ランダムフォレストにおけるノード数と木数の関係

On the trade-off between the number of nodes and the number of trees in Random Forest

熊野 颯^{*1} 阿久津 達也^{*2} So Kumano Tatsuya Akutsu

*1京都大学大学院情報学研究科 Grafuate School of Informatics, Kyoto University *²京都大学化学研究所 Institute for Chemical Research, Kyoto University

Expressibility of machine learning models has been extensively studied. For example, in a Neural Network, it is proved that the efficiency concerning the number of nodes is generated from the depth. On the other hand, it is not clear whether the efficiency exists in Random Forest. Therefore, in this research, we investigate whether the efficiency exists in Random Forest. We first show that Random Forest does not have the same kind of efficiency as Neural Network, and next we show that the efficiency concerning the number of nodes can be generated from the number of trees.

1. はじめに

封筒上の郵便番号の識別から店舗ごとの商品の需要予測、医 療診断に到るまで幅広い分野において機械学習は活用されて いる。いずれの応用においても、まず学習のモデル(例えば、 ニューラルネットワークなら層の数やニューロンの数、ランダ ムフォレストであれば木の本数や葉の数)を決定し、次に損失 関数を最小化するパラメータを発見する。この学習モデルの決 定、すなわち学習によって最適な関数の探索を行う関数族の決 定は非常に重要である。関数族が不必要に大きい場合は過剰な 計算コストや過学習の問題を引き起こし、逆に不十分な場合に は、適切な関数を得ることはできない。課題ごとに適切な関数 族を選択しなければならないのである。

こうした関数族の選択において有効な手がかりを得るべく、学 習モデルの表現能力に関する研究が古くから行われてきた。例え ば、深さ2のニューラルネットワークで任意のボレル可測関数を 近似することができることを主張する universal aproximation theorem が知られている。また、近年ではニューラルネット ワークは層の数によって表現能力が大きく変化することが分 かってきた。例えば、深さ2のニューラルネットワーク (SNN) は任意のボレル可測関数を近似することが可能であるが、必要 となるノード数が膨大となることがある。しかし、深いネット ワーク (DNN)を用いることで、SNNよりも少ないノード数 で同じ関数を表現できる場合がある。こうした DNN が持つ、 深さによって効率良く関数を表現することが可能となる性質は depth efficiency と呼ばれている。

ニューラルネットワークに対しては、表現能力の研究が盛ん におこなわれている一方で、ランダムフォレストにおいては表 現能力に関する研究は少ない。そこで、本研究ではランダム フォレストがニューラルネットワークにおける depth efficiency と同様にノード数に関する効率性を持つかについて考察する。

2. 関連研究

ニューラルネットワークの表現能力は、その深さに対して指 数的に増大すると考えられている。Montúfar らは区分線形関 数を活性化関数として用いた際に、ネットワークが表現する区

連絡先: 熊野颯,京都大学大学院情報学研究科, kumano@kuicr.kyoto-u.ac.jp 分線形関数の線形領域がネットワークの幅に対しては多項式の オーダーでしか増加しないのに対し、深さに対しては指数的に 増大することを示した [1]。 また、Raghu らは、ネットワー クの出力の軌道の長さの観点から、深さに対して表現能力が指 数的に増大することを示した [2]。

DNN は SNN よりも効率的に関数を表すことが出来ること を実際に示した研究も存在する。Telgarsky は、SNN では指数 オーダーのノードが必要だが、DNN では線形オーダーのノー ドで表現できる関数族の存在を示した [3]。Symanzki らは、 DNN において depth efficiency が生じる要因の一つとして入 力の周期性が関与していることを明らかにした [4]。Bengio ら は、Sum-product Network に関して DNN では SNN よりも ノード数の観点で効率的となる関数が存在することを示した [5]。

ランダムフォレストの表現能力に関する研究においては、 Mansour が入力の次元が *d*, ノード数 N の二分決定木の VC 次元の下界が Ω(N)、上界が O(N log *d*) であることを示した [6]。また、VC 次元が *d* の識別器を T 個使用するアンサンブ ル学習の VC 次元の上界は O(*dT* log (*dT*)) である [7]。Oshiro らは、データの密度を定義し密度と最適な木の本数の関係を実 験により確かめた [8]。これらの研究から、ランダムフォレスト は木の本数に対して、表現能力が劇的に増加するといったこと はないと思われる。しかし、これらの結果は depth efficiency と類似の性質の存在を否定する訳ではない。そこで本研究で は、ランダムフォレストにおけるノード数に関する効率性につ いて考察する。

3. ランダムフォレストの深さ

Bengio らは、ブースト木のニューラルネットワークにおけ る深さは3であると述べている。[9]。この結果から、ランダ ムフォレストは深さ3のニューラルネットワーク以下の表現能 力しか持たず、depth efficiency と同種の効率性は存在しない と考えられる。

しかし、この結果は 選言標準形のアナロジーとして述べら れたものである。一方、近年の DNN の depth efficiency に 関する結果の多くは区分線形関数を活性化関数として用いた ニューラルネットワークについて示されている。そこで、本研 究ではまず、ランダムフォレストは区分線形関数を活性化関数 として用いた深さ3のニューラルネットワーク以下の表現能力 しか持たないことを示す。ただし、ニューラルネットワークの 出力ノードは線形関数とする。以後、本研究ではランダムフォ レストを構成する決定木として、各ノードごとに入力空間をあ る変数に対して垂直二分割(他の変数に対して軸並行に分割) する二分決定木を考え、各決定木の出力はラベルのみであると する。つまり、ランダムフォレストが出力するラベルはそのラ ベルを出力する木の本数が最も多いラベルとなる。したがって 2 値分類問題におけるランダムフォレストはすべて奇数本の決 定木から構成される。また、*c* でクラス全体を表し、*H* を

$$H(x) = \begin{cases} 1 \ (x \ge 0) \\ 0 \ (x < 0) \end{cases}$$
(1)

と定義する。また、ニューラルネットワークは出力ノードの値 が最も大きいクラスに入力を分類するものとする。

補題 1. n / -ドの決定木 $DT : \mathbb{R}^d \to c$ は Hを活性化関数に 用いた O(n) / -ド、深さ 3 の NN で表現することができる。

Proof. NN の深さ1のノードとして DT の各エッジが成立す る場合に1、それ以外は0を出力するノードを作成する。す なわち、DT のエッジが $f_1 \ge a$ ならば $H(f_1 - a)$ を作成し、 $f_1 > a$ ならば1 - H(a - x) を作成する。(1 はこのノードの 出力を受け取るノードのバイアスとして与えられる)

次に深さ2のノードとして DT の各葉ノードを表現するノー ドを作成する。各ノードは DT における葉ノードから根まで の枝に対応するノードの和から、DT における葉ノードの深 さ-0.5 を引いた値を活性化関数の入力として受け取る。最後 に、深さ3の出力ノードとして各クラスに対応するノードを 作成し、各クラスに属する葉ノードと対応するノードの出力を 入力として受け取る。この NN が DT を表現することは明ら かである。実際、DT において入力が到達する葉ノードに対応 する深さ2のノードのみが1を出力し、他の深さ2のノード は0を出力する(少なくとも1つの深さ1のノードの出力は0 であるため)

例えば、下図のニューラルネットワークは左図の決定木を模 倣する。



図 1: エッジの数字は重み、ノードの数字はバイアスを表す。 例えば、深さ 1 の左端のノードは *H*(*x*₁ - 1) を表す。

この補題を用いて、多くの区分線形間数に対し、ランダムフォレストはその関数を活性化関数として用いた深さ3のニューラルネットワーク以下の表現能力しか持たないことが示せる。以降、 $S \subset \mathbb{R}^d$ は有限集合であるとする。

定理 2. n ノードの決定木 $DT: S \rightarrow c$ は、非有界な領域に おける傾きが異なる区分線形関数 $g: R \rightarrow R$ を活性化関数と して用いた O(n) ノードの深さ3の NN で表現することがで きる。

Proof. S は有限集合であるから、任意の S の要素が DT の識 別境界上に存在しないと仮定してよい。f は、 $a \neq b$ を満たすあ る定数 a, b に対して、 $\lim_{x\to\infty} f(x) = a$, $\lim_{x\to-\infty} f(x) = b$ を満たすとする。この時、f を用いて ($x \neq 0$ において) H を 表現することができる。実際、

$$H(x) = \lim_{\epsilon \to 0} \left(f(\frac{x}{\epsilon}) - b\right) \frac{1}{a - b} = \frac{1}{a - b} \lim_{\epsilon \to 0} f(\frac{x}{\epsilon}) - \frac{b}{a - b}$$
(2)

である。さらに、gを非有界な領域においてax+b, a'x+b'と表される区分線形関数であるとする。この時、g(x+1) - g(x)は非有界な領域においてそれぞれ、a(x+1)+b-ax-b=a, a'となるため、gを用いて、Hを活性化関数として用いた NNを表現することができる(図 2)。したがって、補題1から非有界な領域における傾きが異なる区分線形関数を活性化関数として用いたO(n)ノード、深さ3の NN で DT を表現することができる。



図 2: g として ReLU 関数を用いた場合、活性化関数として H を用いたネットワーク (左図) は右図のネットワークで表現される

定理 3. n / -ドのランダムフォレスト $RF : S \rightarrow c$ は非有界 な領域における傾きが異なる区分線形関数を活性化関数として 用いた O(n) / -ドの深さ3の NN で表現することができる。

Proof. ランダムフォレスト を構成するそれぞれの木に対して、 対応する NN を作成する。深さ 2 の各ノードはそれぞれ対応 する決定木の葉ノードに入力が到達したときに 1、それ以外の 場合には 0 を出力する。したがって、これらの NN の出力ノー ドを同一のノードとして見做すことで、RF に対応する NN を 得ることができる。□

これらの結果から、ランダムフォレストは深さ3のニュー ラルネットワーク以下の表現能力しか持たないこと、すなわち depth efficiency と全く同種のノードに関する効率性は存在し ないことが分かる。

4. 木の本数と表現能力

4.1 ノード数の下界

前節から、ランダムフォレストにおいてはニューラルネット ワークに対応する深さは固定であること、すなわち DNN と 同様の効率性は存在しないことが分かった。次に、木の本数 によって表現能力がどう変化するかの検証を行う。その為に、 $n \neq O(n)$ /ードのランダムフォレスト を T 本の木からなる ランダムフォレストで表現する際に必要となるノード数を求 める。 入力全体を X で表す。 $x \in X$ とランダムフォレストを構成 する木の各葉ノード leaf に対して、xのラベルを c(x)、leaf に割り当てられているラベルを c(leaf)、x が leaf に到達す るか否かを leaf(x) で表す。

補題 4. 次の性質を満たす集合 X を考える。この時、M = |X|とおくと X を T 本の木からなるランダムフォレストで表現するには $\Omega(M^{\frac{2}{T+1}})$ のノードが必要である。

X を識別する任意の T 本の木からなるランダムフォレスト に対し、下記が成立する。

$$\forall x_1, x_2 \in X, c(x_1) = c(x_2) = c \rightarrow \\ |\{leaf|c(leaf) = c, leaf(x_1), leaf(x_2)\}| < \frac{T+1}{2}$$
(3)

Proof. 上記の性質は任意の同一ラベルを持つ異なる 2 点は、 X を認識する任意のランダムフォストにおける葉ノードのう ち半数以上に同時に正しく認識されることはないということで ある。 $L = \{leaf|c(leaf) = 1\}, l = |L|$ とする。つまり、lは ラベル 1 が割り当てられている葉の数である。同様に L'でラ ベル 0 が割り当てられている葉全体、l'でラベル 0 が割り当て られている葉の数を表す。また、 M_1 で $|\{x \in X | c(x) = 1\}|$ 、 M_0 で $|\{x \in X | c(x) = 0\}|$ を表す。すなわち、 M_i はクラス iのデータの個数を表す。

c(x) = 1を満たす x は少なくとも $\frac{T+1}{2}$ 以上の L の要素に対して leaf(x) が成立する必要がある。一方、(3) より leaf(x) が成立する L の要素のうち、どの $\frac{T+1}{2}$ 個の葉の組み合わせも他の c(x) = 1を満たす x に対して同時に leaf(x)を満たすことはない。したがって、各 $x \in M_1$ に対し他の $x' \in M_1$ では同時に 1 とならない $\frac{T+1}{2}$ 個の L に属する葉の組み合わせが存在するから、

$$\binom{l}{\frac{T+1}{2}} \ge M_1 \tag{4}$$

が成立する。また、ラベル0の場合も同様にして

$$\binom{l'}{\frac{T+1}{2}} \ge M_0 \tag{5}$$

が成立する。したがって、

$$l^{\frac{T+1}{2}} \ge M_1 \tag{6}$$

$$l'^{\frac{T+1}{2}} \ge M_0 \tag{7}$$

が成立する。したがって、このランダムフォレストの葉ノードの数は

$$l + l' \ge M_1 \frac{2}{T+1} + M_0 \frac{2}{T+1} \ge M \frac{2}{T+1}$$
(8)

となる。したがって、必要なノード数の下界は

$$\Omega(M^{\frac{2}{T+1}}) \tag{9}$$

この補題を用いて、 $n \pm O(n)$ ノードのランダムフォレスト を T 本の木からなるランダムフォレストで表現する際に必要 なノード数の下界を求めることができる。下界は次のように なる。 定理 5. $n \neq O(n)$ ノードのランダムフォレスト *RF* : $\{0,1\}^n \rightarrow \{0,1\}$ を *T*本の木からなるランダムフォレスト で表現する際に必要なノード数の下界は $\Omega((\frac{2^n}{\sqrt{T+1}})^{\frac{2}{T+1}})$ である。

Proof. 任意のT (T < n) に対し、補題4の(3)の性質を満たし、 $n \neq O(n)$ ノードのランダムフォレストで識別できるXが存在することを示せばよい。各次元の要素の和が $\frac{n+1}{2}$ (ラベル1を持つ)もしくは $\frac{n+1}{2} - 1$ (ラベル0を持つ)となるn次元の0.1ベクトル全体をXと置く。

X が任意の T (T < n) に対し、(3) の性質を満たすことを 示す。X を識別する T (T < n) 本の木からなるランダムフォ レストが存在するとする。X が (3) の性質を満たさないと仮 定する。このとき、 $c(x_1) = c(x_2) = c$ を満たす、ある x_1, x_2 が存在して、これらに対し $leaf(x_1), leaf(x_2), c(leaf) = c$ を 満たす leaf が $\frac{T+1}{2}$ 個以上存在する。以降では一般性を失う ことなく c = 1 と置く。

一方、 $x_1 \neq x_2$ であるから x_1 で1、 x_2 で0となる要素が少 なくとも1つは存在し、この要素は $x_1 \ge x_2$ の両方が同じ葉 ノードに到達する木の出力に影響を与えない。したがって、 x_1 においてこの要素を反転させた x_1' はこのランダムフォレス トにおいてラベル1と識別される。しかし、 x_1' は要素の和が $\frac{n+1}{2} - 1$ であるから、X に属し、ラベル0を持つはずである。 これは、X が T 本の木からなるランダムフォレストによって 識別されるということに矛盾する。したがって、X は (3) の 性質を満たす。

Xが $n \neq O(n)$ ノードのランダムフォレストによって識別 されるということは容易に示すことができる。それぞれの木に おいて1つの特徴量を評価し、1ならばラベル1、0ならばラ ベル0を出力すれば良い。また、

$$|X| = 2\binom{n}{\frac{n+1}{2}} \ge \frac{2^n}{\sqrt{n}} \tag{10}$$

であるから、必要なノード数の下界は

$$\Omega(\left(\frac{2^n}{\sqrt{n}}\right)^{\frac{2}{T+1}}) \tag{11}$$

この定理から同じ関数を表現する際に、木の本数が n 本で あるランダムフォレストと比較して T 本の木から構成される ランダムフォレストでは多量のノードを必要とすることがある こと、すなわちランダムフォレストは木の本数に対してノード 数に関する効率性を持つことが示された。

5. 結論と展望

本研究では、ランダムフォレストにおいて DNN の depth efficiency のようなノード数に関する効率性が存在するかとい うことについて考察を行った。まず、多くの区分線形関数に対 し、ランダムフォレストはその関数を活性化関数として用いた 3層のニューラルネットワーク以下の表現能力しか持たないこ とを示した。このことから、ランダムフォレストにおいては ニューラルネットワークと全く同種の効率性は存在しないこと が示された。次に、n本の木から構成されるランダムフォレス トを T 本の木から構成されるランダムフォレストで表現する 際には、T が n よりも十分に小さい場合には多くのノードが 必要となることを示した。このことから、木の本数に関しては ノード数に関する効率性が存在することが示された。今後の課

 \square

題としては、定理5の下界の改善、n本の木から構成されらラ ンダムフォレストをT本の木から構成されるランダムフォレ ストで表現する際に必要となるノード数の上界の導出などが考 えられる。

参考文献

- Guido Montufar, F., et al. "On the number of linear regions of deep neural networks." Advances in neural information processing systems. 2014. p. 2924-2932.
- [2] Maithra Raghu, et al. "On the expressive power of deep neural networks." arXiv preprint arXiv:1606.05336 (2016).
- [3] Matus Telgarsky. "Representation benefits of deep feedforward networks." arXiv preprint arXiv:1509.08101 (2015).
- [4] Lech Szymanski, and Brendan McCane. "Deep networks are effective encoders of periodicity." IEEE Transactions on Neural Networks and Learning Systems 25.10 (2014): 1816-1827.
- [5] Olivier Delalleau, and Yoshua Bengio. "Shallow vs. deep sum-product networks." Advances in Neural Information Processing Systems. 2011. p. 666-674
- [6] Yishay Mansour. "Pessimistic decision tree pruning based on tree size." In Press of Proc. 14th International Conference on Machine Learning. 1977. p.195–201.
- Shalev-Shwartz, Shai, and Shai Ben-David. Understanding machine learning: From theory to algorithms. Cambridge university press, 2014. p139
- [8] Mayumi Thais Oshiro, Pedro Santoro Perez, and José Augusto Baranauskas. "How many trees in a random forest?." International Workshop on Machine Learning and Data Mining in Pattern Recognition. Springer, Berlin, Heidelberg, 2012. p. 154-168
- [9] Yoshua Bengio, Olivier Delalleau, and Clarence Simard. "Decision trees do not generalize to new variations." Computational Intelligence 26.4 (2010): 449-467.