

離散化による解釈可能な分類モデルの構築

Learning Interpretable Classification Rules Using Discretization

米田 友花*¹ 杉山 磨人*² 鷲尾 隆*¹
 Yuka Yoneda Mahito Sugiyama Takashi Washio

*¹大阪大学 産業科学研究所

The Institute of Scientific and Industrial Research, Osaka University

*²独立行政法人科学技術振興機構, さきがけ
 JST PRESTO

We propose to learn classification rules from supervised continuous data via a two-step procedure: We first binarize data points and then perform feature selection on the binarized data to obtain minimum necessary binary features that are required to classify given data. Each binary feature represents an interval on the original feature space, hence the user can directly interpret classification rules as intersections of intervals. Although such *interpretability* is indispensable in machine learning applications from biology to economics, we cannot often interpret classification rules obtained by recently emerging highly accurate machine learning methods such as deep learning. We experimentally demonstrate that our method can learn simpler classification rules than decision tree classifiers while keeping reasonable accuracy.

1. はじめに

近年, 連続値データからのクラス分類において, 深層学習と呼ばれるニューラルネットワークに基づく機械学習手法が着目されている [4]. 深層学習では, 膨大な計算量とデータを必要とする一方, 高い分類精度を達成することが可能となってきたが, 学習モデルが複雑すぎるため, 人間が分類規則を解釈することができない, という欠点がある. これに対して, CART, ID3, C4.5 [3] などの決定木 (decision tree) を用いた分類手法に代表される, 特徴空間の矩形分割に基づく教師あり学習手法が知られている. これらの手法では, クラスの分類規則を, 人間が理解することのできる表現形式として明示的に獲得することができるが, 一般的に分類精度が低く, 特に多変数の連続値データに対する簡潔な規則の獲得が難しい. そこで, 多変数の連続値データに対して高い精度を達成ことができ, かつ解釈可能性の高い簡潔な分類規則を獲得する機械学習手法が求められている.

この問題を解決するため, 本研究では, 連続値データから解釈可能かつ簡潔なクラス分類規則を学習する手法を提案する. 提案手法ではまず, データの各特徴を十分に離散化 (discretization) することで, クラスラベルに対して無矛盾な 2 値特徴の集合を生成する. そして, 2 値データに対する教師あり特徴選択手法 Super-CWC アルゴリズム [6] を適用して不要な特徴を削除することで, クラスラベルを判別するために必要最低限の 2 値特徴を獲得する. 各 2 値特徴の組合せが多次元空間上の矩形領域に対応しており, これらの分類規則を用いることで未知データを分類することができる.

本稿は, 以下の様な構成となっている. まず 2 節で提案手法を説明する. 最初に概要を説明したあと, 2.1 節では離散化, 2.2 節では提案手法の一部で使った Super-CWC アルゴリズムについて述べる. 3 節で提案手法の動作例を示したあと, 4 節で提案

手法を実験によって評価し, 結果を考察する. 最後に 5 節で本研究の結論と今後の課題について述べる.

2. 提案手法

提案手法は, 主に 2 つの段階からなる. 第一段階では, 連続値データの各特徴をそれぞれ区間を表現する複数の 2 値特徴へと変換することで離散化し, 2 値データへ変換する. 第二段階では, 第一段階で生成した 2 値データを入力とし, Super-CWC アルゴリズムを適用することで, 不要な特徴を削除しクラスラベルに無矛盾な最小の 2 値特徴集合を獲得する.

提案手法への入力データは, ラベル付きの連続値ベクトルの集合である. 各データ点は, n 次元のベクトル $\mathbf{x}_i = (x_i^1, x_i^2, \dots, x_i^n)$ とクラスのラベル $y_i \in \{0, 1\}$ の組 (\mathbf{x}_i, y_i) ($i = 1, \dots, N$) である. なお提案手法では, 各特徴に対し, **min-max 正規化** (min-max normalization) を用いてあらかじめ区間 $[0, 1]$ に正規化する. すなわち, 特徴ベクトルの各値 x_i^j を $(x_i^j - \min_{1 \leq l \leq N} x_l^j) / (\max_{1 \leq l \leq N} x_l^j - \min_{1 \leq l \leq N} x_l^j)$ へ変換するため, $\mathbf{x}_i \in [0, 1]^n$ となる.

2.1 離散化

離散化では, 連続値の特徴を要素にもつ各特徴ベクトルを離散化して分割することで, 2 値特徴ベクトルへ変換する. 変換された特徴はもとの空間内の区間に対応しており, それぞれの 2 値特徴ベクトルの要素は, その区間への所属を表す. この変換を各区間の幅を狭めていきながら繰り返し, クラスラベルが判別できたところで終了する. ただし, 分割回数の限度をパラメータとして与えることで, 過度の分割を防ぐ.

提案手法における分割 (partition) とは, 区間 $[0, 1]$ をいくつかの区間に等分することである. すなわち, 区間 $[0, 1]$ を k 区間に分割するとき, 各区間は $[0, 1/k], [1/k, 2/k], \dots, [(k-1)/k, k/k = 1]$ となる. 一般に, 関数 g を用いて, 各区間を $g(w_k) = [w_k/k, (w_k + 1)/k]$ と定める. ただし, $w_k \in \{0, 1, \dots, k-1\}$ である.

これらの定義を用いて, 特徴 $j \in \{1, 2, \dots, n\}$ を k 区間に分割した中で, $g(w_k)$ で表される区間に属するか否かによ

連絡先: 米田友花, 大阪大学産業科学研究所,

〒567-0047 大阪府茨木市美穂ヶ丘 8-1

yoneda@ar.sanken.osaka-u.ac.jp

Algorithm 1 離散化アルゴリズム

```
1:  $S \leftarrow \emptyset$ 
2: for all  $k = 2$  to  $\alpha_{\max}$  do
3:   generated by Eq.(1)
4:    $S \leftarrow S \cup Z^k$ 
5:   if  $\text{Bn}(Z^k) = 0$  or  $k = \alpha_{\max}$  then
6:     break
7:   end if
8: end for
9:  $S$  を出力する
```

って2値化した特徴ベクトルを $\mathbf{z}^{(j,w_k)}$ と表し、その特徴を (j, w_k) と書く。具体的には、特徴ベクトル $\mathbf{z}^{(j,w_k)}$ の各要素 $z_i^{(j,w_k)} (i = 1, \dots, N)$ は以下で与えられる。

$$z_i^{(j,w_k)} = \begin{cases} 1 & \text{if } x_i^j \in g(w_k), \\ 0 & \text{otherwise.} \end{cases}$$

そして、すべての特徴を k 分割したときに得られる2値特徴集合を Z^k とおく。すなわち、

$$Z^k = \{(j, w_k) \mid j \in \{1, 2, \dots, n\}, w_k \in \{0, 1, \dots, k-1\}\} \quad (1)$$

となる。このとき、各データ $\mathbf{x}_i = (x_i^1, x_i^2, \dots, x_i^n)$ は $\mathbf{z}_i = (z_i^{(1,0)}, \dots, z_i^{(1,k-1)}, z_i^{(2,0)}, \dots, z_i^{(2,k-1)}, \dots, z_i^{(n,0)}, \dots, z_i^{(n,k-1)})$ という2値ベクトルに変換される。 $k = 2, \dots, \alpha$ まで離散化したときの全特徴からなる集合を S とすると、 $S = Z^2 \cup Z^3 \cup \dots \cup Z^\alpha$ となる。ここで、 S を出力するのは、第二段階の特徴選択後に獲得する区間が、必要以上に細くなることを防ぐためである。

提案手法は、離散化で用いる分割数 k を2からひとつずつ増やしていく。分割数を増やすたびに、その時点で2値化されたデータが、与えられたクラスラベルに対して無矛盾かどうかを判定する。ここで、与えられた特徴集合がクラスラベルに対し無矛盾であるかどうかを、関数 Bn を用いて次のように判定する。ある特徴集合 I に対し $\text{Bn}(I) = 1$ であるとは、 I がクラスラベルに矛盾している状態に対応しており、あるデータ \mathbf{x}_i と $\mathbf{x}_{i'}$ が存在し、 $\mathbf{x}_i = \mathbf{x}_{i'}$ かつ $y_i \neq y_{i'}$ を満たすことをいう。それ以外の場合は $\text{Bn}(I) = 0$ とし、このとき I はクラスラベルに対して無矛盾 (consistent) である。

離散化アルゴリズムの疑似コードを、Algorithm 1 に示す。以下では、この動作について説明する。入力データは、 $[0, 1]$ に正規化された連続値を取る特徴ベクトルとクラスラベルの組からなる集合 $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$ であり、出力は離散化され2値になった特徴ベクトルの集合である。ただし、提案した離散化では、入力パラメータとして、離散化の分割回数の限度である α_{\max} を用いる。2分割から順番に分割数 k を増やしながら、各分割で得られた2値ベクトルが無矛盾かどうか、すなわち $\text{Bn}(Z^k) = 0$ かどうかを調べていく。 $k = 2$ から始まり、 $k = \alpha$ のときの特徴集合 Z^α に対し、 $\text{Bn}(Z^\alpha) = 0$ であれば分割をやめ、離散化後の全特徴からなる集合 $S = Z^2 \cup Z^3 \cup \dots \cup Z^\alpha$ を出力する。 $\text{Bn}(Z^k) = 1$ であれば k に $k+1$ を代入し、初めにもどる。しかし、 $k = \alpha_{\max}$ となっても、 $\text{Bn}(Z^k) = 1$ となる場合、そこで分割をやめ、 $S = Z^2 \cup Z^3 \cup \dots \cup Z^{\alpha_{\max}}$ を出力する。

2.2 Super-CWC による特徴選択

提案手法では、離散化で得た2値ベクトルに対して、Super-CWC アルゴリズムと呼ばれる特徴選択手法を適用することで、クラス分類に不必要な特徴を削除する。この節では Super-CWC アルゴリズムの概要について説明する。

Super-CWC アルゴリズム [6] は、2015年に提案された2値データに対する特徴選択手法である。Super-CWC アルゴリズムでは、 $\text{Bn}(I) = 0$ を満たす特徴の組合せ I を2分探索を用いて効率的に探すことによって、クラスラベルの判別において不要な特徴を削除し、必要最低限の特徴集合を選択することに成功している。

Super-CWC への入力データは、 n 次元の2値ベクトル $\mathbf{z}_i = (z_i^1, z_i^2, \dots, z_i^n) \in \{0, 1\}^n$ とクラスのラベル $y_i \in \{0, 1\}$ の組 $(\mathbf{z}_i, y_i) (i = 1, \dots, N)$ である。また、特徴ごとにデータを表す際は、 N 次元の2値ベクトル $\mathbf{z}^j = (z_1^j, z_2^j, \dots, z_N^j)^T$ 、クラスラベルは、 N 次元のベクトル $\mathbf{y} = (y_1, y_2, \dots, y_N)^T$ のように表現する。また、入力データのすべての特徴からなるインデックス集合を $S = \{j \mid j = 1, 2, \dots, n\}$ とする。出力された特徴集合を $I \subseteq S$ とし、 I を用いた特徴ベクトルの集合は $\{\mathbf{z}^j \mid j \in I\}$ と表される。このとき、Super-CWC は以下の問題を解く。

$$\underset{I \subseteq S, \text{Bn}(I)=0}{\text{argmim}} |I|.$$

つまり、出力は $\text{Bn}(I) = 0$ を満たす無矛盾な $I \subseteq S$ の中で、最小サイズのものである。

また、Super-CWC では、Symmetric Uncertainty [2] (以下では SU と書く) と呼ばれる、正規化された相互情報量を用いて特徴を選択する。SU は、以下の式で定義される。

$$\text{SU}(\mathbf{z}^j, \mathbf{y}) = 2 \frac{I(\mathbf{z}^j; \mathbf{y})}{H(\mathbf{z}^j) + H(\mathbf{y})}.$$

ここで、 $H(\mathbf{z}^j)$ 、 $H(\mathbf{y})$ はそれぞれ特徴ベクトル \mathbf{z}^j 、クラスラベル \mathbf{y} の情報エントロピー、 $I(\mathbf{z}^j; \mathbf{y})$ は相互情報量を表す。定義から、SU は $[0, 1]$ の値をとる。SU の小さい特徴の方が、大きい特徴よりもクラスラベルとの矛盾が大きいため、Super-CWC アルゴリズムでは SU の小さい特徴から、 S より除いていく。

3. 提案手法の動作例

この節では、提案手法の動作例を紹介する。表1に示した $n = 2$ 、 $N = 8$ のデータを入力とし、十分に離散化すると、各特徴を4分割することになる。動作例では簡単のため、出力 S から Z^2 と Z^3 を省略すると、変換した2値データは表2のようになる。その後、2値データを Super-CWC アルゴリズムに入力すると、5つの特徴へと削減され、表3のようになる。獲得した特徴集合から得られる2次元平面上での分類規則は、図1における黄色と青色に塗られた部分で表される。なお、図1には、表1の連続値データもプロットした。

4. 評価実験

この節では、提案手法を人工データに適用して性能評価を行い、その結果について考察する。まず4.1節では、実験に用いるデータ、実験設定、比較手法について述べる。その後、4.2節で実験結果を述べ、結果について考察する。

4.1 実験条件の設定

4.1.1 実験環境

実験に用いた計算機は、Intel Core i7-4790 CPU 3.60 GHz のプロセッサ、実験メモリ 8.00 GB、Windows 8 Enterprise 64

表1: 入力の連続値データ

	x^1	x^2	y
x_1	0.2	0.8	1
x_2	0.3	0.6	1
x_3	0.6	0.4	1
x_4	0.7	0.7	1
x_5	0.1	0.3	0
x_6	0.4	0.2	0
x_7	0.8	0.1	0
x_8	0.9	0.9	0

表2: 表1を離散化した後の2値データ

	$z^{(1,0)}$	$z^{(1,1)}$	$z^{(1,2)}$	$z^{(1,3)}$	$z^{(2,0)}$	$z^{(2,1)}$	$z^{(2,2)}$	$z^{(2,3)}$	y
z_1	1	0	0	0	0	0	0	1	1
z_2	0	1	0	0	0	0	1	0	1
z_3	0	0	1	0	0	1	0	0	1
z_4	0	0	1	0	0	0	1	0	1
z_5	1	0	0	0	0	1	0	0	0
z_6	0	1	0	0	1	0	0	0	0
z_7	0	0	0	1	1	0	0	0	0
z_8	0	0	0	1	0	0	0	1	0

表3: 表2から Super-CWC により獲得した2値データ

	$z^{(1,2)}$	$z^{(1,3)}$	$z^{(2,0)}$	$z^{(2,1)}$	$z^{(2,2)}$	y
z_1	0	0	0	0	0	1
z_2	0	0	0	0	1	1
z_3	1	0	0	1	0	1
z_4	1	0	0	0	1	1
z_5	0	0	0	1	0	0
z_6	0	0	1	0	0	0
z_7	0	1	1	0	0	0
z_8	0	1	0	0	0	0

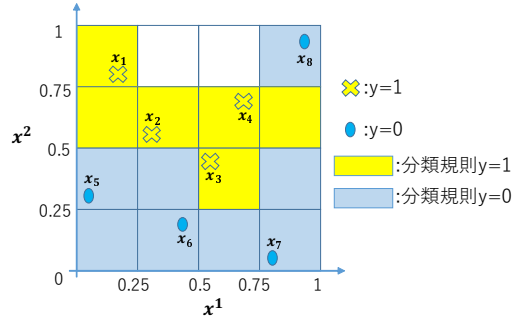


図1: 表3に対応する2次元平面上の領域

bit の OS が搭載された Project White 社製の exComputer である。提案手法および比較対象は、ともに Python 言語 (version 3.5.2) で記述されている。また、Super-CWC 本体は、著者によって公開されている Scala 言語での実装^{*1}を用いた。

4.1.2 データセット

本稿では、人工的に生成したデータ B, H を用いて評価実験を行った。人工データ B は、2つの正規分布からなるデータであり、それぞれの平均 μ と標準偏差 σ を $(\mu, \sigma) = (0, 1)$ と $(\mu, \sigma) = (2, 1)$ に設定した。また、人工データ H は、 $(\mu, \sigma) = (0.5, 0.125)$ の n 次元正規分布に従うデータを1つ目のクラス、領域 $[0, 1]^n$ の一様分布に従うデータを2つ目のクラスとし、それぞれのサイズが4:1の割合となるようにサンプリングした。

実験は、人工データ B, H に対し、分割回数の限度 $\alpha_{\max} = 10$ のとき、特徴数 n 、データ数 N を変化させ、10分割の交差検証法で分類精度を評価した。実験における評価指標は、提案手法と比較手法における閾値数、分類精度である。ここで、閾値数は、提案手法においては、得られた2値ベクトルの特徴数、すなわちベクトルの次元数とし、比較手法である決定木分析 CART においては、生成した木構造における葉以外の全ノード数とした。閾値数は、多いほど分類規則が複雑になっていると言えるため、規則の簡潔さの指標である。

4.1.3 比較手法

本実験では、提案手法を評価するための比較手法として、決定木分析 CART [1] を採用した。その理由は、CART が回帰問題を解くことができ、かつ分類規則を明示的に取得することが可能なためである。Python の機械学習ライブラリ scikit-learn [5] では、決定木分析の標準的手法として採用されており、本実験でもこの実装を用いた。

4.2 実験結果

人工データ B に対する実験では、特徴数を $n = 2$ に固定し、データ数を $N = 1,000, 10,000, 50,000, 100,000$ に変化させた場合と、データ数を $N = 10,000$ に固定し、特徴数を

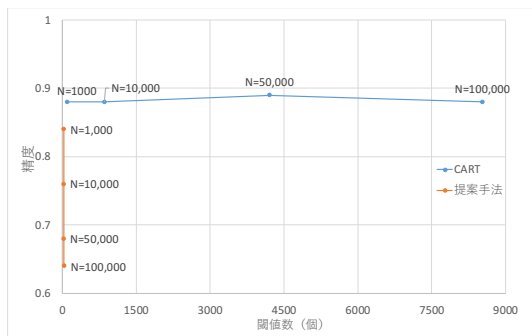
$n = 2, 3, 4, 5, 10, 50, 100$ に変化させた場合の閾値数、分類精度を測定した。結果を図2に示す。横軸は閾値数、縦軸は分類精度である。図2aにおいて、データ数を1,000個から100,000個に変化させたとき、CARTでは、分類精度を0.88付近で保ちつつ、閾値数が94.4から8530.3へ約90倍に増加したが、提案手法では、分類精度は0.84から0.64に低下したものの、閾値数は27.5から35.7へと1.3倍の増加で抑えられた。図2bにおいて、特徴数を2から100に変化させたとき、CARTでは、分類精度が0.88から0.98へ上昇し、閾値数が852.2から60.8へ1/14に減少したことに対し、提案手法の分類精度は、 $n = 2$ から3で0.76から0.85へ上昇し、 $n = 4$ では0.72へ低下したが、その後0.97へ上昇し、閾値数は、 $n = 2$ から3では28.6から41.7へ1.5倍に増加したが、その後 $n = 100$ では12へと1/3.5に減少した。

人工データ H に対する実験では、特徴数を $n = 2$ に固定して、データ数を $N = 1,000, 10,000, 50,000, 100,000$ に変化させた場合と、データ数を $N = 10,000$ に固定して、特徴数を $n = 2, 3, 4, 5, 10, 50, 100$ に変化させた場合の閾値数、分類精度を測定した。結果を図3に示す。横軸は閾値数、縦軸は分類精度である。図3aにおいては、データ数を1,000から100,000へ増やしたとき、CARTにおいては、分類精度を0.85付近で保ちつつ、閾値数が110.2から11233へ102倍に増加したことに対し、提案手法では、 $N = 100$ から10,000では分類精度が0.85から0.9に上昇したあと0.87に低下したものの、閾値数は32.8から54へ1.6倍の増加で抑えられた。図3bにおいては、特徴数を2から100へ増やしたとき、CARTにおいては、精度が0.86から0.95へ上昇し、閾値数は1082.2から154.9へ1/7に減少したことに対し、提案手法では、分類精度は0.9から0.8へ低下し、閾値数は51から21.5へ1/2.4に減少した。

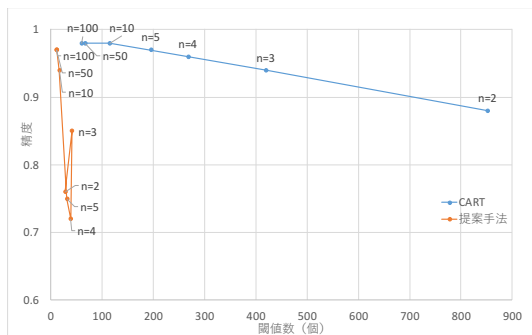
4.3 実験の考察

図2a, 3aからデータ数を増加させたとき、CARTでは閾値数が線形に増加するが、提案手法ではデータ数の増加にはあまり影響されず、CARTより大幅に少ない閾値数で分類を達成し

*1 <https://github.com/tkub/scwc/>



(a) データ数 N を変化させた場合



(b) 特徴数 n を変化させた場合

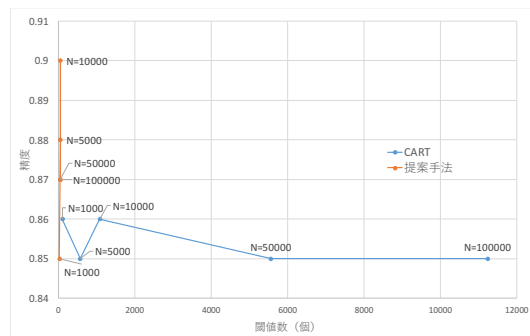
図2: データ B に対する実験結果

ている。これは、提案手法における離散化が、データの疎密によらずに分割し、かつ Super-CWC で不要な特徴を削除しているためであると考えられる。ただし、データ B において、データ数を増やしたとき、提案手法は分類精度が低下する。これは、提案手法がデータ数を増やしても閾値数があまり変わらないため、データ数が多くなるほど正確に分類することが難しくなっているためであると考えられる。それに対しデータ H では、 $N = 10,000$ までは分類精度が上がり、また $N = 1,000$ 以外では CART より精度が高いため、データ H のようなデータに対する分類は、提案手法の方が CART よりも適していると考えられる。

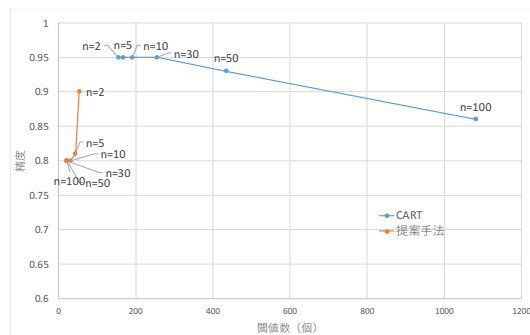
特徴数を増加させたとき、まず図 2b のデータ B に対する実験では、提案手法は CART より閾値数が少なく、分類精度も上昇している。これは、離散化後の区間数が分割回数に特徴数を乗じた値となるため、次元が大きくなると比較的少ない分割で無矛盾な特徴を獲得することができ、離散化後の区間の幅が狭くなりすぎないためであると考えられる。これに対して、図 3b のデータ H における実験では、提案手法は特徴数が増えるほど閾値数が減少し、同時に分類精度も低下している。これは、データ H においては、特徴数が増え区間数が増えると、訓練データにフィットしすぎている規則を学習しているためと考えられる。したがって、同じ特徴数やデータ数でもデータによって過学習となる分割回数が違うため、分割回数の限度を適切に与える必要がある。

5. まとめ

本稿では、連続値データに対し解釈可能性の高いクラス分類規則を構築するため、連続値データを離散化して 2 値データに変換したあと、既存の特徴選択手法である Super-CWC アルゴ



(a) データ数 N を変化させた場合



(b) 特徴数 n を変化させた場合

図3: データ H に対する実験結果

リズムを適用する手法を提案し、その性能の評価を行った。その結果、提案手法を用いることで、特に閾値数において、CART より少ない閾値数で同程度の分類精度を達成することができた。また、データ H の特徴数が比較的小さいときには、データ数が増加すると、低い閾値数の簡潔な分類規則で、テストデータに対し高い精度で分類することができた。

一方、今回の評価実験においては、人工データに対してのみの評価であった。そのため、今後は実データに対して、有益な分類規則を得ることができるかを確かめていく必要がある。

参考文献

- [1] L. Breiman, J. Friedman, C. J. Stone, and R. A. Olshen. *Classification and regression trees*. CRC press, 1984.
- [2] Witten I. H., E. Frank, and M. A. Hall. *Data Mining: Practical Machine Learning Tools and Techniques*. Morgan Kaufmann, 3 edition, 2011.
- [3] J. Han, M. Kamber, and J. Pei. *Data Mining: Concepts and Techniques*. Morgan Kaufmann, 3 edition, 2011.
- [4] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *Nature*, 521:436–444, 2015.
- [5] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [6] K. Shin, T. Kuboyama, T. Hashimoto, and D. Shepard. Super-CWC and super-LCC: Super fast feature selection algorithms. In *Proceedings of 2015 IEEE International Conference on Big Data*, pages 1–7, 2015.