

SHAP値の考え方を理解する

木構造編

概要

- ▶ SHAPライブラリの計算の仕組みを知るために、木構造を対象とした以下の論文を読んだ

Consistent Individualized Feature Attribution for Tree Ensembles

<https://arxiv.org/pdf/1802.03888.pdf>

- ▶ まだ理解できていない部分はあるが、重要そうな考え方は理解できたので本資料を作成してみた



SHAP値って何？

- ▶ SHapley Additive exPlanationsの略
- ▶ ゲーム理論で用いられる考え方をベースにしている
- ▶ 予測の時にどの説明変数が影響したかを示すことができる (Feature importances_の予測バージョン)



例えば、PTRATIOが15.3であることやLSTATが4.98であることは
予測値を増やす要因となっていることが読み取れる

<https://github.com/slundberg/shap> より

SHAPの使い方

- ▶ Pip install shapから簡単に使える。以下のURL等を参照

公式ドキュメント

<https://github.com/slundberg/shap>

<https://shap.readthedocs.io/en/latest/#>

自作ブログ

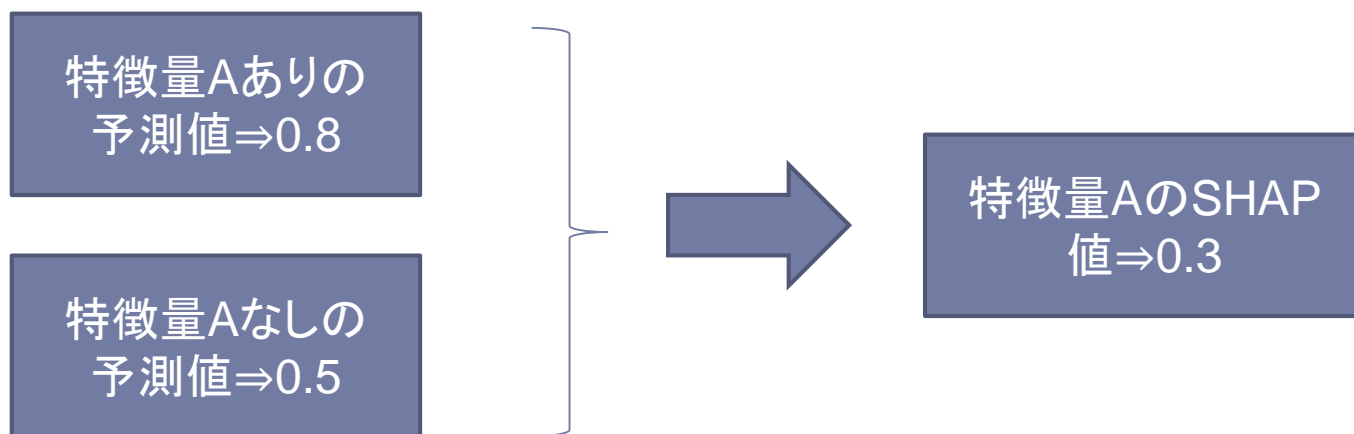
<https://own-search-and-study.xyz/2019/10/05/shap-all-methods/>



SHAP値の計算に必要なこと

- ▶ 「ある特徴量は予測にどう貢献したか」を計算する必要がある

⇒ゲーム理論の考え方 (Shapley value) を使う (次ページ)



この部分をどう算出するか？

SHAP値の基本的な考え方

- ▶ 特徴量*i*があることで、予測がどのように変わったかを全組み合わせについて算出して平均をとる

特徴量*i*のSHAP値 $\phi_i =$ $\sum_{S \subseteq N \setminus \{i\}}$ $\frac{|S|!(M - |S| - 1)!}{M!}$ $[f_x(S \cup \{i\}) - f_x(S)]$, (2)

iを使わないすべての組み合わせ

正規化係数

S+iの特徴量で予測した場合の結果と、Sの特徴量で予測した結果の差

<https://ja.wikipedia.org/wiki/%E3%82%B7%E3%83%A3%E3%83%BC%E3%83%97%E3%83%AC%E3%82%A4%E5%80%A4> より

i: SHAP値を求めたい特徴量

M: 特徴量の総数

N: 特徴量の全組み合わせの集合(順番考慮)

S: 特徴量*i*を含まない全組み合わせの集合の1つ(順番考慮)

f_x : 予測モデル

Sのイメージ

- ▶ 特徴量A,B,Cがあり、CのSHAP値を求めたいとするとSは以下の候補がある(Φ は空集合)
- ▶ Cがどのような状況で追加されたかの全てを考慮するという考え方

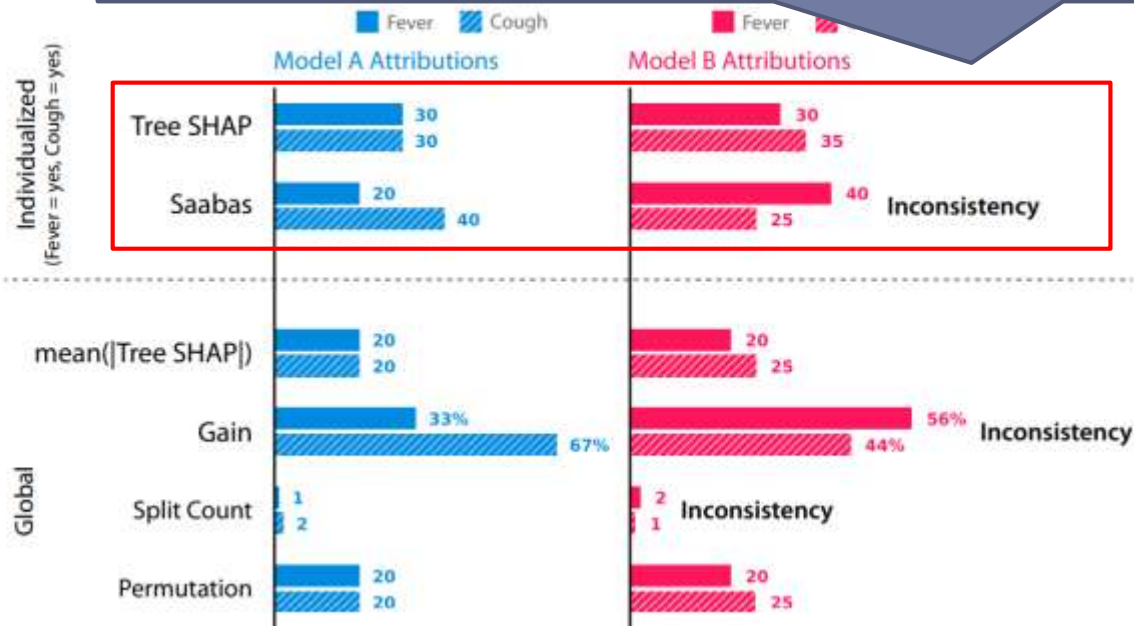
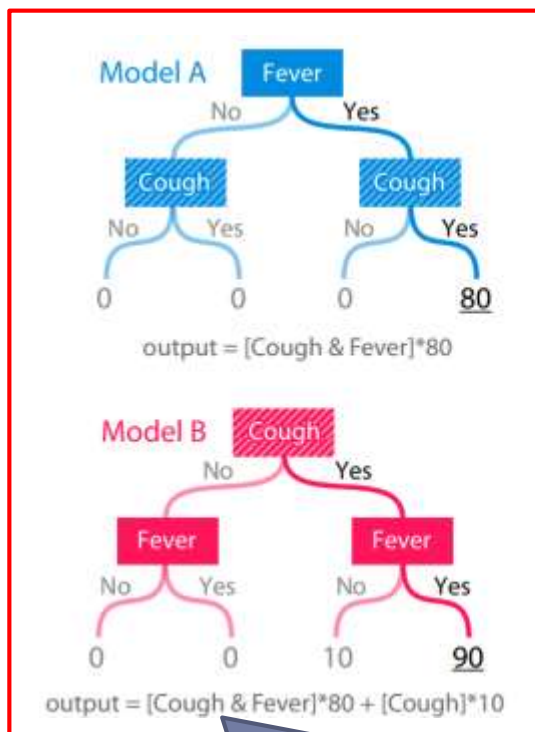
Case	順番	S
1	ABC	{A, B}
2	ACB	{A}
3	BAC	{B, A}
4	BCA	{B}
5	CAB	{ Φ }
6	CBA	{ Φ }

全順列を考えて、
Cより前にあるものを
抽出したのがS

何故Shapley valueの考え方が必要なのか？

- 論文では、Shapley valueの考え方がないと、似たようなモデルで一貫性のある結果が出ないと主張している

影響度の解釈が、モデルA/BでSHAPは一貫しているが、従来ロジック(Saabas)では倍半分違う



ほぼ同じことを表した2つのモデル
(熱と咳が出るなら風邪)

<https://arxiv.org/pdf/1802.03888.pdf> より

ここまでの流れをドラゴンボールで説明

- ▶ 悟空(G)、ピッコロ(P)、ヤムチャ(Y)が交代でフリーザ軍1000人を相手にするシチュエーションにおいてヤムチャの貢献を定量化したい。
- ▶ 単純にはヤムチャが何人倒すかを見ればよさそうだ。
- ▶ ヤムチャが何番目に登場するかで貢献度は変わる。

Case	順番	S	Fx(S)	Fx(SUY)	ヤムチャが倒した人数
1	GPY	{G, P}	1000人	1000人	0人
2	GYP	{G}	1000人	1000人	0人
3	PGY	{P, G}	1000人	1000人	0人
4	PYG	{P}	900人	960人	60人
5	YPG	{Φ}	0人	300人	300人
6	YGP	{Φ}	0人	300人	300人

ヤムチャのSHAP値
=110人
 $(0+0+0+60+300+300)/6$



具体的なSHAP値の算出（その1）

- ▶ $f(x) = E[f(x) | x_S]$ を以下のアルゴリズムで求めるとする
- ▶ 木を辿っていき、最後についたノードの値を返すアルゴリズム。ただし、モデルの分岐の特徴量がSに選ばれていない場合データ量の重みづけて両方を辿った場合を考慮する

Algorithm 1 Estimating $E[f(x) | x_S]$

procedure EXPVALUE($x, S, \text{tree} = \{v, a, b, t, r, d\}$)

procedure G(j, w)

if $v_j \neq \text{internal}$ **then**

return $w \cdot v_j$

else

if $d_j \in S$ **then**

return G(a_j, w) **if** $x_{d_j} \leq t_j$ **else** G(b_j, w)

else

return G($a_j, wr_{a_j}/r_j$) + G($b_j, wr_{b_j}/r_j$)

end if

end if

end procedure

return G(1, 1)

end procedure

x : 特徴量の値

S : 特徴量の部分集合

以下は、ノード数のベクトル

v : ノードの目的変数の値

a : 子ノードのindex(左)

b : 子ノードのindex(右)

t : ノードの分岐の値

r : ノードのデータの量

d : ノードの特徴量のindex

葉にたどり着いた時の処理

葉にたどりつくまで時の処理

モデルの分岐の特徴量がSに選ばれていない場合の処理

具体的なSHAP値の算出（その2）

- ▶ Sの組み合わせは、各特徴量を使う/使わないの組み合わせとなるため、すべての特徴量のSHAP値を求めるためには、 2^M の計算を実施する必要がある
- ▶ 論文では、計算を工夫することで多項式時間で計算するアルゴリズムを検討した

```
Algorithm 2 Tree SHAP
procedure TS(x, tree = {i, u, b, t, r, d})
   $\phi \leftarrow$  array of len(x) zeros
  procedure RECURSE( $l, m, p_1, p_2, p_3$ )
     $m \leftarrow$  EXTEND( $m, p_1, p_2, p_3$ )
    if  $v_l \neq$  internal then
      for  $i \leftarrow 1$  to len( $m$ ) do
         $w \leftarrow$  sum( $\text{UNWIND}(m, i, w)$ )
         $\phi_{m_i} = \phi_{m_i} + w(m_i \cdot a - m_i \cdot x) \cdot c_i$ 
      end for
    else
       $k, x = x_{d_i} \leq t_j ? (a_j, b_j) : (b_j, a_j)$ 
       $i_2 = i_1 = 1$ 
       $k \leftarrow$  FINDFIRST( $m, d_i$ )
      if  $k \neq$  nothing then
         $(i_2, i_3) \leftarrow (m_k, x, m_k, a)$ 
         $m \leftarrow$  UNWIND( $m, k$ )
        end if
      RECURSE( $k, w, i_2/r_k/r_j, i_3, d_j$ )
      RECURSE( $c, m, i_2/r_2/r_j, 0, d_j$ )
    end if
  end procedure
  procedure EXTEND( $m, p_1, p_2, p_3$ )
     $l \leftarrow$  len( $m$ )
     $n \leftarrow$  copy( $m$ )
     $m_{l+1}, (d, z, a, w) \leftarrow (p_1, p_2, p_3, l \in 0 ? 1 : 0)$ 
    for  $i \leftarrow l + 1$  to 1 do
       $m_{i+1}, w \leftarrow m_{i+1}, w + p_3 m_i \cdot w(i/l)$ 
       $m_i, w \leftarrow p_2 m_i, w[(l - i)/l]$ 
    end for
    return  $m$ 
  end procedure
  procedure UNWIND( $m, i$ )
     $l \leftarrow$  len( $m$ )
     $n \leftarrow m_i \cdot w$ 
     $m \leftarrow$  copy( $m_{l-1}, i-1$ )
    for  $j \leftarrow l - 1$  to 1 do
      if  $m_j \cdot a \neq 0$  then
         $t \leftarrow m_j \cdot w$ 
         $m_j, w \leftarrow n \cdot l / (j - m_j \cdot a)$ 
         $n \leftarrow t - m_j \cdot w \cdot m_j \cdot a[(l - j)/l]$ 
      else
         $m_j, w \leftarrow (m_j \cdot w \cdot l) / (m_j \cdot a[(l - j)])$ 
      end if
    end for
    for  $j \leftarrow 1$  to  $l - 1$  do
       $m_j, (d, z, a) \leftarrow m_{j+1}, (d, z, a)$ 
    end for
    return  $m$ 
  end procedure
  RECURSE(1, [1], 1, 1, 0)
  return  $\phi$ 
end procedure
```

長いので、詳細の解読は断念

一度の計算で、すべてのSを考慮しながら計算する作りになっている

今後の予定

- ▶ SHAPライブラリのツリー構造以外の理論を追いかける
- ▶ ツリー構造のアルゴリズム2の中身を理解する

