に関する知識が十分でない場合でも、実験に用いる酵素候補が自動的に選ばれるため、有機合成をより容易に行えることが期待できる。





# 機械学習による EC 番号予測

## § 3.1 計算機上における化学反応の表現法

化学反応を学習器に入力する場合、化合物の構造などを計算機上において扱いやすい形式で表現する必要がある。化学構造を文字列表現したり、化学反応時の物理的・化学的な指標から得られる数値を用いることによって、化学反応の予測や分類を行う。ここでは、計算機上での各種表現法や用いられるライブラリについて説明する。

### MOL ファイル・SDF ファイル

化合物の構造情報を記したテキスト形式のファイル。「.mol」の拡張子で保存されることが多い。ファイル内には結合している原子と各原子の 3 次元座標リストやどの原子同士が結びついているかのリストが記述されている。通常の構造式と mol ファイルを比較したものを図 3.1 に示す。複数の化合物 MOL ファイルを統合したものは、拡張子「.sdf」からなる SDF ファイルとなる。2 つ以上の分子の MOL ファイルをデータベースから同時に入手する際は、SDF ファイルとなることが多い。

### SMILES

化合物の構造を文字列で表したもの。以下の規則に従い、化学構造を文字列に変換していく [20]。

1. 原子は元素記号で表し、2 文字で区別が付きにくい原子 (Nb と NB 等) は [ ] で囲む
2. 水素原子は省略する
3. 隣接する原子は隣に記す
4. 二重結合は =, 三重結合は # で表し、単結合・芳香族結合は省略する (芳香族原子は小文字の c など で表記する)
5. イオンなどで結合がない部分は「.」で分ける
6. 構造が分岐する箇所は ( ) で表記する
7. 環構造は切断して切断箇所を記すとともに (C1 など), 鎖錠構造で表す。

これらに加えて、さらに以下の規則を加えた isomeric SMILES を本研究では用いる。

8. 同位体 (例えば炭素) がある場合 [13C] という表記にする
9. 立体異性体を区別するための絶対配置を「@」または「@@」で表現する
10. 二重結合などで生じる幾何異性を「/」と「\」で表す



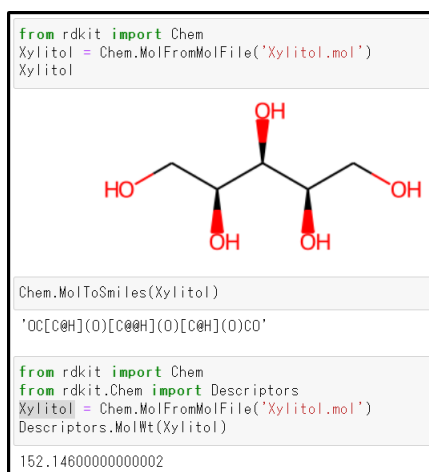


図 3.2: rdkit を用いた化合物の情報

機械学習予測予測 [24] などに用いられている。物理・化学的特性値として、化合物の分子量や油 (水) への溶けやすさを表す MolLogP [27], 分子の静電的な相互作用を表す 14 種の PEOE\_VSA [28], 分子断片の有無に関する指標として、芳香族アミンの数を表す fr\_Ar\_NH, チオール基の数を表す fr\_SH などがある。例えば、化合物の分子量を表す記述子 MolWt を知りたい場合, Descriptors クラスにある MolWt メソッドに生成した構造式オブジェクトを渡すことで, MolWt が計算され出力される。図 3.2 に, rdkit を用いて化合物の構造式と SMILES を出力した様子, および化合物の MolWt を計算した結果を示す。

## § 3.2 EC 番号予測手法

EC 番号は代謝経路の解析や酵素反応の探索において, おおまかな酵素の性質や探索範囲を特定し, 実験コストや時間を短縮するために用いられてきた。機械学習による自動的な EC 番号の予測が行われており, データベース上にある既存のタンパク質配列や酵素反応データを用いて予測モデルを構築し, 精度をより高める予測手法を開発することで, 将来的に未知データに対して最適な EC 番号を提示することが目指されている。既に EC 番号が割り当てられているデータを用い, EC 番号の 1 桁から 3 桁または 4 桁目までのクラスに対する多クラス分類を行うことで, データが正しい EC 番号のクラスに予測されるかを評価する。EC 番号の予測手法として, タンパク質配列に着目した手法や化合物の物理・化学的特性値, 構造的特徴を用いた手法があり, 以下でそれぞれの手法を述べる。

### タンパク質配列を用いた予測

タンパク質配列は文字列で構成されているが, 2 次元上に配列をマッピングし画像の機械学習で用いられる手法などを用いることで, タンパク質配列を正しい EC 番号に割り振ることを目指す。EnzymeNet [3] では, 酵素機能を有する配列約 1,000,000 と有していない配列約 140,000 を用い, EC7 を除く EC 番号の 4 桁全ての予測を行っている。ここではディープラーニングによって文字列を変換し, 配列の長さを 1024 に調整後, 分類器の中で 1024 × 1024 の配列特徴を表すマップを出力している。Convolutional Neural Network (CNN) を用

```

import pubchempy
pubchempy.get_properties([ 'MolecularFormula', 'MolecularWeight',
                           'IsomericSmiles' ], 'glucose',
                           'name', as_dataframe=True)

```

	MolecularFormula	MolecularWeight	IsomericSMILES
CID			
5793	C6H12O6	180.16	C([C@@H]1[C@H]([C@@H]([C@H]([C@H](C(O1)O)O)O)O)O)O

図 3.3: PubChemPy でグルコースの情報を取得した結果

いて配列のマップを深層学習し、配列が酵素機能を有する (非酵素) か否か、および EC 番号 1 桁目クラス (EC1~EC6) に対してどのクラスに属するかのクラス分類が行われた。さらに、転移学習を用いて、1 桁目が予測された配列の 2~4 桁目を予測する手法を開発し、Macro F1 スコア 0.85 の従来の配列ベースの予測で最も高い精度を出力した。

### 化合物の物理・化学的特性値を用いた予測

MOlecular Mapping of Atom-level Properties (MOLMAP) 反応記述子 [5] を用いた手法では、68 種類の構造・物理・化学的指標を用い、EC 番号が判明している化学反応の分類が行われた。初めに、68 種類の化合物の特性値情報を Self-Organizing Map (SOM) [34] によって 2 次元平面上に凝縮マッピングし、 $7 \times 7$  から  $29 \times 29$  個程度のセルで構成される数値の MOLMAP を生成された。次に基質と生成物の MOLMAP の差を取ることで酵素反応が生じた際の化合物の化学変化を表現し、約 7,000 のデータに対する RF の 1~3 桁目分類において、それぞれ予測精度 0.92, 0.85, 0.83 を得ている。

### フィンガープリントの差分を用いた予測

反応物 (基質と試薬) から生成物に変化する際のフィンガープリントの変化に着目した手法として Differential Reaction Fingerprint (DRFP) [4] がある。具体的な実装手順を以下に示す。

1. 化合物内の各分子に着目し、分子単体、分子と 1~2 つ隣に隣接する分子が結び付いた部分構造の SMILES を抽出する。
2. 抽出した全ての部分構造を反応物の集合と生成物の集合に分割し、対称差と取ることで、2 つの集合に重複する部分構造を除外する。
3. 各部分構造を ECFP に基づいて 32 ビットの整数に変換し、10240 で割った余りをフラグ立てるビットとみなし、10240 ビットのフィンガープリントを作成する。

多層パーセプトロンを用いて EC 番号の 3 桁 (EC 7 を除く) まで予測し、F1-Score として Rhea の反応データのみで  $0.87 \pm 0.02$ 、BRENDA、PathBank、MetaNetX を含めたデータでは  $0.77 \pm 0.01$  を得ている。

## § 3.3 機械学習と特徴量エンジニアリング

ここでは、本研究で用いる機械学習手法について述べる。

### 多クラス分類による機械学習予測

2 種類のいずれかの正解ラベルが割り振られたデータを分類する 2 クラス分類に対して、3 種類以上の正解ラベルのいずれかが付与されたデータを分類する問題は多クラス分類と呼ばれる。多クラス分類では、学習データを分類器に入力することでデータの分類体系を近似した分類モデルを構築し、テストデータを入力することで分類モデルがテストデータを正しく分類できるかを評価する。多クラス分類を評価する指標を以下に示す。なお簡便化のため、2 クラス分類を元に説明する。

**Precision, Recall および Accuracy**

$D$  個のデータが2つのクラス正, 負のいずれかに割り当てられているとする. また, 分類モデルが正に分類されるデータを正しく正と予測した数を True Positive (TP), 負に分類されるデータを誤って正と予測した数を False Positive (FP) とする [1]. さらに, 負に属するデータを負と予測した数を True Negative (TN), 正のデータを誤って負と予測した数を False Negative (FN) とする. このとき, 分類器が正に分類したデータのうち, 実際に正に所属するデータであった割合を Precision, 正に所属するデータに対して分類器が正しく正に分類したデータの割合を Recall と呼ばれ, 以下の式で表される.

$$Precision = \frac{TP}{TP + FP} \quad (3.1)$$

$$Recall = \frac{TP}{TP + FN} \quad (3.2)$$

Precision は分類モデルの正確性を測る指標に対して, Recall はデータを正しいクラスであると判断する感度を表す指標であり, トレードオフの関係となっている [26]. また, 全体の予測精度を表す指標である Accuracy は以下の式で表される.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (3.3)$$

$$(3.4)$$

多クラス分類の場合, 1つのクラスに着目し, 着目したクラスを正, それ以外のクラスを負のクラスとみなして Precision, Recall が計算される. そして, 各クラスの平均を取った, Macro Average Precision, Macro Average Recall が用いられる.

### F1-Score

Precision と Recall のバランスを測る指標で, それぞれの調和平均を取った指標となる. F1-Score (F1) を表す式は以下ようになる.

$$F1 = \frac{2Precision \cdot Recall}{Precision + Recall} \quad (3.5)$$

$$(3.6)$$

多クラス分類では, Macro Average Precision, Macro Average Recall の調和平均である, Macro F1-Score が用いられる. クラス間でデータに偏りがある場合, Accuracy ではデータ数の多い多数クラスに影響され, 少数クラスに対する分類精度が低い場合でも, 多数クラスの分類精度が高ければ, 精度が高い分類モデルと評価される. 一方で, Macro F1-Score は各クラスに対する分類精度を評価し, 平均を取っているため, 各クラスの重要度が等しい場合の評価に適している.

### オーバーサンプリング

各クラスのデータ数に偏りが生じているデータを不均衡データと呼ぶ. 多クラス分類において, 少数クラスデータが誤って多数クラスに分類される可能性が高い, もしくは多数ク

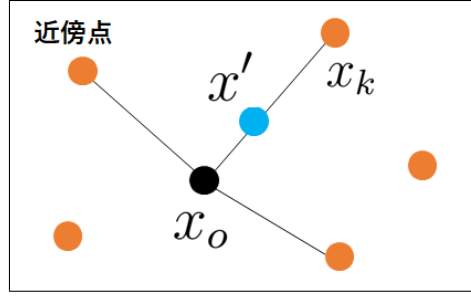


図 3.4: SMOTE によるオーバーサンプリングのイメージ

ラスのデータが少数クラスに誤分類される可能性が低くなるのは必然的である．その対処法として，少数クラスのデータを仮想的に増やすオーバーサンプリングや多数クラスのデータを減らすアンダーサンプリングなどの手法がある．ここでは，Synthetic Minority Over-sampling Technique (SMOTE) [29] を用いたオーバーサンプリングについて述べる．SMOTE はデータが多次元の特徴量を持った特徴空間において，少数クラスに属するデータ間で線分を取り，線分上に同じ少数クラスのデータを新たに生成することでオーバーサンプリングを行う手法である．具体的なアルゴリズムは次のようになる．

1. ある少数クラスの 1 つのデータに着目し，図 3.4 のようにデータの近傍にある  $K$  個の同じ少数クラスのデータを抽出する．
2.  $K$  個の近傍データから 1 つを選択し，以下の式のように，データ点  $x_o$  と選択された近傍データ点  $x_k$  の間に新たなデータ点  $x'$  が生成される．  
ただし， $r$  はランダムに生成される  $[0, 1]$  の値である．

$$x' = x_o + r(x_k - x_o) \quad (3.7)$$

### Random Forests (RF) による多クラス分類

Random Forests (RF) は複数の決定木を用いる機械学習モデルであり，モデルの表現力は高いが，過学習に陥りやすい決定木を組み合わせることで，汎化性能を高めることができる．決定木のイメージを図 3.5 に示す．決定木ではまず根のノードに特徴量を割り当て，特徴量の閾値に応じてデータを下の 2 つのノード内に分割する．それ以降は，根の下にある各ノードに対しても特徴量を割り当て，あるノードに到達したデータを特徴量の閾値に応じて，下位ノードに 2 分割する工程を繰り返していく．

ノードを分割する際の分割基準と用いる特徴量は，以下で定義される情報利得  $IG$  によって決まる [36]．

$$IG(D_P, f) = I(D_P) - \frac{N_{left}}{N_P} I(D_{left}) - \frac{N_{right}}{N_P} I(D_{right}) \quad (3.8)$$

$f$  は分割時に用いられる特徴量， $D_P$  は上位ノード内のデータ， $D_{left}$ ， $D_{right}$  はそれぞれ分割先の下位ノード内のデータを表し， $N_P$ ，および  $N_{left}$ ， $N_{right}$  は上位ノード，下位ノード

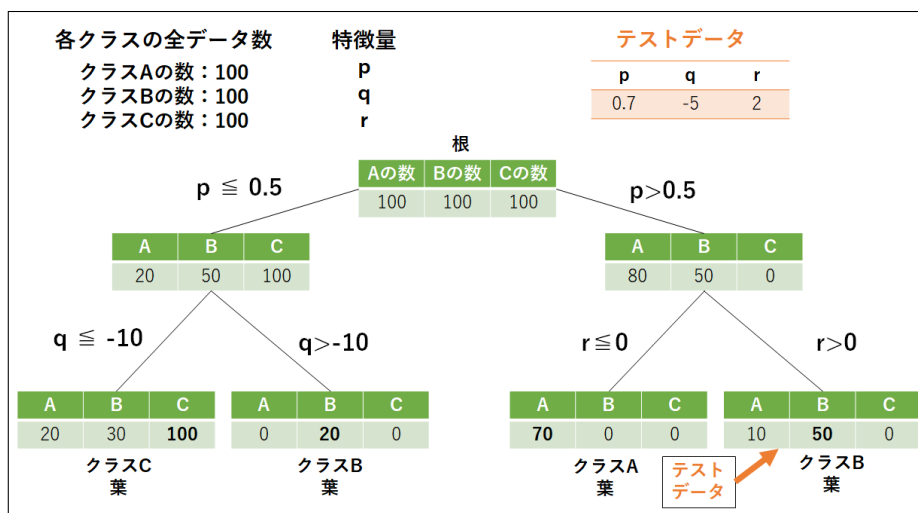


図 3.5: 決定木におけるクラス分類の様子

のデータ数を表す。  $I$  は不純度を表し、ジニ不純度、エントロピーなどが用いられる。ジニ不純度を用いた場合、ノード  $t$  に対するジニ不純度  $I_G(t)$  は以下ようになる。

$$I_G(t) = \sum_{i=1}^c p(i|t)(1 - p(i|t)) = 1 - \sum_{i=1}^c p(i|t)^2 \quad (3.9)$$

$p(i|t)$  はノード  $t$  におけるクラス  $i$  のデータの割合、 $c$  はノード  $t$  内のクラス数を表す。各ノードで  $I_G$  が最大となるように、特徴量  $f$  とデータを分割する閾値が決定される。

基本的に全ての葉ノードで  $I_G(t) = 0$ 、すなわちデータの属するクラスが1種類となるまで分割が進行されるが、過学習を避けるため、決定木の最大深さを(分割されるデータの最大の階層数)を設定することで分割を打ち切ることもある。分割が停止したとき、葉に含まれるデータのクラス割合が、最も大きいクラスを葉ノードのクラスとして決定される。テストデータは葉ノードのいずれかに分類され、葉ノードのクラスとして予測される。

RFでそれぞれの決定木を作成する際には、全データから一部のデータを復元抽出し、決定木に入力される。また、特徴量も一部(または全て)選択されて決定木に用いられる。それによって、異なるデータと特徴量からなる決定木モデルを作成し、組み合わせることで汎化性能を高める。

RF ライブラリである Python の `sklearn.ensemble.RandomForestClassifier` [31] ではパラメータとして、決定木の数 `n_estimators` や、各決定木の最大深さ `max_depth`、用いる特徴量の数 `max_features` などを設定する。また、テストデータの予測クラスは、それぞれの決定木で分類された、葉ノードの各クラス確率の推定値に対する平均値が最も高いクラスとなる。例えば、ある1つのテストデータが3つの決定木において、クラス A,B,C の確率がそれぞれ (0.10, 0.40, 0.50), (0.00, 0.10, 0.90), (0.20, 0.50, 0.30) となる葉ノードに分類された場合、平均値は (0.10, 0.33, 0.57) となり、最終的にクラス C に分類される。

### ラッパー法を用いた記述子選択

必要以上に記述子を用いることは、基質から生成物の化学変化を分類する際の汎化性能低下につながるため、適切な種類の記述子のみ限定することが望ましい。ラッパー法 (Wrapper



Method) では分類・回帰モデルの予測精度を評価し、最も評価の高い特徴の組み合わせを選択するため、今回の記述子選択に適している。

ラッパー法には、指定した特徴数まで特徴を 1 つずつ追加し、追加するたびに最も評価の高い組み合わせを選択する Step Forward 法、全ての特徴を選択した状態で、評価が最も高い組み合わせとなるよう、指定した特徴数まで特徴を 1 つずつ削減していく Step Backwards 法、全ての組み合わせを探索し、最高評価の組み合わせを選択する Exhaustive 法がある。今回は、Step Forward 法を適用し、Python の mlxtend ライブラリ内にある、SequentialFeatureSelector に実装されている、Sequential Forward Selection(SFS) を用いる [35]。特徴選択の手順は以下のようになる。

1.  $n$  個の記述子から 1 つ選択し、 $n$  種類の分類モデルを作成
2. 最も評価の高いモデルに用いられている、記述子を選択する。
3.  $n - 1$  個の記述子から 1 つ選択し、先ほど選択されたモデルに追加することで、新たな分類モデルを作成する。
4.  $n - 1$  個のモデルで最も評価の高いものに用いられている、記述子の組み合わせを選択する。
5. 指定した特徴数になるまで 3 と 4 を繰り返す。

モデルの評価基準としては、各 EC 番号クラス分類時の F1 スコア平均を用いる。



## 提案手法

### § 4.1 特徴ベクトルと特徴エンジニアリングによる EC 番号予測

有機合成では、反応の効率性や環境などの観点で目的の化合物を生成する際に酵素を生体触媒として用いられる機会が増加している。酵素には EC 番号が割り当てられており、酵素の性質の特定や触媒する化学反応の探索効率化のため、機械学習で入力された酵素に対する EC 番号を予測する研究が行われてきた。手法としてタンパク質配列や化合物の構造的特徴、物理・化学的特性値などが用いられてきたが、これらの EC 番号予測手法を有機合成に適用することで、探索する酵素の種類に対して絞り込みができ、実験によって最適な酵素を探索する時間やコストの削減が期待できる。

本研究では、有機合成における酵素探索に焦点を当て、化学反応に用いるべき最適な酵素の候補を EC 番号として予測する機械学習モデルを開発する。提案手法として、従来化学反応の探索や予測の場面で多く用いられてきた 210 種類の RDKit 記述子を用い、Random Forest (RF) に特徴選択とオーバーサンプリングを組み合わせた予測モデルを作成する。有機合成という観点から見れば、従来手法の物理・化学特性値やフィンガープリントを用いた予測手法が容易に適用可能であり、最近の研究ではフィンガープリントの差分を用いた手法で優れた予測精度が得られている。一方で、フィンガープリントは主に化合物の構造や化学反応時の構造変化を捉える指標であるため、物理・化学的な情報を加えることで、酵素を用いた化学反応の特徴をより詳細に捉えることが可能になると考えられる。しかし、フィンガープリントは次元が 1000 や 2000 を超えるものがあり、物理・化学的特性値と組み合わせた場合、学習コストが増大する懸念がある。そこで、RDKit で利用可能な MACCS Keys に類似する 85 種類の分子断片のバイナリ値を用いることで、学習コストを抑えつつ、各 EC 番号における化学反応の特徴をより説明できるようなモデルの作成が期待できる。

#### 特徴ベクトルの作成

RF に入力するデータとして、基質から生成物に変化する際の RDKit 記述子の変化を用いる。EC 番号が割り当てられた酵素反応式において、各化合物の 210 種類の物理・化学的特性値を計算し、基質の特性値の和と生成物の特性値の和に関する差を取った特性値変化量を計算する。例えば、図 4.1 において、基質 A, B と生成物 C, D からなる反応式を考えた場合、A, B, C, D はそれぞれ 210 次元の特性値を持ち、 $(A + B) - (C + D)$  を計算することで 1 つの反応式が 210 次元の特性値変化量を持つことになり、これを特徴ベクトルと定義する。最終的に表 4.1 のような行に EC 番号、列に記述名が記載された反応データを用いる。

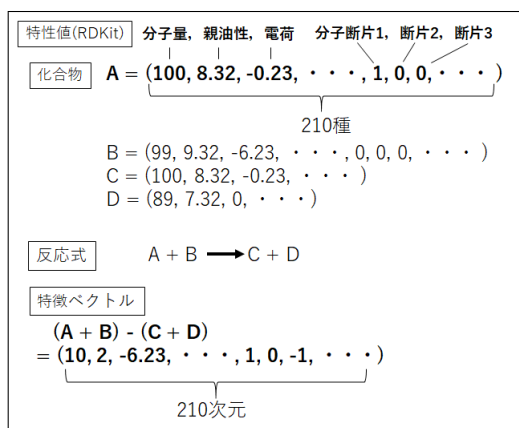


図 4.1: 特徴ベクトルのイメージ

EC	MaxAbsEState Index	MinAbsEState Index	MinEState Index	qed	...
4.1.1.74	-7.449074	-0.064815	-2.62963	-0.360209	...
1.2.1.8	-0.197403	0.405116	1.30787	-0.079826	...
2.5.1.85	0.593569	-2.312624	0.196239	0.488055	...
1.4.1.4	0.234718	0.418052	-0.413194	0.389325	...
1.1.1.3	-0.15593	0.059255	-0.317778	0.016389	...
...	...	...	...	...	...
4.4.1.13	-7.622896	-1.049639	0.2744	-0.411722	...
4.4.1.13	-7.478685	-1.049639	0.25852	-0.440793	...
2.3.1.-	-0.286151	0.454936	0.039395	-0.386083	...
2.3.1.-	-0.306521	0.11352	0.039395	-0.375732	...
2.3.1.-	-0.254718	0.523159	0.028495	-0.386848	...

表 4.1: 実際の特徴ベクトル

## EC 番号予測における記述子選択

RDKit 記述子を全て用いて RF のクラス分類を行った場合, RF は各ノードで 1 つの記述子の閾値によって特徴ベクトル进行分类するため, ある EC 番号の特徴ベクトルを説明するのに不要な記述子やデータに同様の特徴を付与する記述子が複数存在することで, 異なる EC 番号の酵素反応の区別があいまいになる可能性がある. そのため, 記述子選択によって必要な記述子のみ进行分类に用いる. 本研究では Python の `mlxtend.feature.selection` にある `SequentialFeatureSelector` を適用し, Step Forward 法を用いる. また, 記述子数が 30 になるまで記述子選択を継続し, 評価が最も高い記述子の組合せを用いる. 最終的に 210 種類から RF の分類精度が最も高くなる記述子 10~30 程度の組み合わせを用いる.

## オーバーサンプリングによる少数クラスのデータ増大

EC 番号が割り当てられた酵素反応式もクラスごとにばらつきがあり, 不均衡データである. 背景として酵素の汎用性と機械学習の対象となるデータの形式が関わっている. 例えば, EC3.1.1.3 に分類されるリパーゼは様々な反応に作用するため, 産業分野での汎用性が高く, 多くの酵素製品が生産されている [37]. そのため, 提供されている酵素反応のデータが多い. 一方で, その他の EC 番号は酵素反応データが 5 種以下のクラスなども多くあり, 新たな酵素に EC 番号を割り当てる作業にも時間を要する. よって, 汎用性が高い酵素の EC 番号は急増するが, 応用機会が少ないものは年々ほとんど増加しないため, 多数クラスと少数クラスのデータ数の隔たりは拡大していく一方となる.

本研究ではオーバーサンプリング手法の 1 種である SMOTE を適用する. 多数クラスのデータ数を減らすアンダーサンプリングも考えられるが, 多数クラスの重要なデータを削除してしまう可能性が挙げられており [30], 多数クラスを特徴づける説明力が減少する可能性を考慮して今回は用いないこととする. 代わりに少数クラスのデータを仮想的に増やすことで, 少数クラスに対する分類精度の向上を図る. 数値実験では Python の `imblearn.over_sampling` [32] に搭載されている SMOTE を用い, 近傍データ数  $K(k\_neighbors)$  を  $K = 3$  に設定した. また, デフォルトでは少数クラスのデータ数を最多クラスのデータ数に合わせてオーバーサンプリングされるが, 最少クラスに対する最多クラスの比率が 2000 を超えるデータを用いるため, 学習時間が元のデータ数の場合と比べて膨大になることが懸念される. そのため, オーバーサンプルデータ数 (`sampling_strategy`) を合計データ数に合わせて調整するものとする.

rxn	ec	source	rxn	ec	source
<chem>CC(=O)C(=O)[O-]</chem> [4.1.174>> CC=O.O=C=O	4.1.174	brenda_ reaction_smiles	<chem>N[C@@H](CCC(=O)O)C(=O)O.O=C(O)C(=O)Cc1ccccc1</chem> [2.6.157>> N[C@@H](Cc1ccccc1)C(=O)O.O=C(O)CCC(=O)C(=O)O	2.6.157	pathbank_ reaction_smiles
<chem>NC(=O)c1ccc[n+](c1)[C@@H]2O[C@H](COP(=O)(O)OP(=O)(O)OC[C@H]3O[C@H](n4cnc5c(N)ncnc54)[C@H](O)[C@@H]3O)[C@H](O)[C@H]2O)c1.NCCC=O.O</chem> [1.2.18>> NC(=O)C1=CN[C@@H]2O[C@H](COP(=O)(O)OP(=O)(O)OC[C@H]3O[C@H](n4cnc5c(N)ncnc54)[C@H](O)[C@@H]3O)[C@H](O)[C@H]2O)C=CC1.NCCC(=O)O.[H+]	1.2.18	brenda_ reaction_smiles	<chem>O=CC(O)COP(=O)(O)O</chem> [5.3.1.1>> O=C(CO)COP(=O)(O)O	5.3.1.1	pathbank_ reaction_smiles
...	...	...	...	...	...
<chem>NC(=O)CC[C@H](NH3+)[C-](=O)[O-]</chem> Nc1ncnc2c1ncn2[C@@H]1O[C@H](COP(=O)(O)[O-])OP(=O)([O-])[O-][C@H](O)[C@H]1O.O=P([O-])[O-] [6.3.1.2>> Nc1ncnc2c1ncn2[C@@H]1O[C@H](COP(=O)(O)[O-])OP(=O)([O-])OP(=O)([O-])[O-] [C@H](O)[C@H]1O.[NH3+][C@H](CCC(=O)[O-])C(=O)[O-].[NH4+]	6.3.1.2	metanetx_ reaction_smiles	<chem>C[NH+][1CCCC[C@H]1c1ccc(O)nc1.O.O=O</chem> [1.5.3.6>> C[NH2+][CCCC(=O)c1ccc(O)nc1.OO	1.5.3.6	rhea_ reaction_smiles
<chem>NC(=O)CC[C@H](NH3+)[C-](=O)[O-]</chem> [1.4.1.13>> [NH3+][C@H](CCC(=O)[O-])C(=O)[O-].[NH4+]	1.4.1.13	metanetx_ reaction_smiles	<chem>COc1cc(C=C(C(=O)O)CC[N+](C)(C)C)cc(OC)c1O.O</chem> [3.1.1.49>> COc1cc(C=C(C(=O)O-))cc(OC)c1O.[N+](C)(C)CCO.[H+]	3.1.1.49	rhea_ reaction_smiles
...	...	...	...	...	...

図 4.2: SMILE 形式のデータセット (一部抜粋および改変) []

## § 4.2 EC 番号予測モデルの構築と予測

本研究では、酵素反応ベースで EC 番号予測を行った従来手法 [5] [4] と同様に EC 番号の 1 桁目から 3 桁目を予測するモデルとする。データセットとして Rhea, BRENDA, PathBank, MetaNetX の 4 つのデータベースからなる、図 4.2 のような SMILES 形式で記載された酵素反応 [4] を用いる。初めにデータセットを反応式の左辺と右辺に分解し、各化合物を RDKit の構造式オブジェクトに変換する。次に、210 種類の記述子を用いて各化合物の特性値変化量を計算し、EC 番号と 210 次元の酵素反応式の特徴ベクトルを取得する。次に、適切な EC 番号が割り当てられていないものや SMILE から特性値変化量に変換した際に無効値となる要素を含む特徴ベクトルを削除した。また、記述子間とデータ間で相関係数が 1 のものや記述子選択に用いることが不可能な大きな値を要素に持つ特徴ベクトルを削除した。さらに、5 分割交差検証を実施するため、特徴ベクトルが 5 以下のクラスを除外した。最終的に、47090×200 記述子のデータを用いて予測モデルの構築と予測を行う。

データ数 47090 のうち 80% を学習データ、20% をテストデータとし、各クラス均等な割合で行きわたるように分割する。EC 番号予測は記述子選択と学習、予測の 3 段階で構成される。なお、それぞれの段階で RF を用い、決定木における各ノードの分割と用いられる記述子の基準としてジニ不純度を用いる。また、決定木の各ノード分割において、RF で指定する記述子数のみでは適切に特徴ベクトルを分割できない可能性を考慮し、記述子選択された全ての記述子を用いる。実装は Python の `sklearn.ensemble.RandomForestClassifier` を用い、記述子選択の段階ではデフォルトのパラメータ、モデル作成と予測には設定されたパラメータを用いる。図 4.3 にモデル構築と予測の流れを示す。

記述子選択では学習データに対して RF を用い、EC 番号 1 桁目 (EC 1, EC 2, EC 3, EC 4, EC 5, EC 6) の多クラス分類と EC 番号 1 桁目に属する 2,3 桁目 (EC 1.X.X, EC 2.X.X, EC 3.X.X, EC 4.X.X, EC 5.X.X, EC 6.X.X) クラスの多クラス分類に分けて行う。なお EC 7 クラスは特徴ベクトル数が少ないため、本研究の予測対象外とする。このように分割した理由として、2 つ挙げられる。1 つ目は学習に要する時間が関係する。数値実験ではクラス数 148, データ数 37,672 の学習データを用いるため、学習に膨大な時間を要するためである。2 つ目は、EC 番号の分類法によるものである。EC 番号 1 桁目の分類で酵素の性質が

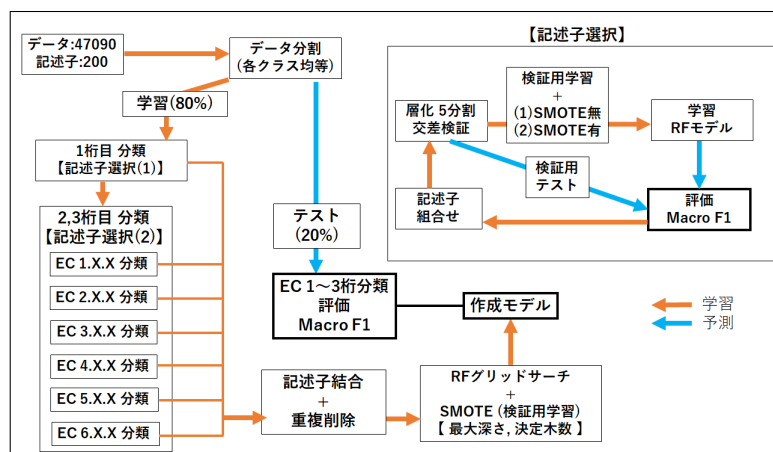


図 4.3: モデル構築と予測の流れ

大きく決まり、2桁目3桁目と進むにつれて、性質が似通った酵素が同じクラスに属する構造となっているためである。そのため、1桁目を分類する記述子と2,3桁目分類する記述子を分割して選択することで、分類粒度が異なる記述子を適切に選択できると考えられる。

7回の記述子選択それぞれに5分割交差検証を実施し、学習データをさらに検証用学習データと検証用テストデータに分割する。そして、各分割の検証用学習データに対してSMOTEを適用する。なお、1桁目分類ではSMOTEを適用せず、素のデータで記述子選択を行っている。表4.2に各クラス分類における学習データの内訳とSMOTEで増加させるデータ数を示す。少数クラスは表内のSMOTE増加数以下のクラスとし、この閾値以下のクラス数を閾値の数になるようにオーバーサンプリングを行う。EC 2.X.Xは最もデータが多く、記述子選択の学習コストを考慮して閾値をデータ数合計の2%とした。EC 1.X.XはEC 2.X.Xの1/4程度のデータだが、クラス数が多いため3%とした。また、EC 3.X.XとEC 4.X.Xは合計の減少に応じて割合を増やし、EC 5.X.XとEC 6.X.Xは最多クラス数に合わせてオーバーサンプリングした。実際は、検証用学習データにSMOTEが適用され、学習データの分割率が4:1の場合を想定して、SMOTE増加数 $\times 0.8$ としたものを閾値に用いている。

記述子の組合せに対して検証用テストデータを適用し、Macro F1-Scoreで分類精度を評価する。新しい記述子が1つ追加されたときの組合せで、5分割分の平均値が最も高いものを選択し、次の記述子選択を行う。30種までの選択で、平均スコアの最高値が4回更新されなくなる直前における記述子の組み合わせを、そのクラス分類で選ばれた記述子とする。最終的に7回分の記述子集合をマージし、重複を削除したものを学習に用いる。

学習では選択された記述子数の次元からなる学習データに対して、RFのパラメータ選択を行う。ここでは、EC番号1~3桁目の多クラス分類を実施する。方法としてグリッドサーチ(`sklearn.model_selection.GridSearchCV` [33])を適用し、5分割交差検証を実施する。記述子選択と同様に各分割の検証用学習データにSMOTEを適用し、Macro F1-Scoreで最適なパラメータの組み合わせを選択する。このとき、オーバーサンプリングの閾値は(最少クラスのデータ数) $\times 100$ とし、調整するパラメータとしては、決定木の数(`n_estimators`)と最大深さ(`max_depth`)とする。

予測ではグリッドサーチで得られたパラメータと選択された記述子で作成されたモデルに対して、データ数9,418のテストデータを入力し、1~3桁目の148クラスに対する多クラ



表 4.2: 学習データの内訳と SMOTE 増加数

EC 1.X.X計	クラス数	SMOTE 増加数	1.1.1	1.14.13	1.2.1	1.14.14	1.3.1	1.13.11	1.14.11	1.4.3	...	1.23.5
6380	64	3%(191)	1745	761	666	408	348	242	173	156	...	5
EC 2.X.X計	クラス数	SMOTE 増加数	2.7.8	2.3.1	2.1.1	2.4.1	2.7.1	2.7.7	2.5.1	2.6.1	...	2.7.3
23160	24	2%(463)	10074	7309	2797	686	602	426	303	280	...	5
EC 3.X.X計	クラス数	SMOTE 増加数	3.1.1	3.1.3	3.6.3	3.2.1	3.5.1	3.6.1	3.1.2	3.1.4	...	3.3.1
5377	27	10%(538)	2277	589	508	448	347	237	169	104	...	5
EC 4.X.X計	クラス数	SMOTE 増加数	4.1.1	4.2.1	4.1.2	4.4.1	4.2.3	4.1.99	4.3.1	4.1.3	...	4.6.1
1878	14	20%(376)	1037	361	106	106	101	44	39	32	...	5
EC 5.X.X計	クラス数	SMOTE 増加数	5.5.1	5.3.1	5.3.3	5.4.99	5.3.2	5.1.3	5.4.2	5.4.3	...	5.1.1
273	12	最多(80)	80	46	44	20	15	13	13	12	...	5
EC 6.X.X計	クラス数	SMOTE 増加数	6.2.1	6.3.2	6.3.4	6.3.5	6.3.1	6.4.1	6.1.2			
604	7	最多(266)	266	233	30	29	22	18	6			

ス分類を行い、予測精度を評価する。評価指標として各クラスに対する Precision, Recall, Macro F1-Score を用い、全体の平均値と Accuracy を出力する。

## § 4.3 提案手法の実装と流れ

SMILES 形式のデータの加工から特徴ベクトルの作成、機械学習モデルの構築から予測までの流れは、Jupyter Notebook 上で実装を行った。本節では、前半で提案システムの実装、後半で提案システムによる流れについて述べる。

### SMILES 反応式の加工

図 4.2 の SMILES 形式の酵素反応式は rxn の項目に記載されている。中央の EC 番号を含んだ文字列で反応式の右辺左辺が分割されており、“.” で第一項、第二項などを区別されている。初めに、正規表現を用いて右辺と左辺の SMILES に分解した。このとき、MetaNetX のデータに含まれている、EC 2.A.3.1 などのデータを削除した。さらに文字列の “.” をもとにして右辺第一項、第二項、..., 左辺第一項、第二項のように、化合物単体の SMILE に切り分け、分割された SMILES の酵素反応式を取得した。次に SMILES 化合物を変換した構造オブジェクトで特性値を計算する際、無効な値を出力する SMILE 化合物の削除を行う。[Co+] や [Ca+2] などの電荷を含む SMILES の一部は記述子の正確な計算が行えない問題がある。そのため、右辺左辺から重複のない化合物 SMILES リストを作成し、特性値計算時に無効値が出る SMILES の除外リストを出力した。そして、先ほど取得した酵素反応式内で除外リストに該当する SMILES を検索し、右辺左辺いずれかに含まれる、反応式を削除した。除外リスト、および削除後の SMILES 反応式を図 4.4, 4.3 に示す。

### 特徴ベクトルの作成

図 4.3 の右辺と左辺の反応式を構造オブジェクトの反応式にそれぞれ変換し、RDKit 記述子と呼び出すことで、全ての化合物に 210 次元の特性値が付与される。左辺の特性値の和と右辺の和に対する差を取ることで、表 4.1 のような 210 次元の特徴ベクトルが得られる。ここで、全てのデータで値が等しくなる記述子の片方 (10 種類) と、同様に全ての記述子で値が等しい特徴ベクトルの片方を削除した。また、記述子選択で用いることができない float32 を超える要素や全て 0 の要素を持つ特徴ベクトルも削除した。

Removal List
*OP(=O)[(O-)]OC[C@H]1O[C@@H](n2cn c3c(=O)n4c(CCC([NH3+])C(=O)OC)c(C)nc4n(C)c32 )[C@H](O)[C@@H]1OP(=O)[(O-)]O*
O[As](O)c1ccccc1
*OP(=O)[(O-)]OCC(C)(C)C(O)C(=O) NCCC(=O)NCCSC(=O)CC(=O)CCCCCCCCCCC CCC
*O[C@@H]1O[C@H](CO)[C@H](O)[C@H] (O)[C@H]1NC(C)=O
*N[C@H](CCCCNC(=O)CCCC [C@H]1CCSS1)C(*)=O ...
*OP(=O)[(O-)]O[C@H]1[C@@H](O) [C@H](*)O[C@@H]1COP(=O)[(O-)] O[C@H]1[C@@H](COP(=O)[(O-)]OP(=O)[(O-)] OP(=O)[(O-)]O- O[C@H](n2cnc3c(=O)[nH]c(N)c32)[C@H]1O

図 4.4: 除外リスト

左辺第一項	左辺第二項	左辺第三項	左辺第四項
CC(=O)C(=O)[O-] NC(=O)c1ccc[n+](C@H] 2O[C@H](COP(=O)(O)OP(=O) (O)OC[C@H]3O[C@@H](n4cnc5c(N)nc nc54)[C@H](O)[C@@H]3O)[C@H](O )[C@H]2O)c1 C=C(C)CCOP(=O)[(O-)]OP (=O)[(O-)]O-]	NCCC=O	O	
N	CC(C)=CCOP(=O) (O)OP(=O)(O)O NC(=O)C1=CN([C@H]2O [C@H](COP(=O)(O)OP(=O) (O)OC[C@H]3O[C@@H](n4cnc5c(N)nc nc54)[C@H](OP(=O)(O)O)[C@H]3O) [C@H](O)[C@H]2O)C=CC1	O=C([O-])CC C(=O)C(=O)[O-]	[H+]
NC(=O)C1=CN([C@H]2O[C @H](COP(=O)(O)OP(=O)(O) OC[C@H]3O[C@@H](n4cnc5c(N)ncnc 54)[C@H](OP(=O)(O)O)[C@H]3O)[C @H](O)[C@H]2O)C=CC1 ...	NC(CC=O)C(=O)[O-]	...	
CC(=O)SCCNC(=O)CCNC(=O) [C@H](O)C(C)(C)COP(=O)[(O-)] OP(=O)[(O-)]O- O[C@H]1O[C@H](n2cnc3c(N)ncn c32)[C@H](O)[C@@H]1OP(=O)[(O-)] O-]	[NH3+]CCSCCC([NH3+]) C(=O)[O-]		

表 4.3: SMILE 反応式 (左辺)

## 機械学習モデルの構築

初めに、特徴ベクトル数が5以下のクラスを削除した。これは学習データにおいて最少クラスのデータ数が5で5分割交差検証を行えるようにし、さらに検証用学習データに分割した際に、最少クラスのデータ数が4で近傍データ数  $K = 3$  による SMOTE を実施できるようにするためである。次に正規表現を用い、3桁まで表示されている (EC X.X.-.-や EC X.-.-を除く) EC 番号を除外し、EC 番号3桁までの分類クラスを作成した。このデータに対して、各クラスで 4:1 の割合になるように学習データとテストデータに分割する。

次に、学習データから EC 番号1桁目のラベルのみ抽出し、1桁目分類に対する記述子選択を行った。ここでは、`mlxtend.feature.selection (SFS)` に RF 分類器を入力し、Macro F1-Score の交差検証平均スコアで 30 種まで選択した。また、EC T.X.X ( $T = 1, 2, \dots, 6$ ) における記述子選択では、EC T.X.X のラベルを抽出し、SFS に SMOTE と RF を組み合わせたパイプラインを入力するように設定した。これにより、検証用学習データにオーバーサンプリングを適用し、同様の記述子選択を実施した。選択記述子による評価値においてスコアが4回更新されなくなった時点で注目し、その直前の記述子の組み合わせを選出する。4回目を選んだ理由として、3回目ではスコアの低い状態で打ち止めとなり、5回目付近ではスコアの増加速度が低下し、データに対して類似の記述子が増加する懸念があったためである。

全ての EC T.X.X 学習データに対し、7回の記述子選択で選ばれた記述子をマージし、重複を削除した記述子リストの要素のみを取り出すことで、表 4.1 のようなデータから次元削減された特徴ベクトルを取得する。最後に `sklearn.model.selection` の `GridSearchCV` を呼び出し、上記と同様に SMOTE と RF を組み合わせたパイプラインを入力し、検証用学習データをオーバーサンプリングしたうえでパラメータが選択された。ここで、`n_estimators` は 10~1000、`max_depth` は 1~30 の 182 通りの組合せでグリッドサーチが行われた。

最後に Macro F1-Score の最も高い、パラメータ組合せで作成されたモデルにテストデータを入力し、分類精度を評価した。`sklearn.metrics` の `classification_report` を用いることで、Precision, Recall, F1-Score, Accuracy を出力した。



## 提案システムの流れ

以下の手順で提案手法を実施する.

1. 200 次元, 特徴ベクトル数 47090 のデータを入力し, 4:1 の割合で学習データと, テストデータに分割する.
2. 5 分割交差検証を実施し, それぞれの分割で学習データを検証用学習データと検証用テストデータに分割する.
3. EC 番号 1 桁目の分類に対する記述子選択を行う.
4. 検証用データに SMOTE を適用し, EC 番号 T.X.X ( $T = 1, 2, \dots, 6$ ) の分類に対する記述子選択を行う
5. 各記述子選択で選ばれた記述子をマージし重複を削除した記述子の組み合わせ (最終記述子) を作成する.
6. 最終記述子次元の学習データに対してグリッドサーチを実行する. (5 分割交差検証で検証データに SMOTE を実行)
7. 決定したパラメータで作成した RF モデルにテストデータを適用し, Macro Average Precision・Recall, Macro F1-Score, Accuracy で EC 番号の 1~3 桁目分類の分類精度を評価する.

## 実験結果並びに考察

### § 5.1 数値実験の概要

本研究では提案手法の数値実験(本実験)の前に予備実験を2つ行う。予備実験1では、従来行われたEC 3の2,3桁目のクラス分類に対してSMOTEを適用し、SMOTE未適用と適用した場合の分類精度を比較することで、オーバーサンプリングの有効性を示す。従来のクラス分類では、Kyoto Encyclopedia of Genes and Genomes(KEGG) [18] から取得した基質に水を加えることで基質が2つの生成物に分解する、左2辺、右2辺の加水分解反応のEC番号分類が行われていた。962種20クラスの特徴ベクトルに対してRFに特徴選択とグリッドサーチを適用し、769種のテストデータを用いて分類モデルを構築した。そして、193種のテストデータに対して予測を行い、全体のMacro F1-Score, Accuracyとしてそれぞれ、0.81, 0.92の精度を得ている。しかし、クラスのデータ分布が6~122となっており、少数クラスに対する分類精度が問題視されていた。5分割交差検証の検証用学習データに対してSFSにSMOTEとRFを組み込んだパイプラインを適用し、閾値を最多クラスとしてオーバーサンプリングを行った。そして、グリッドサーチにもパイプラインを入力することでモデルを構築し、テストデータに対する予測精度の比較を行った。

予備実験2では、本実験のための記述子選択を行うが、記述子選択後に各クラス分類で、グリッドサーチを適用し、EC番号1桁目分類、EC T.X.X (T= 1, 2, ..., 6) クラス単体でのテストデータに対する分類を行う。図 refMLframe2 に実験手順を示す。本実験には、30種までの記述子組合せの中でMacro F1-Scoreが4回未更新状態になる直前の組合せが用いられるが、この実験では、最も高くなる記述子の組合せを用いる。また、いずれもSMOTEとRFのパイプラインを適用し、7回それぞれのグリッドサーチで異なるパラメータが決定されることで、特定のクラスに特化した分類モデルが作成される。

本実験では、記述子組合せを全てマージし、重複削除したものを選択記述子として用いる。EC番号3桁分(EC X.X.X)の多クラス分類に対して、データ数37672の学習データに対して、グリッドサーチにSMOTEとRFのパイプラインを適用した5分割交差検証によるパラメータ調整を行う。ここで、検証用学習データの少数クラスの増量閾値を、最少クラス $\times 100 \times 0.8 = 400$ とし、n\_estimatorsとmax\_depthの調整範囲を、それぞれ100~1000, 10~100の100通りでグリッドサーチが行われた。最後に最適パラメータで作成されたモデルに、データ数9418のテストデータを入力し、平均のPrecision, Recall, F1-Score, およびAccuracyで分類精度を評価した。

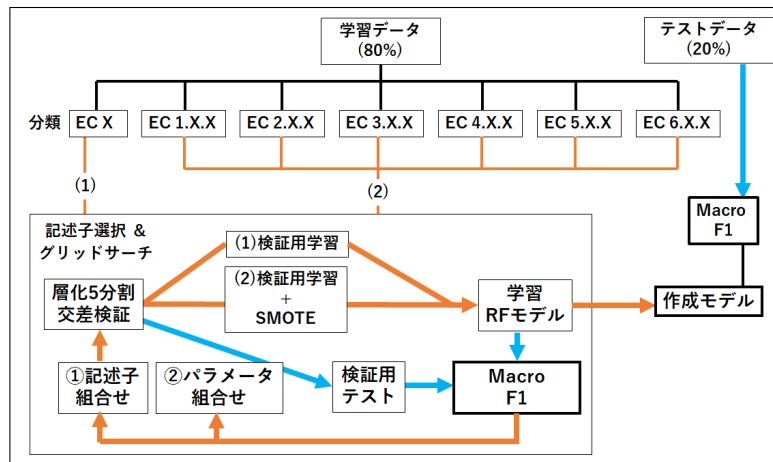


図 5.1: 実験2の流れ

従来モデル					SMOTE適用後				
	precision	recall	f1-score	support		precision	recall	f1-score	support
3.1.1.1.	0.96	0.96	0.96	25	3.1.1.1.	1.00	0.96	0.98	25
3.1.1.2.	0.92	1.00	0.96	12	3.1.1.2.	1.00	1.00	1.00	12
3.1.1.3.	0.91	0.94	0.92	31	3.1.1.3.	0.97	0.94	0.95	31
3.1.1.4.	0.86	1.00	0.92	6	3.1.1.4.	1.00	1.00	1.00	6
3.1.1.6.	1.00	1.00	1.00	3	3.1.1.6.	0.75	1.00	0.86	3
3.1.1.7.	0.00	0.00	0.00	2	3.1.1.7.	1.00	1.00	1.00	2
3.13.1.1.	1.00	0.50	0.67	2	3.13.1.1.	0.50	0.50	0.50	2
3.2.1.1.	0.96	0.96	0.96	26	3.2.1.1.	1.00	0.96	0.98	26
3.2.2.1.	0.83	1.00	0.91	5	3.2.2.1.	0.71	1.00	0.83	5
3.3.2.1.	1.00	1.00	1.00	1	3.3.2.1.	1.00	1.00	1.00	1
3.4.13.1.	0.00	0.00	0.00	1	3.4.13.1.	0.00	0.00	0.00	1
3.4.19.1.	1.00	1.00	1.00	1	3.4.19.1.	1.00	1.00	1.00	1
3.5.1.1.	0.94	0.97	0.95	31	3.5.1.1.	0.94	0.97	0.95	31
3.5.3.1.	0.83	1.00	0.91	5	3.5.3.1.	1.00	1.00	1.00	5
3.5.4.1.	0.89	0.89	0.89	9	3.5.4.1.	0.88	0.78	0.82	9
3.5.5.1.	1.00	1.00	1.00	2	3.5.5.1.	1.00	1.00	1.00	2
3.5.99.1.	1.00	0.50	0.67	2	3.5.99.1.	0.67	1.00	0.80	2
3.6.1.1.	0.86	0.95	0.90	19	3.6.1.1.	0.89	0.89	0.89	19
3.7.1.1.	1.00	0.71	0.83	7	3.7.1.1.	0.88	1.00	0.93	7
3.8.1.1.	1.00	0.67	0.80	3	3.8.1.1.	1.00	0.67	0.80	3
accuracy			0.92	193	accuracy			0.94	193
macro avg	0.85	0.80	0.81	193	macro avg	0.86	0.88	0.87	193
weighted avg	0.91	0.92	0.91	193	weighted avg	0.94	0.94	0.94	193

図 5.2: 予備実験1の結果

## § 5.2 実験結果と考察

2つの予備実験の結果，および本実験の結果と考察を述べる．

### 予備実験1の結果

従来の EC 3 の 2,3 桁目に対する RF×SFS 分類モデルと，新たに SMOTE を適用した分類モデルの比較結果を図 5.2 に示す．Macro F1-Score が従来に比べて 0.06 向上し，SMOTE の有効性が示された．また，少数クラス (EC3.4.13) に関して，SMOTE を適用したモデルでも正しい EC 番号に分類されていない結果となった．

### 予備実験2の結果

EC 番号 1 桁目，EC T.X.X (T= 1, 2, ..., 6) に対する各クラス分類で得られた記述子リスト，およびテストデータに対するクラス分類結果を表 5.1, および図 5.3, 5.4, 5.5, 5.6, に示す．得られた記述子をマージし，重複を削除すると，93 種の記述子が得られた．

表 5.1: 予備実験 2 で選ばれた記述子リスト

	1	2	3	4	5	6	7	8	9	10	11	12	13	14
EC X	MinEState Index	MolWt	NumValence Electrons	FpDensity Morgan1	BCUT2D _MWHI	BCUT2D _MWLOW	BCUT2D _CHGHI	HallKierAlpha	Ipc	Kappa3	PEOE_VSA4	SlogP_VSA2	SlogP_VSA3	FractionCSP3
EC 1.X.X	MinEState Index	BCUT2D _MWLOW	BalabanJ	Chi4v	Kappa1	Kappa3	PEOE_VSA4	SMR_VSA1	SMR_VSA10	SMR_VSA3	SMR_VSA7	SlogP_VSA2	EState_VSA1	EState_VSA3
EC 2.X.X	AvlJpc	Chi1	Ipc	Kappa1	PEOE_VSA2	SlogP_VSA12	SlogP_VSA3	EState_VSA4	VSA_EState1	VSA_EState6	VSA_EState7	VSA_EState8	VSA_EState9	NumAliphatic Carbocycles
EC 3.X.X	BCUT2D _MWHI	BCUT2D _MWLOW	BCUT2D _LOGLOW	Chi1	Chi2n	Chi2v	Chi3v	Chi4v	Kappa1	Kappa3	PEOE_VSA12	PEOE_VSA5	PEOE_VSA9	SMR_VSA10
EC 4.X.X	NumValence Electrons	BCUT2D _MWHI	BCUT2D _MWLOW	BCUT2D _LOGPHI	LabuteASA	SMR_VSA5	SMR_VSA6	VSA_EState2	VSA_EState6	VSA_EState7	RingCount	fr_Al_OH _noTert	fr_COO2	fr_C_O _noCOO
EC 5.X.X	AvlJpc	BalabanJ	Chi1n	Chi2n	Chi3v	Chi4v	HallKierAlpha	VSA_EState5	Num Heteroatoms	fr_C_O	fr_aldehyde	fr_ester	fr_ketone _Topliss	end13
EC 6.X.X	BCUT2D _MWLOW	AvlJpc	HallKierAlpha	Ipc	EState_VSA3	VSA_EState6	Fraction CSP3	NumAromatic Heterocycles	fr_COO	fr_benzene	fr_ether	fr_halogen	fr_piperzine	fr_sulfide
	15	16	17	18	19	20	21	22	23	24	25	26	27	28
EC X	fr_Ar_COO	fr_COO	fr_COO2	fr_allylic_oxid	fr_lactone	end 19								
EC 1.X.X	EState_VSA7	EState_VSA8	VSA_EState6	VSA_EState8	VSA_EState9	NumAliphatic Heterocycles	NumSaturated Carbocycles	fr_Ar_COO	fr_Ar_OH	fr_NH2	fr_SH	fr_aniline	fr_guanido	end27
EC 2.X.X	fr_Al_OH	fr_C_S	fr_epoxide	fr_ether	fr_lactone	fr_phos_ester	end20							
EC 3.X.X	SlogP_VSA8	VSA_EState6	NumAliphatic Rings	NumAromatic Heterocycles	NumRotatable Bonds	NumSaturated Carbocycles	fr_C_O	fr_NH0	fr_guanido	fr_hdrzone	fr_ketone	fr_ketone _Topliss	fr_pyridine	fr_sulfide
EC 4.X.X	fr_NH1	fr_aldehyde	fr_benzene	fr_halogen	fr_methoxy	fr_sulfone	end 21							↑ end28
EC 5.X.X	end13													
EC 6.X.X	fr_unbrch _alkane	end 15												

## 本実験の結果

選ばれた 93 種の記述子を用いてグリッドサーチを行うと、パラメータとして  $n\_estimators=300$ ,  $max\_depth=90$  が得られた。テストデータに対する分類結果は表 5.7, 5.7 のようになった。

## 考察 1

最終的な予測精度が Macro F1-Score 0.79 となり、先行研究 [4] の EC 番号 1~3 桁目の予測精度  $0.77 \pm 0.01$  とほぼ同等の性能となった。しかし、先行研究と比べる記述子選択に多くの時間がかかる点。また、データ加工で削除したデータがあり、先行研究よりもデータ数が少ない点などから、改善すべき課題が多く存在する。一方で、今回選択された記述子について解析し、EC 番号 1 桁目や 2,3 桁目予測で選択された記述子に共通の特徴や傾向を見出し、モデル学習に利用することができれば、酵素による化学反応の特徴をより詳細に捉えることができると考えられる。また、フィンガープリントは構造情報のみだが、特性値は化学反応の情報をより深く説明できるため、記述子選択時に重要な特徴量に重みづけをするなどの方法を用いることで先行研究よりもより説明力のあるモデル構築が可能になると考えられる。

## 考察 2

今後の展望として、2 つ挙げられる。1 つ目は、EC 番号の 4 桁目を予測する手法の開発である。4 桁目の分類は 3 桁目の分類よりもさらに細かく、酵素同士の違いを区別するのがより難しい、酵素反応ベースの研究では、4 桁目を予測した研究例が少ないため、より最適な酵素候補を見つけることにつながる。2 つ目は実際の有機合成に用いることである。実際に使ってもらった結果をモデルに反映することができればより、応用性のあるモデル構築が期待できる。

	precision	recall	f1-score	support		precision	recall	f1-score	support
1.1.1	0.93	0.96	0.94	440	1.2.99	0.67	1.00	0.80	2
1.1.3	0.86	0.76	0.81	33	1.20.1	1.00	1.00	1.00	1
1.1.5	0.79	0.94	0.86	16	1.21.1	1.00	1.00	1.00	2
1.1.98	0.67	0.67	0.67	3	1.21.3	0.50	0.67	0.57	3
1.1.99	0.64	0.78	0.70	18	1.23.5	1.00	1.00	1.00	1
1.10.3	0.86	0.86	0.86	7	1.3.1	0.92	0.92	0.92	85
1.10.5	0.00	0.00	0.00	1	1.3.3	0.87	0.81	0.84	16
1.10.99	0.50	0.50	0.50	2	1.3.5	1.00	0.92	0.96	13
1.11.1	0.81	0.85	0.83	20	1.3.8	1.00	0.92	0.96	13
1.11.2	0.64	0.70	0.67	10	1.3.98	0.33	1.00	0.50	1
1.12.98	1.00	1.00	1.00	1	1.3.99	1.00	0.33	0.50	3
1.13.11	0.88	0.95	0.91	61	1.4.1	0.90	1.00	0.95	18
1.13.12	0.69	0.60	0.64	15	1.4.3	0.92	0.85	0.88	39
1.13.99	0.00	0.00	0.00	1	1.4.5	1.00	1.00	1.00	2
1.14.11	0.82	0.84	0.83	44	1.4.99	1.00	1.00	1.00	5
1.14.12	0.96	0.77	0.85	30	1.5.1	0.92	0.85	0.88	27
1.14.13	0.80	0.90	0.85	192	1.5.3	0.67	0.57	0.62	7
1.14.14	0.91	0.71	0.80	102	1.5.5	0.33	0.50	0.40	2
1.14.15	0.50	0.25	0.33	4	1.5.99	0.86	1.00	0.92	6
1.14.16	1.00	0.60	0.75	5	1.6.3	0.67	1.00	0.80	2
1.14.17	0.75	1.00	0.86	3	1.6.5	0.85	0.81	0.83	27
1.14.18	0.73	0.67	0.70	12	1.7.1	1.00	0.85	0.92	13
1.14.19	0.91	1.00	0.95	10	1.7.2	1.00	1.00	1.00	1
1.14.20	0.78	0.82	0.80	17	1.7.3	1.00	1.00	1.00	4
1.14.21	1.00	0.50	0.67	4	1.7.5	1.00	0.50	0.67	2
1.14.99	0.95	0.86	0.90	22	1.7.99	0.50	1.00	0.67	1
1.17.1	1.00	0.50	0.67	8	1.8.1	1.00	0.57	0.73	7
1.17.3	0.50	0.67	0.57	6	1.8.3	1.00	1.00	1.00	2
1.17.5	1.00	0.50	0.67	2	1.8.5	0.88	1.00	0.93	7
1.18.1	0.00	0.00	0.00	2	1.97.1	1.00	1.00	1.00	2
1.2.1	0.96	0.96	0.96	170					
1.2.3	0.69	0.69	0.69	16	accuracy			0.88	1601
1.2.4	1.00	0.89	0.94	9	macro avg	0.79	0.78	0.77	1601
1.2.5	0.50	1.00	0.67	1	weighted avg	0.88	0.88	0.88	1601

図 5.3: EC 1.X.X に対する分類

	precision	recall	f1-score	support		precision	recall	f1-score	support
2.1.1	1.00	1.00	1.00	694	3.1.1	0.97	0.98	0.97	570
2.1.2	1.00	1.00	1.00	4	3.1.2	0.93	0.91	0.92	43
2.1.3	1.00	0.60	0.75	5	3.1.3	0.87	0.91	0.89	148
2.1.4	0.00	0.00	0.00	2	3.1.4	0.90	0.69	0.78	26
2.2.1	0.80	0.89	0.84	18	3.1.6	1.00	1.00	1.00	2
2.3.1	1.00	1.00	1.00	1827	3.1.8	1.00	0.83	0.91	6
2.3.2	0.75	1.00	0.86	3	3.13.1	0.00	0.00	0.00	2
2.3.3	1.00	1.00	1.00	5	3.2.1	0.93	0.98	0.95	113
2.4.1	0.97	0.98	0.97	171	3.2.2	0.92	1.00	0.96	12
2.4.2	0.97	0.88	0.92	41	3.3.1	0.50	1.00	0.67	1
2.4.99	1.00	0.50	0.67	2	3.3.2	1.00	0.88	0.93	16
2.5.1	0.97	0.92	0.95	78	3.4.11	0.78	0.90	0.84	20
2.6.1	0.86	0.92	0.88	71	3.4.13	0.82	0.78	0.80	18
2.7.1	0.96	0.97	0.96	153	3.4.17	1.00	0.83	0.91	6
2.7.2	1.00	0.62	0.77	8	3.4.19	1.00	1.00	1.00	6
2.7.3	0.50	1.00	0.67	1	3.4.21	1.00	1.00	1.00	2
2.7.4	0.75	0.87	0.81	31	3.5.1	0.89	0.84	0.86	87
2.7.6	1.00	0.57	0.73	7	3.5.2	0.71	0.71	0.71	7
2.7.7	0.96	0.99	0.98	107	3.5.3	1.00	1.00	1.00	4
2.7.8	1.00	1.00	1.00	2518	3.5.4	0.86	0.80	0.83	15
2.7.9	1.00	0.50	0.67	2	3.5.5	0.95	1.00	0.97	18
2.8.1	1.00	0.50	0.67	4	3.5.99	1.00	0.33	0.50	3
2.8.2	0.95	0.95	0.95	21	3.6.1	0.73	0.62	0.67	58
2.8.3	0.93	0.88	0.90	16	3.6.3	1.00	1.00	1.00	126
					3.6.4	0.33	1.00	0.50	1
accuracy			0.99	5789	3.7.1	0.84	0.94	0.89	17
macro avg	0.89	0.81	0.83	5789	3.8.1	0.94	0.94	0.94	18
weighted avg	0.99	0.99	0.99	5789	accuracy			0.93	1345
					macro avg	0.85	0.85	0.83	1345
					weighted avg	0.93	0.93	0.93	1345

図 5.4: EC 2.X.X, 3.X.X に対する分類

	precision	recall	f1-score	support		precision	recall	f1-score	support
4.1.1	0.98	0.98	0.98	260					
4.1.2	0.93	0.96	0.95	28					
4.1.3	0.44	0.50	0.47	8					
4.1.99	0.78	0.70	0.74	10	5.1.1	0.00	0.00	0.00	1
4.2.1	0.93	0.95	0.94	88	5.1.3	0.75	0.75	0.75	4
4.2.2	1.00	1.00	1.00	1	5.2.1	0.00	0.00	0.00	1
4.2.3	0.95	0.87	0.91	23	5.3.1	0.92	1.00	0.96	11
4.3.1	1.00	0.89	0.94	9	5.3.2	1.00	1.00	1.00	4
4.3.2	0.50	0.50	0.50	2	5.3.3	0.92	1.00	0.96	11
4.3.3	1.00	0.33	0.50	3	5.3.99	1.00	0.67	0.80	3
4.4.1	0.92	0.92	0.92	24	5.4.2	0.60	1.00	0.75	3
4.5.1	1.00	0.75	0.86	4	5.4.3	1.00	0.67	0.80	3
4.6.1	1.00	1.00	1.00	1	5.4.4	0.67	1.00	0.80	2
4.99.1	0.33	1.00	0.50	1	5.4.99	0.83	1.00	0.91	5
					5.5.1	1.00	0.89	0.94	19
accuracy			0.94	462	accuracy			0.90	67
macro avg	0.84	0.81	0.80	462	macro avg	0.72	0.75	0.72	67
weighted avg	0.95	0.94	0.94	462	weighted avg	0.89	0.90	0.88	67

図 5.5: EC 4.X.X, 5.X.X に対する分類

	precision	recall	f1-score	support		precision	recall	f1-score	support
6.1.2	0.50	0.50	0.50	2	1	0.98	0.98	0.98	1601
6.2.1	0.99	0.99	0.99	67	2	0.99	0.99	0.99	5789
6.3.1	0.86	1.00	0.92	6	3	0.97	0.98	0.98	1345
6.3.2	0.98	0.98	0.98	59	4	0.95	0.94	0.94	462
6.3.4	0.86	0.75	0.80	8	5	0.83	0.72	0.77	67
6.3.5	0.71	0.71	0.71	7	6	0.99	0.92	0.95	154
6.4.1	1.00	1.00	1.00	5					
accuracy			0.95	154	accuracy			0.98	9418
macro avg	0.84	0.85	0.84	154	macro avg	0.95	0.92	0.94	9418
weighted avg	0.95	0.95	0.95	154	weighted avg	0.98	0.98	0.98	9418

図 5.6: EC 6.X.X, 1 桁目に対する分類

EC Class	Num of Data	precision	recall	f1-score	EC Class	Num of Data	precision	recall	f1-score	EC Class	Num of Data	precision	recall	f1-score
1.1.1	440	0.95	0.95	0.95	1.2.1	170	0.97	0.96	0.97	1.8.1	7	0.80	0.57	0.67
1.1.3	33	0.89	0.76	0.82	1.2.3	16	0.68	0.81	0.74	1.8.3	2	1.00	1.00	1.00
1.1.5	16	0.84	1.00	0.91	1.2.4	9	0.88	0.78	0.82	1.8.5	7	0.78	1.00	0.88
1.1.98	3	1.00	0.67	0.80	1.2.5	1	1.00	1.00	1.00	1.97.1	2	1.00	1.00	1.00
1.1.99	18	0.76	0.72	0.74	1.2.99	2	1.00	1.00	1.00	2.1.1	694	1.00	1.00	1.00
1.10.3	7	0.57	0.57	0.57	1.20.1	1	1.00	1.00	1.00	2.1.2	4	1.00	1.00	1.00
1.10.5	1	0.00	0.00	0.00	1.21.1	2	1.00	1.00	1.00	2.1.3	5	1.00	0.80	0.89
1.10.99	2	1.00	0.50	0.67	1.21.3	3	0.67	0.67	0.67	2.1.4	2	0.00	0.00	0.00
1.11.1	20	0.80	0.80	0.80	1.23.5	1	1.00	1.00	1.00	2.2.1	18	0.88	0.78	0.82
1.11.2	10	0.70	0.70	0.70	1.3.1	85	0.93	0.92	0.92	2.3.1	1827	0.99	0.99	0.99
1.12.98	1	0.50	1.00	0.67	1.3.3	16	0.93	0.81	0.87	2.3.2	3	0.00	0.00	0.00
1.13.11	61	0.92	0.93	0.93	1.3.5	13	1.00	0.85	0.92	2.3.3	5	0.83	1.00	0.91
1.13.12	15	0.71	0.67	0.69	1.3.8	13	0.92	0.92	0.92	2.4.1	171	0.97	0.97	0.97
1.13.99	1	1.00	1.00	1.00	1.3.98	1	0.25	1.00	0.40	2.4.2	41	0.94	0.83	0.88
1.14.11	44	0.90	0.84	0.87	1.3.99	3	1.00	0.33	0.50	2.4.99	2	1.00	1.00	1.00
1.14.12	30	0.96	0.80	0.87	1.4.1	18	0.89	0.94	0.92	2.5.1	78	0.92	0.86	0.89
1.14.13	192	0.83	0.87	0.85	1.4.3	39	0.95	0.95	0.95	2.6.1	71	0.84	0.86	0.85
1.14.14	102	0.88	0.71	0.78	1.4.5	2	1.00	1.00	1.00	2.7.1	153	0.95	0.95	0.95
1.14.15	4	0.33	0.25	0.29	1.4.99	5	1.00	1.00	1.00	2.7.2	8	0.83	0.63	0.71
1.14.16	5	1.00	0.60	0.75	1.5.1	27	0.86	0.93	0.89	2.7.3	1	0.50	1.00	0.67
1.14.17	3	0.75	1.00	0.86	1.5.3	7	1.00	0.57	0.73	2.7.4	31	0.76	0.90	0.82
1.14.18	12	0.56	0.75	0.64	1.5.5	2	0.50	0.50	0.50	2.7.6	7	0.80	0.57	0.67
1.14.19	10	0.83	1.00	0.91	1.5.99	6	0.86	1.00	0.92	2.7.7	107	0.97	0.99	0.98
1.14.20	17	0.74	0.82	0.78	1.6.3	2	0.67	1.00	0.80	2.7.8	2518	1.00	1.00	1.00
1.14.21	4	0.67	0.50	0.57	1.6.5	27	0.79	0.85	0.82	2.7.9	2	1.00	0.50	0.67
1.14.99	22	0.90	0.86	0.88	1.7.1	13	1.00	0.85	0.92	2.8.1	4	1.00	1.00	1.00
1.17.1	8	1.00	0.63	0.77	1.7.2	1	1.00	1.00	1.00	2.8.2	21	1.00	1.00	1.00
1.17.3	6	0.56	0.83	0.67	1.7.3	4	1.00	1.00	1.00	2.8.3	16	0.79	0.94	0.86
1.17.5	2	1.00	0.50	0.67	1.7.5	2	0.50	0.50	0.50					
1.18.1	2	0.00	0.00	0.00	1.7.99	1	0.00	0.00	0.00					

図 5.7: EC 番号 1 桁～3 桁目までの分類結果 1

EC Class	Num of Data	precision	recall	f1-score	EC Class	Num of Data	precision	recall	f1-score
3.1.3	148	0.89	0.90	0.89	3.1.1	570	0.97	0.96	0.96
3.1.4	26	0.95	0.81	0.88	3.1.2	43	0.91	0.91	0.91
3.1.6	2	1.00	1.00	1.00	4.2.2	1	0.50	1.00	0.67
3.1.8	6	1.00	1.00	1.00	4.2.3	23	0.95	0.78	0.86
3.13.1	2	0.00	0.00	0.00	4.3.1	9	1.00	0.89	0.94
3.2.1	113	0.91	0.96	0.94	4.3.2	2	1.00	1.00	1.00
3.2.2	12	1.00	1.00	1.00	4.3.3	3	0.67	0.67	0.67
3.3.1	1	1.00	1.00	1.00	4.4.1	24	0.77	1.00	0.87
3.3.2	16	0.88	0.94	0.91	4.5.1	4	0.50	0.50	0.50
3.4.11	20	0.76	0.80	0.78	4.6.1	1	1.00	1.00	1.00
3.4.13	18	0.76	0.89	0.82	4.99.1	1	1.00	1.00	1.00
3.4.17	6	0.71	0.83	0.77	5.1.1	1	1.00	1.00	1.00
3.4.19	6	1.00	1.00	1.00	5.1.3	4	0.50	0.50	0.50
3.4.21	2	1.00	1.00	1.00	5.2.1	1	0.00	0.00	0.00
3.5.1	87	0.84	0.91	0.87	5.3.1	11	0.75	0.82	0.78
3.5.2	7	1.00	1.00	1.00	5.3.2	4	0.60	0.75	0.67
3.5.3	4	0.67	1.00	0.80	5.3.3	11	0.82	0.82	0.82
3.5.4	15	0.87	0.87	0.87	5.3.99	3	0.67	0.67	0.67
3.5.5	18	0.94	0.94	0.94	5.4.2	3	0.75	1.00	0.86
3.5.99	3	0.25	0.33	0.29	5.4.3	3	1.00	0.67	0.80
3.6.1	58	0.68	0.69	0.68	5.4.4	2	0.33	0.50	0.40
3.6.3	126	1.00	1.00	1.00	5.4.99	5	0.71	1.00	0.83
3.6.4	1	0.33	1.00	0.50	5.5.1	19	0.83	0.79	0.81
3.7.1	17	0.81	1.00	0.89	6.1.2	2	1.00	0.50	0.67
3.8.1	18	1.00	0.89	0.94	6.2.1	67	0.99	0.99	0.99
4.1.1	260	0.97	0.99	0.98	6.3.1	6	1.00	0.83	0.91
4.1.2	28	0.75	0.96	0.84	6.3.2	59	0.95	0.97	0.96
4.1.3	8	1.00	0.50	0.67	6.3.4	8	0.78	0.88	0.82
4.1.99	10	1.00	0.70	0.82	6.3.5	7	1.00	0.29	0.44
4.2.1	88	0.93	0.93	0.93	6.4.1	5	1.00	0.80	0.89
					Total	9418			
					macro average		0.812184	0.802785	0.792879
					weighted average		0.956119	0.954767	0.954222
					accuracy		0.954767		

図 5.8: EC 番号 1 桁～3 桁目までの分類結果 2

### おわりに

近年、新型コロナウイルスになどの影響で、新薬開発の需要が高まり、化学反応の設計や予測を行う研究が発展を続けている。一方で、反応の効率化と環境の面から、酵素の生体触媒を用いて合成が行われる機会が増えており、目的の反応に対して最適な酵素を予測することが重要視されている。しかし、基質特異性などの酵素の性質は生物分野にかかわるため、有機合成の知識のみでは解決が難しく、酵素研究の専門家と協力する、または、酵素データベースを参照するなどして最適な酵素候補は探索されていた。

目的とする反応に対して、酵素候補を予測するシステムがあれば、次のステップである、1つの酵素に絞るスクリーニングまでスムーズに進めることができる。本研究では、ターゲット反応式と EC 反応式を比較し、類似する EC 反応式の EC 番号を、最適な酵素候補として予測する方法を提案した。化合物の物理・化学的な特性値を計算し、反応物から生成物への特性値変化量を要素とする 210 次元の特徴ベクトルをもとに、EC 番号の 3 桁目まで予測する手法を開発した。

本研究では特徴選択とオーバーサンプリングを組合せて分類精度に寄与する記述子を選択し、約 47,000 あるデータの予測を行った。結果として従来手法を上回ることが出来なかったが、選択された記述子組合せを分析することで、より詳細な予測ができると考えられる。以上の結果から、ターゲット反応式に最適な酵素候補を提示するシステムとしての性能を高めるためには、構造変化の特徴抽出をより詳細にする、あるいは EC 番号に対して、さらに別の分類規則を加えるなどの工夫が必要であることが明らかになった。また、最適な酵素候補と予測された EC 番号に対して、その酵素が実際に優れた反応を示すかどうか、検証する必要があると考えられる。

今後の課題として、特徴ベクトルのクラスタリングなどの、EC 番号以外の分類法の追加などが挙げられる。





# 謝辞

本研究を遂行するにあたり，多大なご指導とご鞭撻を賜りました，富山県立大学工学部 電子・情報工学科情報基盤工学講座 奥原浩之教授，António Oliveira Nzinga René 講師に深く感謝の意を表します．また，同大学工学部生物工学科酵素化学工学講座 浅野泰久教授，同大学くすりのシリコンバレー TOYAMA 研究拠点化プロジェクトディレクター補佐 岩崎源司博士には，有機化学・酵素分野の立場から大変貴重なご意見，および当該分野に関して，一からご指導を賜りました．心よりお礼申し上げます．さらに、ケモインフォマティクスにおける機械学習や化学構造の表現法についてご助言を賜りました、国立研究開発法人医薬基盤・健康・栄養研究所上級研究員・プロジェクトリーダー 荒木通啓博士，神戸大学大学院工学研究科応用化学専攻 渡邊直暉氏に深く感謝申し上げます。最後になりましたが，多大な協力をしていただいた研究室の同輩諸氏に感謝致します．

2024 年 2 月

武藤 克弥



## 参考文献

- [1] Tamas Benkovics, John A. McIntosh, Steven M. Silverman, Jongrock Kong, Peter Maligres, Tetsuji Itoh, Hao Yang, Mark A. Huffman, Deeptak Verma, Weilan Pan, Hsing-I Ho, Jonathan Vroom, Anders Knight, Jessica Hurtak, William Morris, Neil A. Strotman, Grant Murphy, Kevin M. Maloney, and Patrick S. Fierl, “Evolving to an Ideal Synthesis of Molnupiravir, an Investigational Treatment for COVID - 19”, *ChemRxiv*, 2020.
- [2] “Enzyme Nomenclature”, <https://iubmb.qmul.ac.uk/enzyme/>, 閲覧日 2024.1.21.
- [3] Naoki Watanabe, Masaki Yamamoto, Masahiro Murata, Yuki Kuriya, Michihiro Araki, “EnzymeNet: residual neural networks model for Enzyme Commission number prediction”, *Bioinformatics Advances*, Vol. 3, No. 1, 2023.
- [4] Probst Daniel, “An explainability framework for deep learning on chemical reactions exemplified by enzyme-catalysed reaction classification”, *Journal of Cheminformatics*, Vol. 15, No. 1, pp. 113, 2023.
- [5] Probst, Daniel, “An explainability framework for deep learning on chemical reactions exemplified by enzyme-catalysed reaction classification”, *Journal of chemical information and modeling*, Vol. 49, No. 7, pp. 1839-1846, 2009.
- [6] 北川勲, 磯部稔, “天然物化学・生物有機化学 I”, 朝倉書店, 2008.
- [7] 西村淳, 樋口弘行, 大和武彦, “有機合成化学入門 -基礎を理解して実践に備える”, 丸善株式会社, 2010.
- [8] “日本化学会・ケモインフォマティクス部会”, <https://cicsj.csj.jp/>, 閲覧日 2022.1.23.
- [9] Yu Chenggang, Zavaljevski Nela, Desai Valmik and Reifman Jaques, “Genome-wide enzyme annotation with precision control: Catalytic families (CatFam) databases”, *Proteins: Structure, Function, and Bioinformatics*, Vol. 74, No. 2, pp. 449-460, 2009.
- [10] 中野裕太, 瀧川一学, “化学反応ネットワークにおける最適反応経路候補の列挙”, 情報処理学会研究報告, Vol. 122, No. 16, 2019.
- [11] 佐藤寛子, “化学情報学 - 化学反応の系図と反応予測 -” 国立情報学研究所, 2003.
- [12] 藤波 美起登, 清野 淳司, “量子化学計算情報を記述子とした機械学習に基づく反応予測手法の開発”, *Journal of Computer Chemistry, Japan*, Vol. 15, No. 3, pp. 63-65, 2016.
- [13] “特異なタンパク質進化 Circular permutation による酵素の機能改変”, <https://www.amano-enzyme.co.jp/corporate/foundation/pdf/19/pg09.pdf>, 閲覧日 2022.1.25.

- [14] “酵素の化学”, <http://www.sc.fukuoka-u.ac.jp/~bc1/Biochem/biochem5.htm>, 閲覧日 2022.1.31.
- [15] “酵素基質とは”, <https://bizcomjapan.co.jp/iris-biotech/knowledge/substrate/>, 閲覧日 2022.1.31.
- [16] “新設された酵素分類 EC7 の和名提案について”, [https://www.jbsoc.or.jp/notice/ec\\_translocase.html](https://www.jbsoc.or.jp/notice/ec_translocase.html), 閲覧日 2022.1.15.
- [17] 白兼孝雄, “酵素の分類と命名法”, JAS 情報, 2017.
- [18] “KEGG: Kyoto Encyclopedia of Genes and Genomes”, <https://www.genome.jp/kegg/kegg-ja.html>, 閲覧日 2022.1.17.
- [19] “BRENDA The Comprehensive Enzyme Information System”, <https://www.brenda-enzymes.org/index.php>, 閲覧日 2022.2.1.
- [20] “SMILES 記法は化学構造の線形表記法”, <https://future-chem.com/smiles-smarts/>, 閲覧日 2022.1.27.
- [21] Joseph L. Durant, Burton A. Leland, Douglas R. Henry, and James G. Nourse, “Re-optimization of MDL Keys for Use in Drug Discovery”, *American Chemical Society*, Vol. 7, No. 12, 2012.
- [22] “The RDKit Documentation”, <https://www.rdkit.org/docs/GettingStartedInPython.html#list-of-available-descriptors>, 閲覧日 2022.2.6.
- [23] Wenbo Sun et al., “Machine learning-assisted molecular design and efficiency prediction for high-performance organic photovoltaic materials”, *Science advances*, Vol. 5, No. 1, eaay4275, 2019.
- [24] Sakiyama, Hiroshi, Motohisa Fukuda, and Takashi Okuno., “Prediction of blood-brain barrier penetration (bbb) based on molecular descriptors of the free-form and in-blood-form datasets”, *Molecules*, Vol. 26, No. 24, 7428, 2021.
- [25] Rogers David and Hahn Mathew, “Extended-connectivity fingerprints”, *Journal of chemical information and modeling*, Vol. 50, No. 5, pp. 742-754, 2021.
- [26] “Python: 機械学習で分類問題のモデルを評価する指標について”, <https://blog.amedama.jp/entry/2017/12/18/005311>, 閲覧日 2022.1.27.
- [27] Scott A. Wildman and Gordon M. Crippen., “Prediction of Physicochemical Parameters by Atomic Contributions”, *Journal of chemical information and computer sciences*, Vol. 39, pp. 868-873, 1999.
- [28] 小林誠一., “化合物活性予測におけるベイズモデルの利用”, *CICSJ Bulletin*, Vol. 27, No. 4, 2009.

- [29] Chawla Nitesh V., Bowyer Kevin W., Hall Lawrence O., Kegelmeyer W. Philip, “SMOTE: synthetic minority over-sampling technique”, *Journal of artificial intelligence research*, Vol 16, pp. 321-357, 2002.
- [30] Fernández Alberto, Garcia Salvador, Herrera Francisco, Chawla, Nitesh V, “SMOTE for learning from imbalanced data: progress and challenges, marking the 15-year anniversary”, *Journal of artificial intelligence research*, Vol 61, pp. 863-905, 2018.
- [31] “sklearn.ensemble.RandomForestClassifier — scikit-learn 1.4.0 documentation”, <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>, 閲覧日 2024.1.31.
- [32] “SMOTE — Version 0.12.0”, [https://imbalanced-learn.org/stable/references/generated/imblearn.over\\_sampling.SMOTE.html](https://imbalanced-learn.org/stable/references/generated/imblearn.over_sampling.SMOTE.html), 閲覧日 2024.1.31.
- [33] “sklearn.model\_selection.GridSearchCV — scikit-learn 1.4.0 documentation”, [https://scikit-learn.org/stable/modules/generated/sklearn.model\\_selection.GridSearchCV.html](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html), 閲覧日 2024.1.31.
- [34] Teuvo KOHONEN, “Self-organized formation of topologically correct feature map”, *Biological Cybernetics*, Vol. 43, pp. 59–69, 1982.
- [35] “Mlxtend.feature selection ”, [http://rasbt.github.io/mlxtend/api\\_subpackages/mlxtend.feature\\_selection/#sequentialfeatureselector](http://rasbt.github.io/mlxtend/api_subpackages/mlxtend.feature_selection/#sequentialfeatureselector), 閲覧日 2022.3.8.
- [36] Sebastian Raschka, Vahid Mirjalili 著, 株式会社クイープ訳, 福島真太郎監訳, “[第3版] Python 機械学習プログラミング 達人データサイエンティストによる理論と実践”, 株式会社インプレス, 2020.
- [37] Chandra, Prem and Enespa and Singh, Ranjan and Arora, Pankaj Kumar “Microbial lipases and their industrial applications: a comprehensive review”, *Microbial cell factories*, Vol 19, pp. 1-42, 2020.