

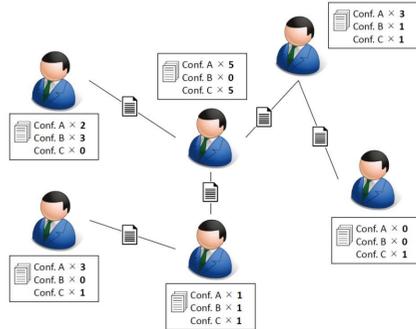
ノンパラメトリックなリンク予測に関する理論

奥野 彰文 統計思考院 助教

※本研究は矢野恵佑准教授(統計数理研究所・数理推論系)との共同研究であり <https://arxiv.org/abs/2012.13106> にてプレプリントを公開済です。

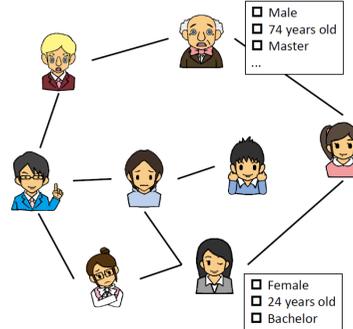
1. 背景

グラフのノードの特微量(共変量)を $\{x_i\} \subset \mathcal{X} \subset \mathbb{R}^p$, ノード間のリンクの重みを $\{y_{ij}\}_{1 \leq i < j \leq n} \subset \mathcal{Y} (\subset \mathbb{R})$ と書きます:



例) 共著ネットワーク:

$y_{ij} \in \mathbb{N}$: 研究者 i, j が2人で書いた論文の本数
 $x_i \in \mathbb{N}^p$: 研究者 i が各雑誌に乗せた論文の本数を並べたもの



例) 友達ネットワーク:

$y_{ij} \in \{0,1\}$: 個人 i, j が友達であるかどうか
 $x_i \in \mathbb{N}^p$: 個人 i の性質(年齢など)を並べたもの

リンク回帰問題

組 (x_i, x_j) から 重み y_{ij} を予測したい: $y_{ij} \approx f_\theta(x_i, x_j)$

※類似度学習などはリンク回帰に含まれる。

2. 提案法

ノンパラメトリックリンク回帰

$$\hat{f}_{n,h}(x, x') := \frac{\sum_{1 \leq i < j \leq n} y_{ij} K_h(x - x_i) K_h(x - x_j) + \lambda_{n,h}}{\sum_{1 \leq i < j \leq n} K_h(x - x_i) K_h(x - x_j) + \lambda_{n,h}}$$

(ただし $h = h_n > 0$ はバンド幅, K_h はカーネル関数, $\lambda_{n,h} > 0$ は正則化パラメータ.)

$h = h_n \rightarrow 0$ のとき推定量は一致性を持つ: $\hat{f}_{n,h}(x, x') \rightarrow f_*(x, x') = \mathbb{E}[Y | X = x, X' = x']$.

3. 予測値の漸近分散に関する発見

ノンパラメトリック統計解析では共変量 x_i に以下のどちらかのデザインを仮定することが多い:

Random design: 共変量 x_i が実験ごとにある分布からランダムに生成されると仮定する

Fixed design: 共変量 x_i はすべての実験で常に同じ値をとると仮定する

古典的なノンパラメトリック回帰では, 上記のどちらのデザインを仮定しても漸近的に同等である (Brown and Low (AoS1996), Reiss (AoS2008) など).

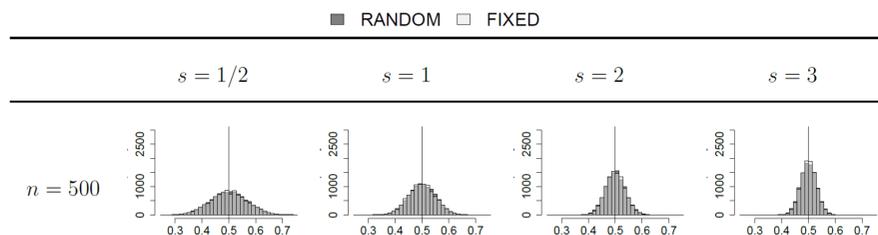


表1: 古典的なノンパラメトリック回帰での予測値のヒストグラム. デザインによらずヒストグラムはほぼ同じ. バンド幅 $h = n^{-\frac{1}{s+d}}$.

では, 提案したノンパラメトリックリンク回帰ではどうなるか?

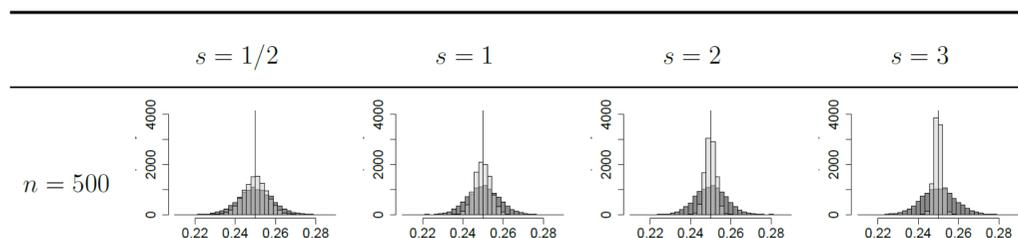


表2: ノンパラメトリックリンク回帰での予測値のヒストグラム. バンド幅によらず, ヒストグラムがデザインに依存する. バンド幅 $h = n^{-\frac{1}{s+d}}$.

ノンパラメトリックリンク回帰の予測値の漸近分散は共変量 x_i のデザインに依存することを発見し, 理論的にも評価した.

これは既存のノンパラメトリック解析で一般的に知られていない特殊な現象であり, 類似度学習などで予測値の(漸近的な)統計的性質を調べるには特別な注意を払う必要があることが分かった.