

大規模データのための時空間回帰モデリング

村上 大輔 データ科学研究系 助教

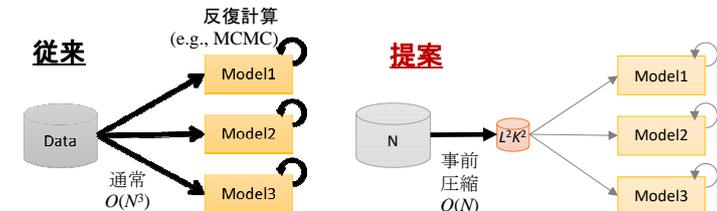
【地理空間データのオープン化】

- Google Earth EngineやOpenSteeMapといったオープンデータが急増
- 大規模な地理空間データを効率よく解析する方法が求められている

大規模な地理空間データの回帰・予測等を計算効率よく行うための実用手法の開発

【方針:大規模データモデリングの冗長性排除】

- 従来:モデル推定毎に大規模データの反復処理が必要→**重い**。
- 大規模データを圧縮してからモデルを推定するようになればよいのでは?



【提案手法】

- 線形混合効果モデルを仮定

$$y = \sum_{k=1}^K x_k \circ \beta_k + E y + \varepsilon \quad y \sim N(0, \tau^2 \Lambda^\alpha) \quad \varepsilon \sim N(0, \sigma^2 I)$$

ガウス過程の E: 空間近接行列のL個の固有ベクトル(の近似)
低ランク近似 A: L個の固有値からなる対角行列

空間可変パラメータ(今回はこれ)

場所毎に回帰係数を変える(場所以外も可)

$$\beta_k = b_k \mathbf{1} + E \gamma_k \quad \gamma_k \sim N(0, \tau_k^2 \Lambda^{\alpha_k})$$



グループ効果

小地域推定等にも使えそう

$$\beta_k = b_k \mathbf{1} + G v_k \quad v_k \sim N(0, \tau_k^2 I)$$

自己回帰・スプライン回帰等にも応用可能

Theta = {tau_1^2, ..., tau_K^2, alpha_1, ..., alpha_K} の推定手順

- (1) 事前圧縮→サイズがNの行列を消す

$$- M_{XX} = X'X, M_{XE} = X'E, m_{Xy} = X'y, m_{EY} = [E, \dots, E]y, m_{yy} = y'y$$

- (2) 事後確率 p(theta|y) を以下のように書き換える→推定の計算量がNに依存しない

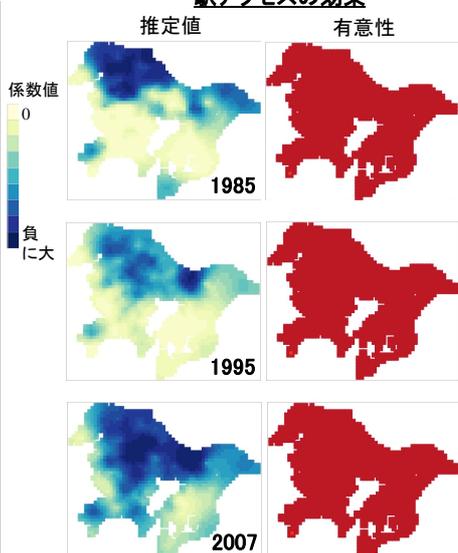
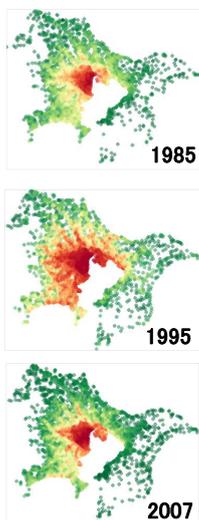
$$p(\theta|y) = -\frac{1}{2} \ln \left| \begin{bmatrix} M_{XX} & M_{XE} \tilde{V}(\theta) \\ \tilde{V}(\theta) M'_{XE} & \tilde{V}(\theta) M_{EE} \tilde{V}(\theta) + I \end{bmatrix} \right| - \frac{N-K}{2} \left(1 + \ln \left(\frac{2\pi d(\theta)}{N-K} \right) \right)$$

- (3) p(theta|y) をk毎に逐次最大化(省略)→Kが大きくても高速推定可

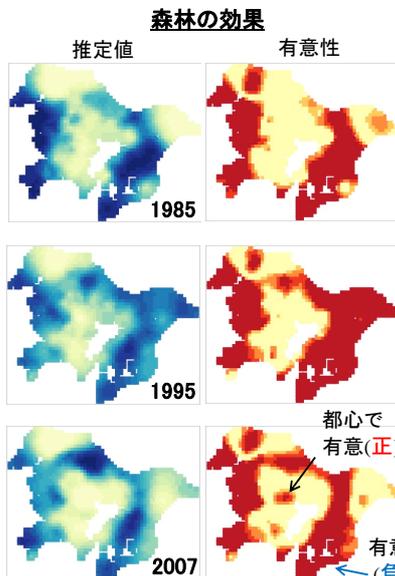
【実データへの適用】

- 住宅地公示地価の回帰に適用した後、係数を補間

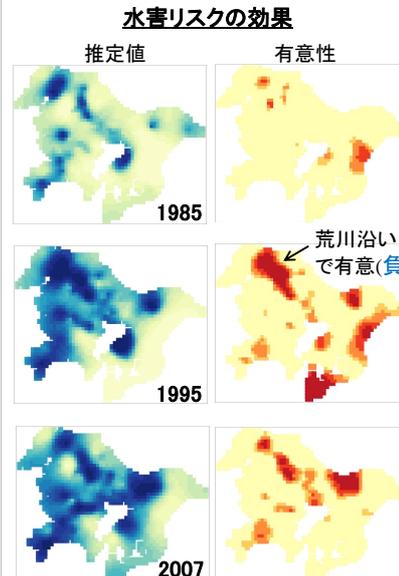
住宅地価データ



駅アクセスの効果



森林の効果



水害リスクの効果

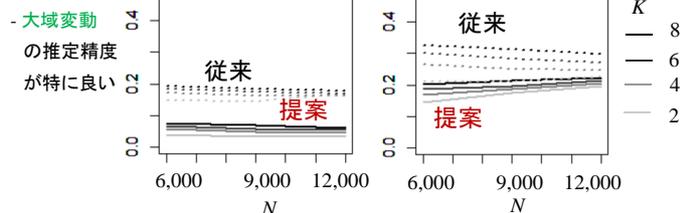
【シミュレーション実験】

- N ∈ {6,000, 9,000, 12,000}、空間可変パラメータ数 K ∈ {2, 4, 6, 8} を想定
- 緯度経度は標準正規分布からそれぞれ生成
- 以下から生成されるデータへのあてはめを200回繰り返す ※K=8の場合

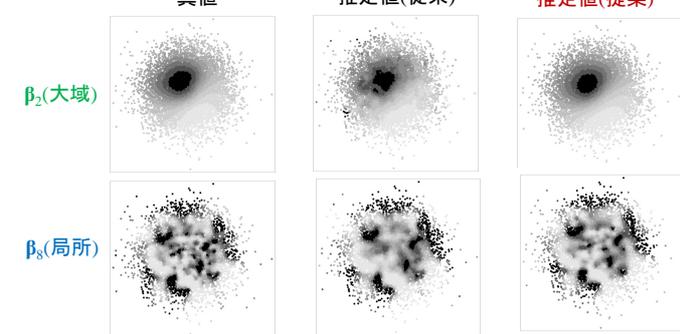
$$y = \sum_{k=1}^4 x_k \circ \beta_k + \sum_{k=5}^8 x_k \circ \beta_k + \varepsilon \quad \varepsilon \sim N(0, \sigma^2 I)$$

大域的 (alpha=3) 局所的 (alpha=0.5) R^2=0.71に なるよう設定

推定精度

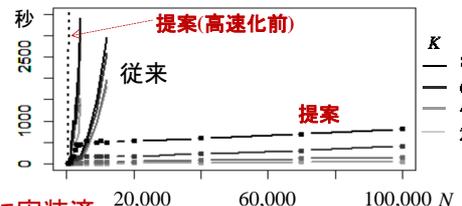


結果の例



計算時間の比較

- 計算時間を大幅に短縮
- 事前圧縮の計算量のみ Nに応じて増加



→ Rパッケージ spmoran に実装済

$$d(\theta) = m_{yy} - 2[\hat{b}, \hat{u}] \begin{bmatrix} m_{xy} \\ \tilde{V}(\theta) m_{EY} \end{bmatrix} + [\hat{b}, \hat{u}] \begin{bmatrix} M_{XX} & M_{XE} \tilde{V}(\theta) \\ \tilde{V}(\theta) M'_{XE} & \tilde{V}(\theta) M_{EE} \tilde{V}(\theta) \end{bmatrix} \begin{bmatrix} \hat{b} \\ \hat{u} \end{bmatrix} + \|\hat{u}\|^2$$

誤差分散 回帰係数の分散

$$\begin{bmatrix} \hat{b} \\ \hat{u} \end{bmatrix} = \begin{bmatrix} M_{XX} & M_{XE} \tilde{V}(\theta) \\ \tilde{V}(\theta) M'_{XE} & \tilde{V}(\theta) M_{EE} \tilde{V}(\theta) + I \end{bmatrix}^{-1} \begin{bmatrix} m_{xy} \\ \tilde{V}(\theta) m_{EY} \end{bmatrix} \quad \tilde{V}(\theta) = \frac{1}{\sigma^2} \begin{bmatrix} \tau_1 \Lambda^{\alpha_1} & & \\ & \ddots & \\ & & \tau_K \Lambda^{\alpha_K} \end{bmatrix} \quad \hat{u} = \begin{bmatrix} u_1 \\ \vdots \\ u_K \end{bmatrix} \quad u_k \sim N(0, \sigma^2 I)$$