

事前圧縮による地理的加重回帰の高速化

村上 大輔 データ科学研究系 助教

【背景】

場所毎の回帰係数 (SVC: Spatially varying coefficients) を推定する方法である Geographically weighted regression、(GWR) は応用分野で幅広く用いられてきた。しかしながら GWR の推定に必要な Leave-one-out cross-validation (LOOCV) の計算量は標本数 N の2乗のオーダーであり、大規模な空間データには不向きである。

高速なSVC推定手法の比較

SVC を高速推定する数多くの手法が存在するが、過度な平滑化を招かず (非縮退) かつ高速に SVC を推定する方法が存在しない

手法	推定法	計算量	非縮退
Fast GWR (Li et al., 2019)	局所	$O(N^2)$	○
Expansion method (Casetti, 1972)	大域	$O(N)$	×
Bivariate spline method (Mu et al., 2018)			
ME approach (Murakami and Griffith, 2019)	局所	$O(N)$ だがMCMCが必要。遅い。	○
Predictive process (Finley et al., 2011)			
Nearest-neighbor GP (Finley et al., 2019)	局所	$O(N)$	○

→ そこで本研究では GWR を高速化 (+ R パッケージ `scgwr` で公開)

提案型 GWR	局所	$O(N)$	○
---------	----	--------	---

【提案手法】

以下の手順で LOOCV を高速化した (詳しくは右図):

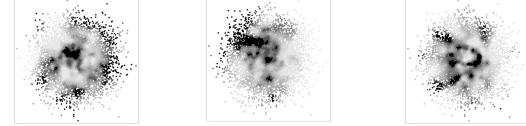
- (1) カーネル関数を線形形式で近似
- (2) パラメータを含まない行列を事前処理することで、サイズが N に依存する行列・ベクトルを消去
- (3) LOOCV を実行

【シミュレーション実験】

- $N \in \{500, 1000, 1500, 2000, 3000, 5000, 10000, 20000, 80000\}$
- 下記から生成した擬似データに繰り返しあてはめることで (各ケース 200 回反復)、従来型 GWR と提案型 GWR を比較

$$y = \beta_0 + x_1 \beta_1 + x_2 \beta_2 + \varepsilon \quad \varepsilon \sim N(0, I) \quad x_k \sim N(0, I)$$

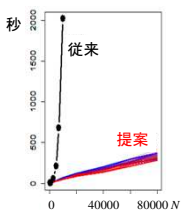
$$\beta_0 \sim N(1, 0.5^2 \hat{G}) \quad \beta_1 \sim N(1, 2^2 \hat{G}) \quad \beta_2 \sim N(1, 0.5^2 \hat{G})$$



\hat{G} : 第 (i, j) 要素を $g(d_{i,j}; b = 1)$ とするカーネル行列の Nystrom 近似

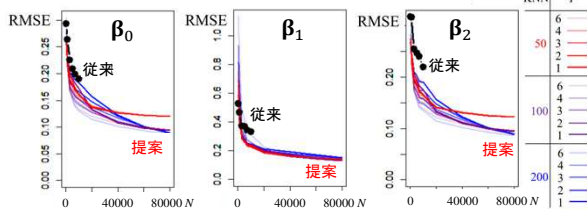
結果: 計算時間

高速化したことを確認



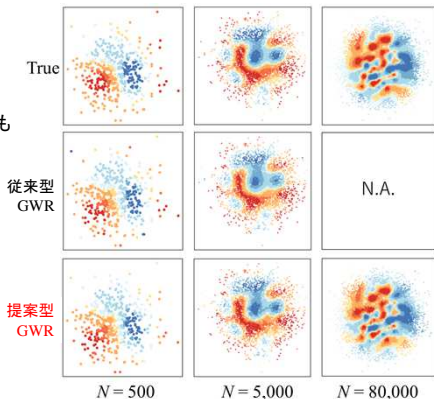
結果: SVC推定精度

カーネルをマルチスケールにしたため精度改善



推定された SVC の例

標本が大きくなっても縮退せずに精度良く SVC が推定されたことを確認



【実データへの応用】

東京23区内の非侵入窃盗件数の要因分析に応用

非侵入窃盗件数 (2012; 町丁目別)

すり、自動車盗、万引き、置き引き、...

説明変数

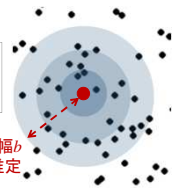
- 先月の犯罪件数
- 昼間人口密度 (国勢調査)
- 夜間人口密度 (国勢調査)

出典: 大東京防犯ネットワーク (<https://www.bouhan.metro.tokyo.lg.jp/>)

カーネル $g(d_{i,j}; b)$ のイメージ

$d_{i,j}$: 地点 i と j の間の直線距離
 b : バンド幅。影響の空間範囲を決める

バンド幅 b LOOCV で推定



【GWR】

- 地点 i の回帰係数を下式で推定する
- LOOCV でバンド幅を最適化した後に回帰係数を推定

$$y_i = \sum_{k=1}^K x_{i,k} \beta_{i,k} + \varepsilon_i$$

$$E[\varepsilon_i] = 0 \\ \text{Var}[\varepsilon_i] = \sigma^2$$

地点 i の回帰係数

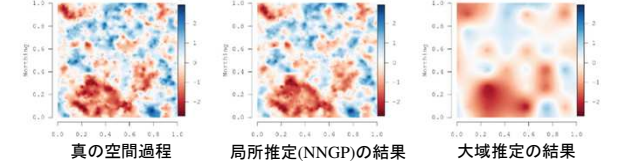
- 近隣の標本に大きな重みを与えて推定

$$\hat{\beta}_{i,k} = \underset{\beta_{i,k}}{\text{argmin}} \sum_{j=1}^N g(d_{i,j}; b) \left(y_j - \sum_{k=1}^K x_{j,k} \beta_{i,k} \right)^2$$

推定法と縮退

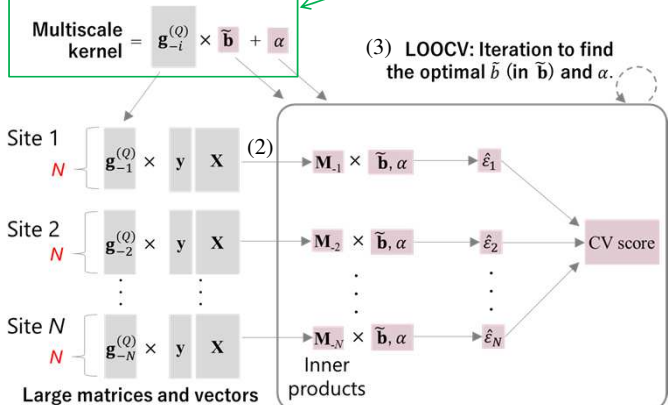
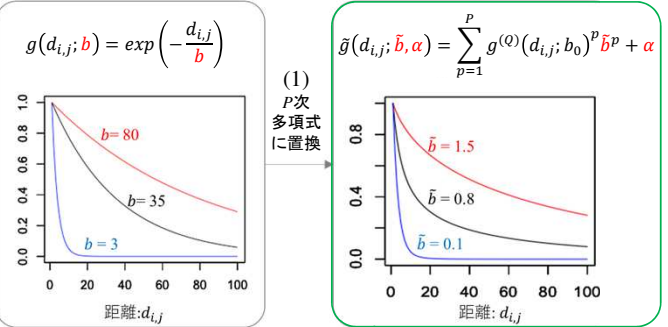
- **大域:** 空間基底関数 L 個 ($L \ll N$) の線形和で SVC をモデル化 (i.e., $\beta = E\gamma$)
→ 計算負荷は小さいが、縮退する (SVC が過度に平滑化される)
- **局所:** 場所毎の局所モデルを推定
→ 縮退しない。並列化しやすい

Datta et al. (2016)



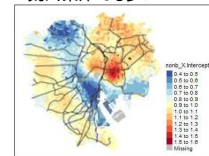
提案する LOOCV の手順

バンド幅 b は予め b_0 で固定。代わりに局所 (Q 近傍) と大域に対する重み (\bar{b}, α) を推定



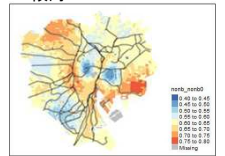
場所毎の定数項

- 都心・下町で多い
- 荒川東岸でも多い



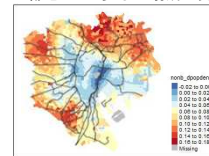
先月の犯罪件数の影響

- 山手線沿線で繰り返される傾向



昼間人口密度の影響

- 郊外: 人口密集地で発生
- 都心: 人が少ない場所で発生



夜間人口密度の影響

- 沿岸部の人が少ない場所は夜危険

