

2021年10月中旬 白田由香利

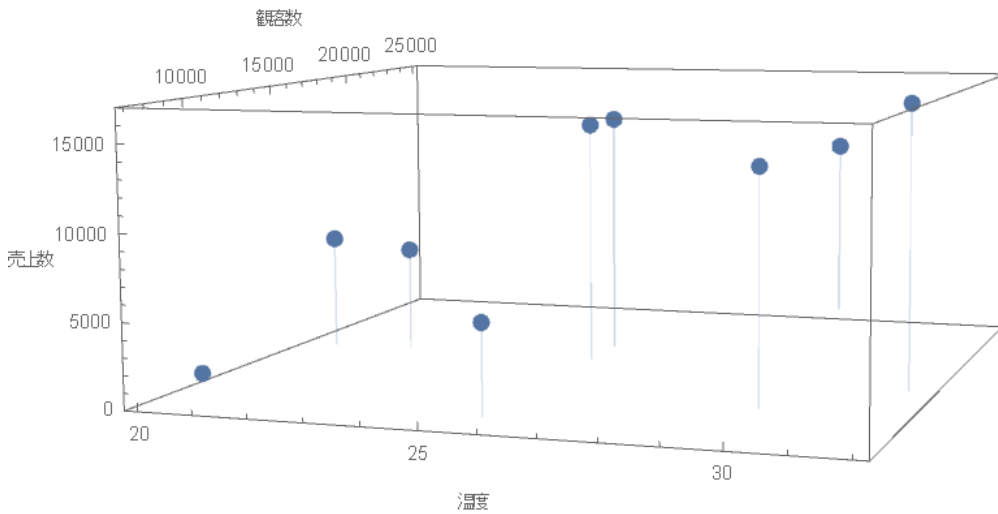
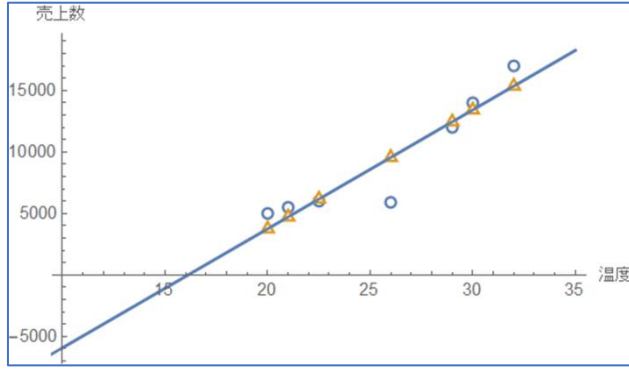
隣町まで足を伸ばして古本屋さんをあさっていて、奥の細道に関するすごいテキストを見つけてしまった。「身体感覚で芭蕉を読み直す『おくのほそ道』謎解きの旅」[1]。著者の安田氏はプロの能楽師であり、つまり能の謡は体に染みわたるように入っている人である。「芭蕉一門はおそらく全員が謡曲を習っていたし日常に口ずさんでいた。おくのほそ道の中でも謡曲の中からの引用が多く、その引用した能のイメージが脳裏に沸いてこそ、その芭蕉の気持ちを理解できる」ということが書いてある。例えば、「行春や鳥啼 魚の目は泪」の句から、季節は春であるし、能「道成寺」の「月落ち鳥ないて。。。」の謡が連想されて、道成寺のシテのへびに変身した女性が泳ぎ渡った日高川が連想され、芭蕉たちの眼前にある実物の隅田川が冥界への橋のように感じられる、というような連想ゲームであると。高校生のときに「隠された十字架」を読んだとき以来の驚き、スカイフォールであった。「ああ、そういうことだったのですね」と、膝を打ってしまう。機械学習でもそれは言える。サーキットランなどのパイソンのライブラリは世の中に流布しているので、それをインストールして使ってみることは容易い。しかし、そこにある数学的理論、例えば、Shapley 値の意味を定義から理解していないと、自分のやりたい分野で使いこなせない。評価が表面的で浅くなってしまうのだ。また、理論やモデルの拡張をしようとする場合には、絶対に、その数学的理論を理解していないとできない。ですからライブラリを使うには、その数学を理解したほうが良い。それゆえに、「まずは数学理論が何を意味しているのかグラフィクスで見てください。それから、数式を見ると理解が容易になりますよ」と、私はグラフィクス教材で数学公式を教えている。勉強すると深いところが分かるようになる。

能と芭蕉の句は密接な関係があると言われても、現代社会で謡曲に通じている人は殆どいない。私にしても、お稽古した演目は少なく、暗記している部分はさらに非常に非常に限定されていて、芭蕉の句をみても、どのお能のどこを引用したのか全く分からない。誰か謡曲のテキストデータベースを青空文庫のように作ってくると検索が容易になるのだが、そのようなデータベースは聞いたことがない。それに著作権問題も複雑そう。私が中高のとき、仕舞部で謡のお稽古をしたときは、拍子や抑揚を口伝えに習うだけで、意味の解説は無かった。つまり「何を言っているのか全く分かりません」のまま耳で覚えたのである。ですから、安田氏が上述した本の中で、能と謡曲の解説までしてくださり、ここがこう引用されていて、能のこのような感覚をもったのだろう、と解説してくださると、能と謡曲の解釈までも学べて、大変勉強になる。そしてお能を見に行きたくなる。コロナで映画ひとつ見に行かなかったが、久々にお能のチケットを買った。

さて、SHAP の話である。SHAP の理論を理解しないまま使うのでは、分析結果も浅くなってしまう。難しい内容を承知で、やはりゼミの皆さんには概要を理解して頂きたいと思う。そこで、SHAP に関してよく出る質問を以下にまとめてみた。分かりやすくするために少し脚色してある。

Q: 機械学習の手法で回帰モデルを作る、ということですが、モデルとはテスラのモデル3のような機種のことをいうのですか？

A: 回帰モデルとは、予測関数  $f(X)$  のことを言います。例えば、温度からビールの売上数を回帰によって予測した場合、温度  $temp$  に対して売上予測値  $f(temp)$  が求まります。

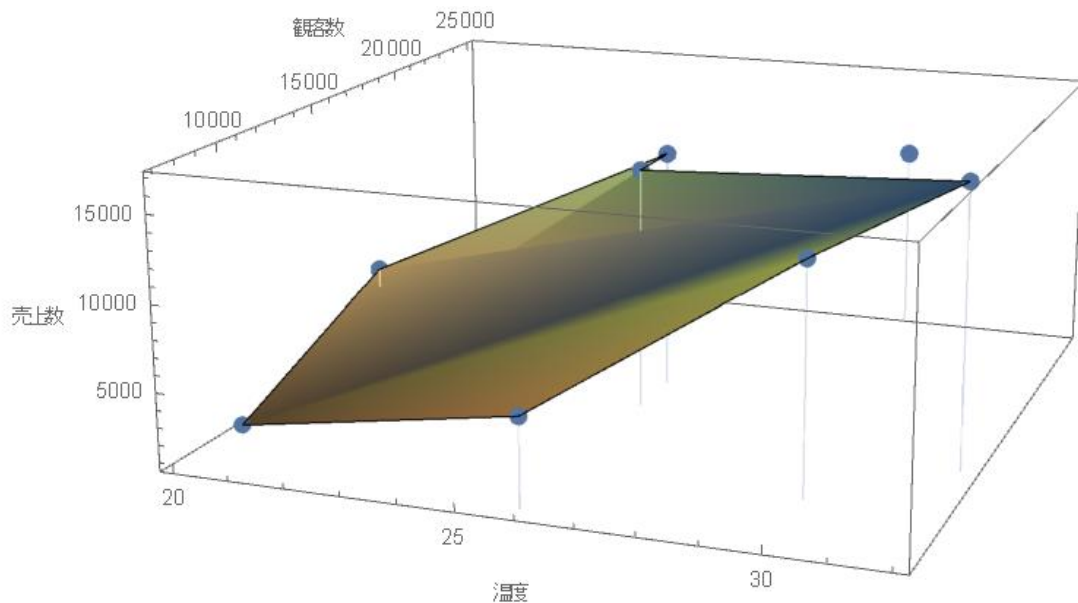


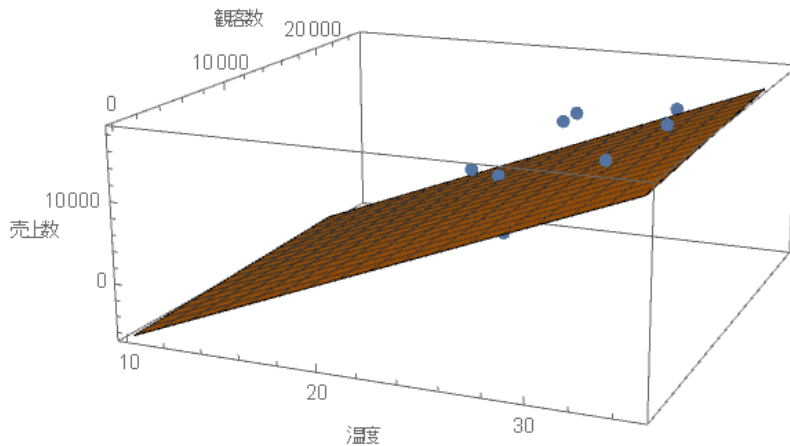
温度と観客数の2つの説明変数からベールの売上数を回帰によって予測した場合、温度 temp と観客数 spect の2つの引数に対して売上予測値  $f(\text{temp}, \text{spect})$  が求まります。左図を見てください。

温度の軸、観客数の軸、そして垂直方向に売上予測値

$f(\text{temp}, \text{spect})$  があります。{温度が30, 観客数が9800, 売上予測値は14000} という点が見えますね。

この観測値にできるだけフィッティングするようにして回帰モデル  $f(\text{temp}, \text{spec})$  を例えば以下のように作ります。平面ではなく、亀の甲羅のように形が複雑になっています。





伝統的な「線形重回帰分析」では左図のような平面で近似しますので、フィッティングがよくないです。複雑形状の関数  $f$  を使って近似をする機械学習回帰のほうがフィッティングがよくなります。

Q: 複数のプレイヤーというのが、複数の説明変数に対応するというのは分かりますが、儲けた金額というのは、回帰分析の何が対応するのですか？

A: ターゲット値が対応します。例えば、ビールの売上数 14000 という値が対応します。

Q: 複数プレイヤーは人間なので、人数が増えたら利益金額が増加する、というの分かりますが、説明変数が増えるとターゲット値が増える、というのをおかしいと思います。

A: 気持ちは分かりますが、「どの説明変数が重要でしょうか？」と考えて、説明変数「温度」のほうが「観客数」のほうよりも重要である、というような結果を出したいのです。それで、**特性関数**を求めて、この7月30日の特性として、温度だけだと売上数は10000であるが、温度&観客数の2つの説明変数の効果を評価すると売上数は15000に増加する、というような関係性を見たいのです。

Q: 説明変数1個の場合、予測値がいくらになるかなんてわからないと思います。そのような関数を求めるなんて絶対無理だと思います。

A: その通りです。それで、**疑似的に**、参加していない説明変数の値には、**全体の平均値**を入れて計算します。例えば、7月30日(温度は30度だった)の特性関数を求めるとして、温度だけの貢献をみたいとき、 **$f(30, \text{観客数の平均値}14088.89)$** で計算します。 $f(\text{temp}, \text{spec})$ の形状は機械学習の回帰で求めた亀の甲羅のような複雑形状の関数を使います。7月30日(観客数は17000人のだった)の特性関数の、温度と観客数を合わせた貢献を計算する場合  **$f(30, 17000)$** と計算します。17000人は平均観客数の14088.89人よりも大きいですから、温度だけの貢献よりも値はきっと増えるでしょう。しかし、反対に、7月30日の観客数が12000人と少なかった場合、 **$f(30, 12000)$** の値は  **$f(30, \text{観客数の平均値}14088.89)$** よりも小さくなりそうです。

Q: その場合、7月30日特性関数において観客数の貢献はマイナスとなるのですか？

A: 疑似特性関数の計算では、全て平均値を代入した値  **$f(\text{温度平均値}, \text{観客数平均値})$** をベースとしているので、その値よりも小さい予測値もできます。

Q: 参加していない説明変数には、平均値を入れる、という荒っぽい感じのやり方で、本当に正しく分析できるのでしょうか？

A: 経営分析では、**その会社が業界平均よりも上か下かを重要視**します。その説明変数に業界平均値を入れて、他の説明変数の貢献を測る、というやり方は納得できるやりかただと思います。実際、このやり方で納得できる結果がでています。

Q: Shapley 値の以下の公式が理解できません。 
$$\phi_i = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(|N|-|S|-1)!}{|N|!} [v(S \cup \{i\}) - v(S)]$$

A: ここが最大の難関かと思います。 $\phi_i$ (ファイの*i*番目)は、*i*番目のプレイヤー(説明変数)の Shapley 値のことです。その値をどのように定義したかを表しています。この定義式では、特性関数を  $v$ (ギリシャ語でニュー)で表しています。グループ  $S$  にプレイヤー*i*さんが参加したら いくら利益が増えるかは、 $v(S \cup \{i\}) - v(S)$  と表現できます。プレイヤー全体の集合を  $N$ 、 $S$ はその部分集合、 $|S|$ を集合  $S$  のメンバー数とします。

$(S \subseteq N \setminus \{i\})$  という表現は **プレイヤー*i*さんを含まない全ての部分集合  $S$** を表します。  $i=2$  として以下説明します。つまりプレイヤーが5人いた場合 ( $|N|=5$ )、プレイヤー2番さんを含まない全ての集合は、 $2 \times 2 \times 2 \times 2 = 16$ 通りです。

プレイヤー 2さんの Shapley 値は以下のように定義されます。

$$\phi_2 = \sum_{S \subseteq N \setminus \{2\}} \frac{|S|!(|N|-|S|-1)!}{|N|!} [v(S \cup \{2\}) - v(S)]$$

始めグループ  $S$  だけで協同作業していたところにプレイヤー2番さんが参加してくれて、おかげで利益が  $v(S \cup \{2\}) - v(S)$  だけ増えました。集合  $S$  の取り方は **16通り**ありました。では、 $S = \{1,3\}$ としてみましょ。これは16通りの中の1つです。仕事に参加した順番で順に並んでもらうとすると、 $S$ の2人の順列は  $2! = 2 \times 1$  となります。そこにプレイヤー*i*さんが参加します。その後、 $\{4,5\}$ の2人が参加します。この2人の人数は全体の  $5 - 2 - 1 = 2$ で計算できます。この部分を変数でかくと、 $(|N| - |S| - 1)$ となります。縦棒でくるむと  $N$ の人数  $|N|=5$ 、 $S$ の人数  $|S|=2$ を意味します。全体では、ならばかたの総数は以下の掛け算になります。

$$|S|! \times 1 \times (|N| - |S| - 1)!$$

$|N|$ 人のプレイヤーが並ぶ順列は  $|N|!$  通りです。ですから、 $\{1,3\} \rightarrow \{2\} \rightarrow \{4,5\}$  という順番に並ぶ確率は、 $\frac{|S|!(|N|-|S|-1)!}{|N|!}$  となります。プレイヤー2の Shapley 値は、集合  $S$  の取り方の16通りについてこれを計算して合計したものです。

Q: 難しく理解できません。

A: 了解です。 $\phi_i$ は、*i*番目のプレイヤー(説明変数)の Shapley 値ということだけ覚えておいてください。計算方法は理解できなくてもよし、とします。

Q: 企業分析のときのソニーの特性関数、パナソニックの特性関数とかの話はどう関連するのでしょうか?

A: 特性関数  $v$  は、各企業に対して別関数として定義されます。ソニーの特性関数とパナの特性関数は異なります。ですから、Shapley 値は、ソニーにおける説明変数「棚卸資産回転率」の Shapley 値は  $x$   $x$  です、というように企業ごとに計算されます。

終わり

引用元: 行春や鳥啼 魚の目は泪 芭蕉

参考文献

[1] 安田登:「身体感覚で芭蕉を読み直す『おくのほそ道』謎解きの旅」、春秋社、2012年。