

深層ディジションフォレストによる疎密混合データの解析

大北 剛^{1,a)} 井上 創造^{1,b)}

概要：ランダムフォレストは、弱学習器によりアンサンブル学習を実現するが、この機構を用いて目的変数を設定すれば説明変数の重要度を求めるタイプの学習に用いることができる。深層ディジションフォレストは、決定木におけるユニークなルーチングをベルヌーイ確率変数により確率的に表現し、微分可能とする。このことにより、決定木のディープ化においてよく見られる局所的でグリーディな最適化ではなく、グローバルな最適化を行なう版であり、精度が上昇したことが報告されている。実験では、介護施設におけるインシデント報告というマルチモーダルなデータを用いて、なぜ事故が起こったのかという原因究明の解析を深層ディジションフォレストとランダムフォレストを比較する。

1. イントロダクション

ランダムフォレスト [1] は、弱学習器によりアンサンブル学習を実現するが、この機構を用いて目的変数を設定すれば説明変数の重要度を求めるタイプの学習に用いることができる。深層ディジションフォレスト [5] は、決定木におけるユニークなルーチングをベルヌーイ確率変数により確率的に表現し、微分可能とする。このことにより、決定木のディープ化においてよく見られる局所的でグリーディな最適化ではなく、グローバルな最適化を行なう版であり、精度が上昇したことが報告されている。実験では、介護施設におけるインシデント報告という疎密の混合するデータを用いて、なぜ事故が起こったのかという原因究明の解析を深層ディジションフォレストとランダムフォレストを比較する。

2. 深層ディジションフォレスト

Kontschieder[5] にしたがって、深層ディジションフォレストを記述する。入力を \mathcal{X} 、出力を \mathcal{Y} とする分類問題を考える。決定木は、分割ノード（内部ノード）と推定ノード（葉ノード）よりなる木構造の分類器である。分割ノードと推定ノードをそれぞれ \mathcal{N} 、 \mathcal{L} によりインデックスづける。推定ノード $\ell \in \mathcal{L}$ は \mathcal{Y} 上の確率分布 π_ℓ をもつ。分割ノード $n \in \mathcal{N}$ は θ によりパラメータ化された決定関数 $d_n(\cdot; \theta) : \mathcal{X} \rightarrow [0, 1]$ を割り当てられる。

決定木のノードを θ によりパラメータ化したとする。このとき、木 T からのサンプル x に対する確率分布の推定は

$$P_T[y|x, \theta, x] = \sum_{\ell \in \mathcal{L}} \pi_{\ell_y} \mu_\ell(x|\theta) \quad (1)$$

と表現される。なお、 $\pi = (\pi_\ell)_{\ell \in \mathcal{L}}$ とし、 π_{ℓ_y} をサンプルが葉ノード ℓ に到達してクラス y となる確率とする。また、 $\mu_\ell(x|\theta)$ をルーチング関数と呼び、

$$\mu_\ell(x|\theta) = \prod_{n \in \mathcal{N}} d_n(x; \theta)^{1_{\ell \leftarrow n}} \bar{d}_n(x; \theta)^{1_{n \rightarrow \ell}} \quad (2)$$

と表現されるとする。なお、 $\bar{d}_n(x; \theta) = 1 - d_n(x; \theta)$ とし、1 を指示関数とする。

決定ノードにおいて確率的ルーチングを行なうことが決定木と異なる点である。 σ をシグモイド関数とし、 $f_n(\cdot; \theta)$ を \mathcal{X} 、 θ を変数とする関数とする。すると、確率的ルーチング $d_n(x; \theta)$ を

$$d_n(x; \theta) = \sigma(f_n(x; \theta)) \quad (3)$$

と表現できる。

さらに、フォレストを k 本の決定木のアンサンブル $\mathcal{F} (= \{T_1, \dots, T_k\})$ とする。すると、入力 x に対する確率分布の推定は、 \mathcal{F} における各々の木の出力を平均化することで得られる。つまり、

$$P_{\mathcal{F}}[y|x] = \frac{1}{k} \sum_{h=1}^k P_{T_h}[y|x] \quad (4)$$

と表現される。

データ集合 $T \subset \mathcal{X} \times \mathcal{Y}$ に対して、ログ損失関数 $L(\theta, \pi; x, y)$ を用いて、

$$R(\theta, \pi; T) = \frac{1}{|T|} \sum_{(x,y) \in T} L(\theta, \pi; x, y) \quad (5)$$

という経験損失を最小化する。なお、ログ損失関数

¹ 九州工業大学
〒804-8580, 北九州市戸畑区仙水町 1-1
^{a)} tsuyoshi.okita@gmail.com
^{b)} sozo@mns.kyutech.ac.jp

$L(\theta, \pi; x, y)$ をトレーニング集合 $(x, y) \in T$ に対するロ
グ損失項としており、

$$L(\theta, \pi; x, y) = -\log(P_T[y|x, \theta, \pi]) \quad (6)$$

として定義される。この最適化は、決定ノードの学習と葉
ノードの学習という二つの要素を学習するものとなる。各
エポックにおいて、決定ノードの学習（ミニバッチ分）と葉
ノードの学習を交互に両方行なう。

決定ノードの学習は以下の通りである。確率的勾配降下法
を用いて θ に対して経験損失最小化法を適用する。すると、

$$\theta^{(t+1)} = \theta^{(t)} - \eta \frac{\partial R}{\partial \theta}(\theta^{(t)}, \pi; B) \quad (7)$$

となる。 θ に関する損失関数 L の勾配は以下のように表現
され、

$$\frac{\partial L}{\partial \theta}(\theta, \pi; x, y) = \sum_{n \in N} \frac{\partial L(\theta, \pi; x, y)}{\partial f_n(x; \theta)} \frac{\partial f_n(x; \theta)}{\partial \theta} \quad (8)$$

決定木に依存する f_n に関する損失関数 L の勾配は

$$\frac{\partial L(\theta, \pi; x, y)}{\partial f_n(x; \theta)} = d_n(x; \theta) A_{n_r} - \bar{d}_n(x; \theta) A_{n_\ell} \quad (9)$$

と表現される。なお、 A_m を一般的なノード $m \in \mathcal{N}$ に対して

$$A_m = \frac{\sum_{\ell \in L_m} \pi_{\ell_y} \mu_\ell(x|\theta)}{P_T[y|x, \theta, \pi]} \quad (10)$$

と定義されるとする。なお、中間層のない版を用いている。

次に葉ノードの学習は式 5 を θ を固定して、 π に関して
最小値を求める。これも Kotschieder[5] にしたがっている。
つまり、

$$\min_{\pi} R(\theta, \pi; T) \quad (11)$$

となる。これは

$$\pi_{\ell_y}^{(t+1)} = \frac{1}{Z_l^{(t)}} \sum_{(x, y') \in T} \frac{\mathbf{1}_{y=y'} \pi_{\ell_y}^{(t)} \mu_\ell(x|\theta)}{P_T[y|x, \theta, \pi]} \quad (12)$$

を逐次的に解くこととなる。なお、Kotschieder[5] が、この
値が収束することを証明している。

以上は、一つの深層決定木の学習であったが、これを k 木
のフォレスト F にする。これも Kotschieder[5] にしたがう。
 k 木のフォレストは θ を共有することができるが、異なる
決定関数と構造、葉ノードにおける確率分布の推定 π が独
立でなければならない。このためには π に問題はなく（フォ
レスト F が独立した π をもっていれば、エポックが進んで
も確率分布の推定 π は独立のままである）、 θ は各ミニバ
ッチにおいてランダムに選択すればよい。

3. 実装と重要度の計算

実装に際しても Kotschieder[5] にしたがったが、重要度
に関する記述はなかったため、これに関しては補った。

決定木の決定関数 d_n を $f_n(x; \theta_n) = \theta_n^T x$ として k 本の

木を並列に並べた層を決定木の層とした。したがって、全体
のネットワーク構造は、入力層、ReLU を活性化関数とする
全結合層、シグモイドを活性化関数とする全結合層、そして、
決定木の層からなる。

ランダムフォレストにおける重要度は、フォレストのそ
れぞれの木において、特定の変数に対して、その項目の中身
をシャッフルした場合にそれらによりどれだけ推定値が悪
くなったかを測定した値である [1]。この値が大きくなれば悪
くなれば重要度が高いと考える。それぞれの変数ごとにフォ
レスト内のすべての木上でこの重要度を平均する。深層ディ
ジションフォレストにおいても同様のことを行なったが、
特定の変数は決定木のノードに縛られていないため、まず、
特定の変数に関してシャッフルしたテスト集合を用意し、
単純に k 個の深層ディジションフォレスト内の各々の木に
対して、クラスごとに精度がどう落ちるかを測定した。つ
まり、シャッフルしたテスト集合に対して

$$A_\ell = \frac{\pi_{\ell_y} \mu_\ell(x|\theta)}{P_T[y|x, \theta, \pi]} \quad (13)$$

という量がどう変化するかをクラス毎に測定し、シャッ
フルしないテスト集合に対して得られる値との差の絶対値を
クラス上で平均したものを重要度と考えた。なお、シャッ
フルしたテスト集合はそれぞれのクラスに対し N 種類用意
し、それらに対して上記を測定した。その上で平均値を取
り、これとシャッフルしないテスト集合に対して得られる
値との差の絶対値を考えた。以下の実験においては $N = 10$
を用いた。

4. 実験結果

4.1 データ

実験においては、介護施設におけるインシデントデー
タを用いた [6]。項目は、発生時間、年齢、介護度、認知度、利用
区分、事故歴、事故の種類、発生場所、目的、詳細、重症度、診
断よりなる。詳細のみが自然言語の自由形式の記述で行な
われる。その他は離散値による記述もしくはいくつかのカ
テゴリに属する名詞として記述される。

「詳細」項目における日本語による記述は、1 件は複数行
を含むことがあり、件数の数は 5, 190 件となり、265, 397 語
を含む。また、数字、漢字、アルファベット、許容範囲の記号
(., 「」など) 以外の文字をチェックした。半角コードと全
角コードが混在し、一方、表 1 に示すような記号が存在する
ことがわかった。このためこれらは全角コードに統一して
正規化した。なお、これら以外にも表記の揺れは存在したが、
これらに関する正規化は行っていない。

データをトレーニング集合 4,106 タプル、ディベロッ
プメント集合 456 タプル、テスト集合 507 タプルに固定した
やり方で分割した。データ 5 は利用区分を目的変数として、
その他を説明変数とした。同様にデータ 7、データ 10、デー
タ 11 は事故の種類、重症度、診断をそれぞれ目的変数とし、

その他を説明変数とした。分類におけるクラス数は表 2 に示した通りである。

4.2 アルゴリズムに関する詳細

実験 1 はデータの大きさが $m = 11$ で、すべての項目を離散値の ID に変換した。実験 2 はデータの大きさが $m = 12$ で、上記の実験 1 で用いたデータに自然言語による表現を追加した。

ランダムフォレストは sklearn^{*1}を用いた。一方、深層ディビジョンフォレストは Lasagne^{*2}/Theano^[8]^{*3}ベースの実装^{*4}をベースとして用いたが変更を加えている。

ランダムフォレストにおいては、トレーニング集合を用いてクラス分類器を構築し、テスト集合に対して性能(精度と OOB スコア)を測定した。つまり、ランダムフォレストにおいては、ディベロップメント集合を用いたパラメータなどのチューニングは行っていないが、深層ディビジョンフォレストと条件を同じにするため、ディベロップメン

ト集合を設定した。したがって、トレーニング集合の大きさを若干犠牲にしている。

一方、深層ディビジョンフォレストにおいては、ディベロップメント集合を用いて、ハイパーパラメータのチューニングを行なった。チューニングは以下のパラメータを用いたベイジアン最適化 [7] による。関係するパラメータは以下の通りであり、各々の選択は経験により行なった。したがって、総当りでチューニングを行なったわけではない。

depth	[2, 3, 4, 5]
n__estimators	[10, 20, 30, 40]
n__hidden	[64, 128, 256]
learning__rate	[0.1, 0.01, 0.001, 0.0001]
num__epochs	[500, 2000]
batch__size	[100, 200]
momentum	[0.0, 0.8]
dropout	[0.0, 0.2, 0.25, 0.3, 0.4, 0.8]
update	[<i>adagrad</i> , <i>sgd</i> , <i>adam</i> , <i>RMSprop</i>]

表 3 ハイパーパラメータのチューニングで用いた探索範囲

シンボル	正規化後の表現
リットル記号 (半角)	リットル
『	度
』	「
...	」
cm (半角と全角)/センチ	(消去)
mm (半角と全角)/ミリ	センチ
右手大 (囲み文字 4) 指の痛み	センチ
皮膚剥離を (囲み文字 1) 発見する	ミリ
本人様による事後報告	右手大四指の痛み
2 2 の表皮剥離	皮膚剥離を発見する
第 II 指の爪	本人様による事後報告
ベット トイレ間	2 x 2 の表皮剥離
て転倒 外傷なし	第二指の爪
1 リットル 1.25 リットル	ベットとトイレ間
1 1 : 4 5 下痢が続いており	て転倒したが外傷なし
S P 2/SpO2	1 リットルから 1.25 リットル
Dr	1 1 : 4 5 から下痢が続いており
NS, Ns	S P O 2
NC, Nc	医師
	看護師
	ナースコール

表 1 インシデントデータにおける記号データに関する正規化。

データ 5	15 クラス
データ 7	16 クラス
データ 10	4 クラス
データ 11	15 クラス

表 2 四種の分類タスクにおける目的変数のクラスの数

ディベロップメント集合に対して高い精度をあげたモデルに対し、固定テスト集合に対して、精度を測定した。損失関数はクロスエントロピーを用いた。

形態素解析に関しては JUMAN を用いた [3]。

4.3 実行環境

両者ともインテルの CPU i7-7920HQ 上で実行し、さらに深層ディビジョンフォレストに関しての実行は NVIDIA GTX1080 上で行なった。

4.4 実験 1 の詳細

上述したように実験 1 は離散値のみのデータからなる。実験 1 の実験結果のテスト集合に対する実験結果は以下の通りである。

	ランダムフォレスト		深層ディビジョンフォレスト	
	テスト	ディベロップ	テスト	ディベロップ
データ 5	0.619	0.678	0.635	0.560
データ 7	0.688	0.735	0.692	0.699
データ 10	0.714	0.789	0.716	0.663
データ 11	0.481	0.548	0.487	0.458

表 4 テスト集合/ディベロップメント集合に対する実験結果 (精度) を示す。

いずれの項目においても、深層ディビジョンフォレストはランダムフォレストの精度をそれぞれ 2.6%, 0.6%, 0.3%, 1.2% 上回った。興味深いことは深層ディビジョンフォレストはディベロップメント集合に対し、あまり良い性能は示していないが、テスト集合では良い性能を示している。これは汎化性能が良いことを示していると思われる。なお、前述したがランダムフォレストはディベロップメント集合をト

*1 <http://scikit-learn.org>

*2 <https://lasagne.readthedocs.io/en/latest/>

*3 <http://deeplearning.net/software/theano/>

*4 <https://github.com/SkidanovAlex/ShallowNeuralDecisionTree>

レーニングにおいて使っていない。

チューニングにおける最適化アルゴリズムの選択を眺めると、データ 5、データ 7、データ 11 は adagrad [2] の性能が良く、データ 10 は adam [4] の性能が良かった。また、データ 10 においては誤差サーフェイスが起伏に富んでいるため、学習率の小さいものが良い性能を上げた。

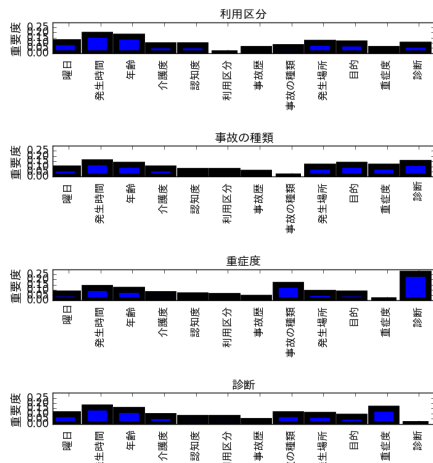


図 1 リレーショナルデータのみに対するランダムフォレストによる重要度のプロット

4.5 実験 2 の詳細

実験 2 の実験結果のテスト集合に対する実験結果は以下の通りである。自由文を形態素解析 [3] を行ない、ストップワードを除くと、7,745 項目を得た。各々の項目に対し、頻度を素性としたものを、実験 1 と同様にトレーニング集合、ディベロップメント集合、テスト集合に分割した。これを元に、同様にランダムフォレスト、深層ディシジョンフォレストで分類タスクを解かせた。

ランダムフォレストは、表 4 と表 5 を比較して明らかのように、データ 5, 7, 10, 11 のいずれにおいても混在データの方が 10%, 3%, 2%, 22% 精度が良い。自然言語という記述を 7,745 項目という疎な形のデータが密な離散データの精度を補っている。同様の傾向が深層ディシジョンフォレストでも見られる。

	ランダムフォレスト		深層ディシジョンフォレスト	
	テスト	ディベロップ	テスト	ディベロップ
データ 5	0.690	0.662	0.694	0.660
データ 7	0.736	0.728	0.742	0.722
データ 10	0.715	0.704	0.730	0.707
データ 11	0.570	0.561	0.581	0.570

表 5 マルチモーダルデータに対するテスト集合/ディベロップメント集合に対する実験結果 (精度) を示す。

5. 結論

本論文においては、分類タスクにおけるランダムフォレストを深層化した深層ディシジョンフォレストをランダムフォレストと比較した。最初の実験においては離散値のみからなるデータを用いた。すべての項目においてわずかながら深層ディシジョンフォレストの精度がランダムフォレストを上回った。二つ目の実験においては自然言語の記述からストップワードを除いた語をという疎なデータを最初の実験の密な離散値のデータに加えた。ここにおいても同様にすべての項目においてわずかながら深層ディシジョンフォレストの精度がランダムフォレストを上回った。

参考文献

- [1] L. Breiman, Random forests, Machine Learning, 45:5-32, 2001.
- [2] John Duchi, Elad Hazan, Yoram Singer. Adaptive Subgradient Methods for Online Learning and Stochastic Optimization, Journal of Machine Learning Research 12, page 2121-2159, 2011.
- [3] 河原大輔, 黒橋禎夫: 自動構築した大規模フレームに基づく構文・格解析の統合的確率モデル, 自然言語処理, Vol.14, No.4, pp.67-81, (2007.7).
- [4] Diederik P. Kingma, Jimmy Ba, Adam: A Method for Stochastic Optimization, the 3rd International Conference for Learning Representations, San Diego, 2015
- [5] Peter Kotschieder, Madalina Fiterau, Antonio Criminisi, Samuel Rota Bulo, Deep Neural Decision Forests, The IEEE International Conference on Computer Vision (ICCV), pp. 1467-1475, 2015.
- [6] 峯崎 智裕, 井上 創造, "介護サービス向上に向けた介護事故事例テキストの分析", マルチメディア, 分散, 協調とモバイル (DICOMO2016) シンポジウム, pp. 1663-1669, 2016.
- [7] Jasper Snoek, Hugo Larochelle and Ryan Prescott Adams. Practical Bayesian Optimization of Machine Learning Algorithms. Neural Information Processing Systems, 2012.
- [8] Theano Development Team, "Theano: A Python framework for fast computation of mathematical expressions", arXiv e-prints, abs/1605.02688, 2016.