

ニューラルネットを用いた 解釈可能な非平行連続順序回帰モデルの提案

奥野 彰文 統計思考院 助教

※本研究は原田和治博士(東京医大)との共同研究であり <https://arxiv.org/abs/2303.17823> にてプレプリントを公開済です。

1. 背景

例えば家の築年数(共変量 X)から、家の値段(応答変数 H)が u [ドル] 以下である確率

$$\mathbb{P}(H \leq u | X = x)$$

を予測する手法を「**順序回帰**」と呼びます。このとき、例えば

$$\mathbb{P}(H \leq u | X = x) = \sigma(a(u) + \langle b, x \rangle)$$

のようにし、単調増加関数 $a(u)$ と係数ベクトル b をうまく推定できれば、係数 b は共変量 x が応答 H に与えるある種の影響力を表現していると考えられます。この単純なモデルでは閾値 u の値に依らず係数が一定であることを仮定しており、これを順序回帰では**平行性仮定**と呼びます。

現実的な問題として、例えば値段の高い家に求められる要素と値段の安い家に求められる要素は異なりますから、平行性仮定はやや強い仮定です。そこで係数が閾値 u に依存して変わる**非平行連続順序回帰モデル**

$$\mathbb{P}(H \leq u | X = x) = \sigma(a(u) + \langle b(u), x \rangle)$$

を用いることがより自然と考えられますが、真の累積分布が u に関して単調増加するのに対し、 u が単調に増加しても非平行連続順序回帰モデルの予測累積分布は単調増加しない場合があります。つまり非平行連続順序回帰モデルは**確率モデルとして不適切ではないか?**という問題がありました。

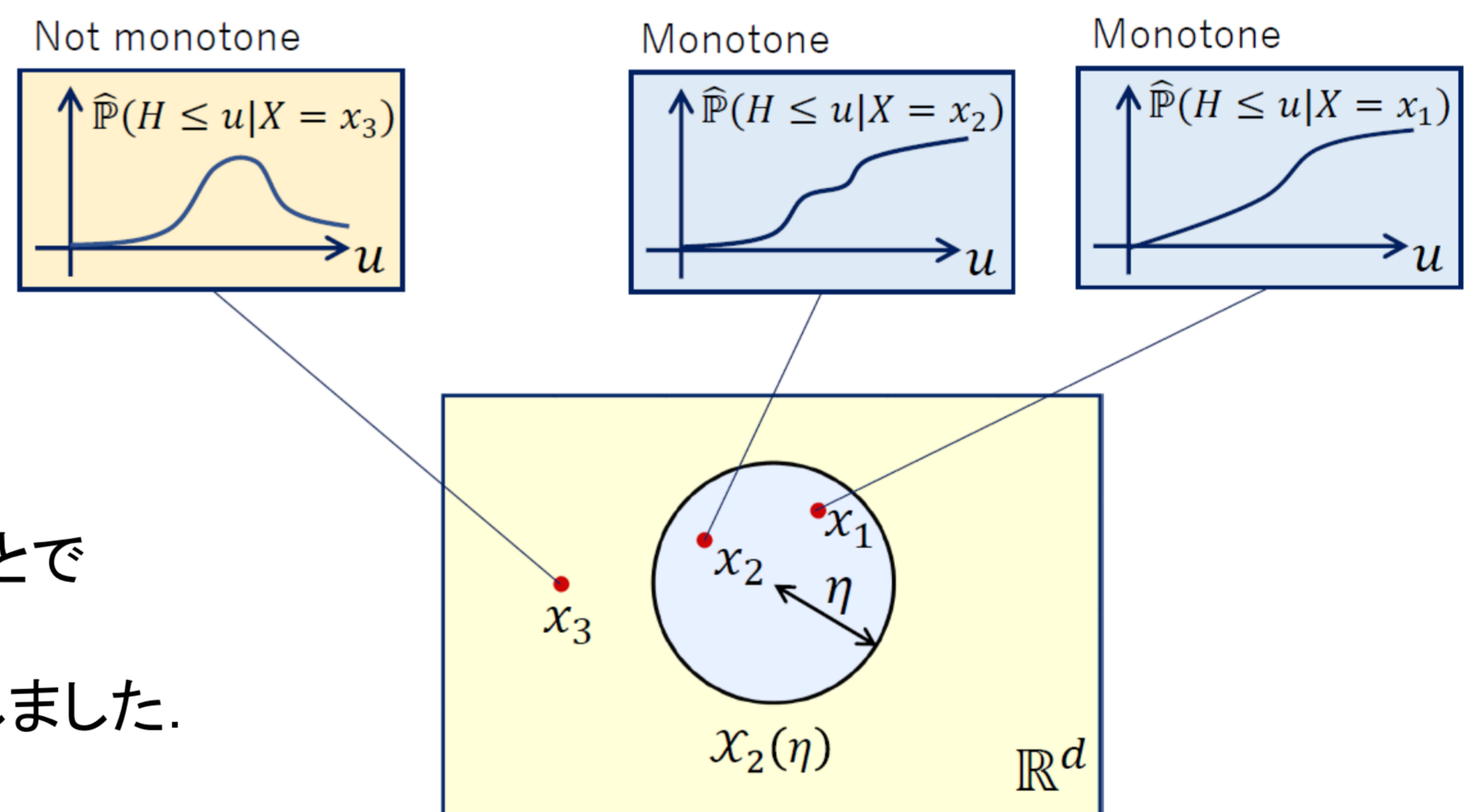
2. 我々の結果

理論的な進展: 単調性を保証した最適化アルゴリズム

我々の研究ではまず、共変量 X が \mathbb{R}^d 全域に値を取るとき、**予測累積分布関数が常に単調である(確率モデルとして常に適切である)ような非平行連続順序回帰モデルが存在しない**ことを示しました(Proposition 1)。このPropositionは既存法を含めたすべての非平行連続順序回帰モデルの限界を示しています。

次に、非平行連続順序回帰モデルの係数ベクトルを表す関数 $b(u)$ として隠れ層1層の**ニューラルネットワーク**を利用し、予測累積分布関数が任意の $x \in \mathcal{X}_2(\eta) := \{x \in \mathbb{R}^d : \|x\|_2 \leq \eta\}$ に対して**単調であるための十分条件**を示しました(Proposition 2)。

提案した十分条件を利用し、ニューラルネットの重みを更新することで単調性を担保させながら非平行連続順序回帰モデルを学習する**Monotonicity-Preserving Stochastic (MPS) アルゴリズムを提案**しました。



実際のデータへの適用と考察

ここで、例えば家の価格が u [ドル] 以上である確率

$$\text{logit}(\mathbb{P}(H > u | X = x)) = r(u) + \langle s(u), x \rangle$$

を予測してみますと、 $s(u) = -b(u)$ がある成分が正の値をとれば、対応する共変量は家の価格の向上に寄与すると考えられます。

提案法をreal-estate datasetに適用した右のグラフを見てください。House_age (築年数) の係数は、House_price (家の価格) が低いとき -0.8 でネガティブな影響を持ちますが、**家の価格が高いとき係数は0に近付き、つまり家の価格に影響を及ぼしにくくなる**、という仮説を読み取ることができます。

本発表では家の価格を例に挙げましたが、提案法は例えば何が車の燃費効率を上げるかとか、何がレストランの点数を上げるかとか、様々なデータの解釈に利用できます。

Dataset: real-estate
住宅価格が高いところでは係数がゼロつまり築年数が家の価格に影響を及ぼしにくそうである。

住宅価格が低いところでは係数がマイナス、つまり築年数が家の価格に悪影響を及ぼしそうである。

