3 調 査 研 究

3・1 報 文

1)ノンパラメトリックデータの回帰:一般化加法モデルによる 大気環境データの非線形回帰モデル構築に関する検討

古澤 尚英

要旨

ノンパラメトリックデータの回帰・推論・予測を行う場合、一般化加法モデルを用いると、 線形回帰モデルに比べて予測精度が高く、モデルの説明性が高いと言われている。本報では、 ノンパラメトリックデータである光化学オキシダント濃度について、一般化加法モデルで日射 量、気温及び風速データを用いた回帰・推論・予測を行い、最適なモデルを得ることができ た。また、一般化加法モデルは線形回帰モデルより高い精度で回帰・推論・予測を行うことが でき、環境データを解析する際の重要な選択肢の一つとなることが示された。

キーワード:一般化加法モデル,光化学オキシダント,予測,推論,R言語

はじめに

大気汚染物質は、大気汚染防止法第 22 条に基づき都道 府県知事が監視することとなっており¹⁾、大気汚染防止 法第 22 条の規定に基づく事務の処理基準²⁾に基づき常 時監視項目や観測局数等の体制整備に関することが定め られている。大気汚染物質のうち、自動測定機によりモ ニタリングされている窒素酸化物、硫黄酸化物、微小粒 子状物質等の物質は、令和4年度の全国調査結果におい て環境基準をほぼ100%達成していたが、一方で、光化学 オキシダント(以下「Ox」という。)は環境基準達成率 が 0.1%)であった。

このような状況から, Weather Research and Forecasting model (WRF) や (CMAQ) 等の大気環境シミュレーショ ン等を用いて,大気汚染物質の空間濃度分布の予測,大 気汚染物質の生成,飛来予測などの研究が進められてい る⁴⁾⁻⁶⁾。また,大気汚染物質の空間濃度分布予測は,大 気汚染常時監視測定局の配置検討や光化学オキシダント 注意報等発令地域区分の検討等,行政的な活用がしやす い有用なデータであるが,一方で,大気環境シミュレー ションには専用のソフトウェアの準備が必要であり,こ れらのソフトウェアの構築には多大なコストがかかるた め地方自治体職員が単独で実施することは難しく,シミ ュレーションモデルを用いた検討は一部の研究者が先行 して実施している状況である。

そこで筆者を含む検討メンバーは、国立環境研究所と 地方環境研究所とのⅡ型実施共同研究⁷⁾の枠組みにおい

て、常時監視測定局の観測体制見直し及び発令地域区分 等の検討に使用する目的で、大気汚染物質の空間濃度分 布を統計・機械学習的に求めるライブラリの開発を行っ た^{8),9)}。 このライブラリに導入した空間濃度分布を計 算する Regression Kriging 法 (以下「RK 法」という。) は、土地利用回帰モデル (Land Use Regression: LUR) と 地球統計学手法である Ordinary Kriging (OK) を組み合わ せた手法であり, 大気汚染物質の空間濃度分布を精度良 く予測することができると報告されている^{10),11)}。 筆者 らは既報¹¹⁾ で微小粒子状物質を対象とした検討を行った が、このライブラリでは予測できる大気汚染物質の対象 を拡張して構築した。 また、LUR の計算方法に、線形 回帰モデル以外にもランダムフォレスト、サポートベク トルマシン、一般化線形モデル、一般化加法モデル及び ニューラルネットワークの全6モデルを利用できるよう にした⁸⁾。 このライブラリを利用して Ox の高精度な空 間濃度分布を計算する条件検討を進めているが、これら 6モデルのうち一般化加法モデル (以下「GAM」とい う。)を実装するための日本語情報が少なく、また、大 気環境データのような時間-空間情報を持つノンパラメト リックデータに適応させた事例は更に情報が少ない状況 であった。 そのため、本報では Ox 空間濃度分布を予測 する GAM の構築・精度検証の検討を行い, GAM の有用 性について実証するとともに、今後の環境データ解析に おける基礎検討資料とすることを目的とする。

方 法

1. 線形回帰モデルについて

従来,環境分野における事象の解析や予測には線形回帰 モデルが広く利用されてきたため,GAMと比較するため にもまずは線形回帰モデルの構築と解釈の方法を解説す る。

式1に線形回帰式を示す¹²⁾。 p>1 である場合は重回 帰式とも呼ばれる。

$$\hat{Y} = \hat{\beta}_0 + \sum_{j=1}^p X_j \hat{\beta}_j \qquad (\not\preccurlyeq 1)$$

式1において、 \hat{Y} は予測変数、 $\hat{\beta}_0$ は切片、 X_i は説明 変数, $\hat{\beta}_i$ は各説明変数の重み (直線の傾き) を示してい る。この線形回帰モデルは、説明変数の効果が重み $\hat{\beta}_{1}$ と して表現され,事象に対する説明性が高く解釈も容易であ ることから広く利用されている。 R による線形回帰モデ ルを構築する lm() 関数では, 直線以外にも対数, 累乗等 の曲線を当てはめたモデルを構築することもできる。 一 方で,回帰直線・曲線に当てはまりが悪い事象は予測精度 が著しく悪くなることが知られており,正規分布に従わな い、大気汚染物質の観測データのようなノンパラメトリッ クデータはその傾向が大きい。 このことから,近年はラ ンダムフォレストやニューラルネットワーク等の機械学 習モデルを用いた解析により, 推論・予測が行われている 事例が増えている¹³⁾。高い精度が得られる機械学習モデ ルであるが, ランダムフォレスト, サポートベクトルマシ ン及びニューラルネットワーク等のモデルは、モデル内部 はブラックボックスと呼ばれ,説明変数がどのように作用 して予測を行っているのかは説明が難しいとされている 14)。 そのため,解釈などの説明性を求める場合は線形回 帰,精度の高い回帰や予測を行う合は機械学習を使う等, 目的に応じて使い分けられている。

2. 一般化加法モデルについて

ー般化加法モデル (Generalized Additive Model: GAM) は、もともと一般化線形モデル (Generalized Liner Model: GLM) で非線形性を扱うことができるよう改良されたモ デルであり、説明性と予測精度の両方を備えたモデルと言 われている¹⁴⁾。 式2に GAM のモデル式を示す。

$$Y = \alpha + \sum_{j=1}^{p} f_j(X_j) + \epsilon \qquad (\not \exists 2)$$

式2で、 ϵ は平均0の誤差項を表す。式1と異なるの は、説明変数がそれぞれ非線形な平滑化関数 f_j に適用さ れていることである。また、GAMでは1つの関数を全デ ータに対してそのまま適応させるのではなく、はじめにデ ータを複数の区間に分割し、それぞれの区間にパラメータ を調整した平滑化関数を設定し、各平滑化関数の合計値が 回帰モデルとなる回帰スプラインが代表的な手法となっ ている¹⁴⁾。区間数や規定関数を適切に選択することで、 どのようなデータに対しても回帰が可能と言われている。 GAM に関する詳細は参考文献^{12),15)}を参照してほしい。

3. 利用するデータの選定及び計算条件

Ox 生成の起源として,日射量,気温,風速及び VOC 等の前駆物質の存在が挙げられる^{16),17)}。九州では春先 4~ 5月にかけてこれらの気象条件が揃うことが多く,熊本県 で過去に光化学スモッグ注意報等を発令したのもこの時 期であった¹⁸⁾。本報で構築する Ox の回帰モデル構築の 説明変数には,信頼性の高いデータが手に入りやすい日射 量,気温及び風速を用いることとした。また,GAM では 時空間属性を含めてモデル化することが可能であり¹⁹⁾, 本報ではこれらの検討も行うことから,緯度,経度及び時

衣下 使用したナーダの一見					
利用データ	目的変数 • 説明変数名	概要			
常時監視測定局観測データ	0x	光化学オキシダント1時間値濃度 [ppb]			
地上付近気象庁気象データ(MSM-S)	temperature	気温 1 時間値[℃]			
	wind-speed	風速 1 時間値 [m/s]			
地上付近気象庁気象データ(MSM-r1h)	solar-radiation	日射量 1 時間値 [W/m^2]			
位置情報	lon	経度[°]			
	lat	緯度[°]			
時間	date	通算時間 [-]			
	time	0 ~ 23 [hour]			

1 は 田 ト キーゴー ク の 一 監

間 (通算時間及び日内変動時間 0~23 時) も説明変数と して準備した。 データの詳細は表1に示す。

学習及び予測を行う期間は、九州の広い範囲で Ox 濃度 が高濃度となり、熊本県でも注意報の発令が行われていた 2019 年 5 月 23 日 ~ 5 月 24 日を含む 10 日間とした。 また、予測地点には九州域を 5 × 5 [km]の格子領域に分 割した地点を設定した。 表 2 に計算条件の一覧を示す。

4. 解析を実行するためのツールの選択

解析には R version 4.4.1 (2024-06-14) 及び mgcv パッ ケージ (ver. 1.9.1)を用いた。 R は回帰モデルの構築が 容易であり、かつ、構築できるモデルの種類が多いことか ら、経済、医療、環境など様々な分野の解析に利用されて いる。 また、 mgcv パッケージは R で GAM を実行す るための標準的なパッケージであり、利用実績も多いこと から採択した。

以降では、GAM を用いて非線形モデルによる一般的な モデル (BASE),時間効果モデル (temporally vary coefficient model: TVC),空間効果モデル (spatially varying coefficient model: SVC)及び時空間効果モデル (Spatially and Temporally Varying Coefficient models: STVC)による GAM を構築する。 なお,比較として線形回帰モデル (Linear Regression model: LM)による回帰モデルも構築す る。

回帰モデルの構築と予測

1. 一般的な回帰モデル

時間や空間効果を検討する前に、本章では一般的な回帰 モデルの構築について解説する。

1.1	線形回帰モデルによる予測	



図1 宇土運動公園における 0x 濃度と日射量,気温及び 風速の関係。各変数の散布図及び時系列図。

R での回帰モデル構築について簡単に紹介する。回帰 モデル解説には、分かりやすいように宇土運動公園1局の データを使ってモデルを構築し、目的変数 (Ox) と説明変 数 (日射量、気温及び風速)の関係を確認する。宇土運動 公園局の Ox 濃度、日射量、気温及び風速については図1 に示す。図1によると、Ox 濃度と日射量及び気温は正の 相関関係が見られている。また、時系列データから日射 量、気温及び Ox 濃度が日内変動で関連した挙動を示して いることが示されている。

続いて, リスト 1 に線形回帰によるモデル化の結果 (m1)を示す。以下に,出力結果について簡単に解説する。

計算対象	概要	
学習データ	対象期間:	2019-05-21 0h ~ 2019-05-27 23h
	予測変数:	九州域の常時監視測定局 0x 1 時間値データ (0x)
	説明変数∶	日射量(solar-radiation)
		気温 (temperature)
		風速(wind-speed)
予測データ	対象期間:	2019–05–23 6h \sim 2019–05–24 23h
	計算領域:	5 × 5 [km](九州域)
	説明変数∶	日射量(solar-radiation)
		気温 (temperature)
		風速(wind-speed)

表 2 計算条件



図2 宇土運動公園における 0x 濃度の線形回帰分析結果 青破線:回帰直線 赤曲線:平滑化曲線

Call: には線形回帰モデルの構築方法が表示されている。 formula が実際のモデル式を表し,予測変数 Ox を説明変 数である日射量 (solar_radiation),気温 (temperature)及 び風速 (wind_speed) で説明している。

Coefficients: には線形回帰式の切片 (Intercept) と各説 明変数の重み (Estimate) が示されており,標準誤差 (Std.Error), t 値 (t value) 及び p 値 ($\Pr(>|t|)$) が更に続 いている。また,最後列には p 値の優位性 (Signif.codes) が記号 (*,.,) で示されている。 それ以降は決定係数 (*R*² 値: Multiple R-squared) や *F* 検定の結果から求めた *p* 値 (p-value) が示されている。

リスト1から,線形回帰では R² 値が 0.67 であり,予 測に用いた説明変数のうち気温及び風速の影響が大きか った。なお,日射量の影響は小さかった。

図2に線形回帰で得られたOxと各説明変数との関係を 示す。 図のx軸は説明変数の値,y軸は各説明変数によ るOxの予測値を示している。 また,青破線は線形回帰 により得られた回帰直線を示したもので,赤実線は平滑化 曲線を示す。 図2から,説明変数のうち気温は回帰直線 と平滑化曲線がほぼ一致しており,当てはまりが良い結果 であった。 一方で,日射量と風速はデータのばらつきが 大きく,直線性が見られず十分にモデル化されていなかっ た。

1.2 一般化加法モデルによる予測

次に,同じデータを用いて GAM による回帰の方法を示 す。GAM を使った計算には, mgcv パッケージの gam()

リスト2 宇土運動公園における 0x 濃度の一般化加法モデル結果

```
(m2 < -gam(formula = 0x \sim s(solar_radiation) + s(temperature) + s(wind_speed), data = na.omi
t(uto))) |> summary()
Family: gaussian
Link function: identity
Formula:
0x \sim s(solar_radiation) + s(temperature) + s(wind_speed)
Parametric coefficients:
           Estimate Std. Error t value Pr(>|t|)
(Intercept) 57.1617
                      0.9612 59.47 <2e-16 ***
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '' 1
Approximate significance of smooth terms:
                    edf Ref.df
                                   F p-value
s(solar_radiation) 2.210 2.732 1.676 0.227818
                 6.949 8.040 22.072 < 2e-16 ***
s(temperature)
s(wind_speed)
                  3.645 4.528 4.979 0.000663 ***
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '' 1
R-sq. (adj) = 0.743 Deviance explained = 76.3%
GCV = 168.19 Scale est. = 154.28 n = 167
```



図3 宇土運動公園における 0x 濃度の一般化加法モデル 結果

黒破線:回帰曲線 赤色領域:95%信頼区間



図 4 宇土運動公園における 0x 濃度の交互作用項を導入 した GAM 予測結果

黒破線:回帰曲線 赤色領域:95%信頼区間

関数を用いる。また、GAMではスプライン関数を適応させるために s() 関数で説明変数を囲ってモデル式を構築 する。モデル化した GAM の結果 (m2) をリスト2 に示 す。線形回帰と似た出力となっているが、切片や説明変 数の重みを示す coefficients: はパラメトリックとノンパ ラメトリック (Approximate significance of smooth terms:) に分けて表示されている。また、結果に逸脱度の説明性 (Deviance explained) [%] と一般化クロスバリデーション (Generalized Cross Validation: GCV) 誤差が表記されている 点は、線形回帰の結果の表示と異なっている。ここで、 逸脱度の説明性 (Deviance explained) とは、説明変数を用 いずに計算した際の逸脱度 (Null Deviance) と説明変数を 用いて計算した際の逸脱度 (Residual Deviance) との比を とっており、値が大きいほどモデルの適合度は高いと言え る (式 3)。

表3 説明変数間の相関係数

	solar_radiation	wind_speed	temperature
solar_radiation	1	0. 4575	0. 7593
wind_speed	0. 4575	1	0. 4405
temperature	0. 7593	0. 4405	1

Deviance Explained =

$$\left(1 - \frac{Null \ Deviance}{Residual \ Deviance}\right) \times 100 \qquad (\equiv 3)$$

m2 の結果では, R^2 値は 0.74 であり,線形回帰より も良い結果であった。 また,逸脱度の説明性は 76.3 [%] であり,総合的には m1 よりも m2 が良い結果となって いた。

次に,説明変数の効果を図3に示す。 図の黒実線は回 帰曲線,赤色領域は95%信頼区間を表している。 これに よると,気温は25[℃],風速は2[m/s]付近でOx濃度 が高い関係性が見られた。 説明変数の効果についても線 形回帰の結果と異なっていた。

1.3 説明変数の選択と交互作用について

交互作用とは,説明変数間で何らかの関係性が見られる ことを指し,この交互作用をモデルに組み込み事で予測精 度が向上することがある。GAM では説明変数の交互作用 を適応させた GA2M が考案されており,mgcv パッケー ジでは説明変数の交互作用を標準で扱うことができる²⁰⁾。 本章では,GAM の交互作用の確認方法と説明変数の選択 方法について確認する。

はじめに,交互作用の確認のため説明変数間の相関係数 を確認する(表3)。変数間に相関が見られる場合,交互 作用がある可能性は高い。表3によると、日射量、気温 及び風速の間には一定以上の相関があり,特に日射量と気 温の間に強い相関が見られた。 このように,説明変数間 で交互作用が見られる可能性が高いことから, m2 モデ ルにすべての変数同士の交互作用を追加したモデル (m3) を定義した。 m3 の結果をリスト3及び図4に示す。 交 互作用を定義していない m2 モデルと比べ, 交互作用を 定義した m3 は R² 値及び逸脱度の説明性が向上してお り, m3 は m2 より明らかに良い結果となった。一方で, m3 の説明変数の中には m2 より p 値が低く, 精度が悪 くなる変数も見られていた。 本検討では上記手順のよう に全ての説明変数間の組み合わせに相関が見られたため 全ての組み合わせを適応させた計算を行ったが、本来、交 互作用を検討する場合は、計算負荷軽減のためにも、事前 に関連性の高い説明変数の組み合わせを検討し、その中か

> ら特に効果が期待される組み合わせを取捨選択 して計算させるべきであり、このようにすること で回帰モデルの予測精度向上が期待できるよう になると思われる。

一方で,上記のように説明変数が複数存在する 場合に説明変数間の関係を個別に確認すること は,時間に余裕がある場合を除き非効率的であり, 期待どおりの結果が得られない場合もある。そ リスト3 交互作用を追加した GAM の構築

```
(m3 \leftarrow gam(formula = 0x \sim s(solar_radiation) + s(temperature) + s(wind_speed) +
            s(solar_radiation, temperature) + s(temperature, wind_speed) + s(solar_radiatio
n, wind_speed), data = na.omit(uto))) |> summary()
Family: gaussian
Link function: identity
Formula:
0x \sim s(solar_radiation) + s(temperature) + s(wind_speed) + s(solar_radiation,
   temperature) + s(temperature, wind_speed) + s(solar_radiation,
   wind_speed)
Parametric coefficients:
           Estimate Std. Error t value Pr(>|t|)
(Intercept) 57.1617 0.7618 75.03 <2e-16 ***
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '' 1
Approximate significance of smooth terms:
                                    edf Ref.df
                                                 F p-value
s(solar_radiation)
                              5.585e+00 5.604 0.022 0.99994
                              1.000e+00 1.000 5.224 0.02401 *
s(temperature)
s(wind_speed)
                              5.526e+00 6.630 3.864 0.00121 **
s(solar_radiation, temperature) 1.071e-06 27.000 0.000 0.22728
s(temperature, wind_speed)
                             1.279e+01 27.000 2.800 < 2e-16 ***
s(solar_radiation, wind_speed) 1.951e+01 27.000 2.540 < 2e-16 ***
____
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '' 1
R-sq. (adj) = 0.839 Deviance explained = 88.2%
GCV = 133.12 Scale est. = 96.923 n = 167
```

のため、ステップワイズ法のように説明変数の取捨選択を 自動化する手法も広く利用されており、計算時間の短縮化 も期待できる。 R の mgcv パッケージでは、GAM のス テップワイズ法は定義されていないため、本報では単純な GAM の前方ステップワイズ法を新たに定義し、説明変数 の自動化による取捨選択の手法を確認することとする。 なお、説明変数の取捨選択の判断は BIC 基準²¹⁾として実 行した (リスト4;全出力は文末に記載)。

リスト4の結果を確認すると、気温及び風速を説明変数 としたモデルが BIC 基準で最も当てはまりの良い結果 である結果となった。一方で、 R² 値及び逸脱度の説明 性を基準とすると、気温、風速及び日射量と風速の交互作 用項を説明変数とすると, m3 よりもモデルの当てはま りが良くなっていた。 これは BIC 基準がパラメータの 数 (説明変数の数) に対してペナルティを取る方法に起 因していると思われ,ステップワイズ法を適用する場合は 基準とする統計量の選択にも注意を払う必要がある。 な お, R の step() 関数にも使われている AIC 基準で試し たところでは,全ての説明変数及び交互作用を含めること が,最も良い当てはまりとなっていた。

1.4 交差検証法による精度管理

続いて,モデルの精度検証に広く利用されている交差検 証法 (クロスバリデーション)を用いた検証を行う。 交 リスト4 前方ステップワイズ法による交互作用項を導入した GAM の構築(結果のみ表示)

```
Family: gaussian
Link function: identity
Formula:
0x \sim s(temperature) + s(wind speed)
Parametric coefficients:
           Estimate Std. Error t value Pr(>|t|)
(Intercept) 57.1617
                     0.9672
                                59.1 <2e-16 ***
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '' 1
Approximate significance of smooth terms:
               edf Ref.df
                             F p-value
s(temperature) 7.468 8.429 38.812 < 2e-16 ***
s(wind_speed) 3.870 4.797 5.983 8.29e-05 ***
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '' 1
R-sq. (adj) = 0.74 Deviance explained = 75.8%
GCV = 168.67 Scale est. = 156.21
                                  n = 167
```

差検証法では、予測に用いるデータを学習用と予測用に分け、学習用データのみを使って構築したモデルで予測用デ ータを再現する作業を行う。 一般的な交差検証法では、 データを10分割し、 9/10 のデータを用いて回帰モデル を構築し、残り 1/10 のデータを予測する。この計算を全 ての組み合わせ (10 パターン) に対して行う (10-foldcross-validation)。また、分割数を k で表して k-fold-crossvalidation とも呼ばれている。

リスト 5 に新たに GAM 用の交差検証法を定義した。 この関数を用いて,ここまでに構築したモデル ml, m2 及び m3 の予測精度評価を行った。 交差検証法の評価に は *R²* 値及び *RMSE* 値を用いた。

図 5 に交差検証法による精度評価の結果を散布図で示 す。 この図より、3 モデルの中では m3 で R² 値が最も 大きく、 RMSE 値が最も小さい値となり、交互作用項を 導入したモデルの予測精度が最も良い結果となった。 一 方で、 線形回帰モデルの m1 は、3 モデルの中で最も予 測精度が悪い結果となった。 このように、交差検証法を 用いると予測精度を簡単に比較することができ、ここまで に構築した m1 ~ m3 では、GAM は線形回帰よりも予測 精度が高いことが示された。

1.5 モデルの再構築

前節までの検討により, 宇土運動公園 1 局の観測データ を用いた回帰分析では, 線形回帰モデルよりも GAM の予 測精度が高く, かつ, 回帰の説明性が高いことが示された。 次節からは, 九州全域の観測データを用いて回帰モデルを 構築・比較していくため, すでに構築したモデルの m1 及 び m2 を, 九州全域のデータを使ってモデルを再構築し た (m1:LM, m2:BASE, リスト 6)。 次節以降では新しい m1 (LM) 及び m2 (BASE) を用いて, 時間及び空間属性を 導入した GAM との比較を行う。

なお, m3 で構築した交互作用項ありのモデルについ ては,次節からは扱うデータの次元が増えることで回帰モ デルの構築の難易度が上がることが予想されるため,これ 以上の検討は本報では扱わないこととする。検討方法は 他の研究論文等を参照してほしい²²⁾。

時間効果を考慮した GAM

前章では、一般的な事象である、無次元のデータを用いた GAM の構築について検討した。本章からは、対象期間の九州全域の常時監視測定局 Ox 濃度 1 時間値を用いて、時間-空間変動を考慮した GAM の構築について検討

```
リスト5 交差検証法の定義
```

```
cross_validation <- function(data = NULL,</pre>
                                  formula = NULL,
                                  nfold = 10,
                                  rep = 10, ...) {
  stopifnot (nfold > 0)
  nfold <- floor(nfold)</pre>
  num \langle - \text{ seq}(1, \text{ rep, by } = 1)
  apply(do.call(cbind, lapply(num, function(n) {
    id <- sample(nrow(data))</pre>
    fold_list <- split(id, ceiling(seq_along(id) / floor(length(id) / nfold)))</pre>
    ret \langle - \text{ matrix}(\text{NA}, \text{ nrow} = \text{length}(\text{id}), \text{ ncol} = 1)
    for (fn in fold_list) {
       fit model <- mgcv::gam(formula = formula, data = data[-fn, ], ...)
       ret[fn] <- predict(fit_model, newdata = data[fn, ])</pre>
    }
    return(ret)
  })), 1, mean)
RMSE \langle - function(x = double(),
                   y = double()) \{
  round(sqrt(mean((x - y) \hat{2}, na.rm = TRUE)), digits = 2)
R2 \leq function(obs = double()),
                 predict = double()) {
  round(1 - mean((obs - predict) ^ 2, na.rm = TRUE) / var(obs, na.rm = TRUE), digits = 2)
```

リスト6 m1 (LM) 及び m2 (BASE) モデルの再構築

```
m1 <- lm(formula = 0x ~ solar_radiation + temperature + wind_speed, data = Kyushu)
m2 <- gam(formula = 0x ~ s(solar_radiation) + s(temperature) + s(wind_speed), data = Kyushu)</pre>
```

リスト7 時間効果を導入した GAM(m4)の構築

```
m4 <- gam(formula = 0x ~ s(time, bs = "gp") + s(date, bs = "gp"),
data = Kyushu)
```



図5 宇土運動公園のデータを用いた交差検証法による精 度管理の結果 青線:回帰直線

する。本節では時間効果を考慮した GAM (TVC) の検討 を行う。

GAM の構築には、引き続き mgcv パッケージの gam() 関数を用いる。 mgcv パッケージで時間変動を扱うため には説明変数に時間情報を追加する必要があり、本報で は通算時間 (date) による s(date, bs = "gp"), または、日 内変動時間 (0~23時: time) による s(time, bs = "gp") を モデル式に追加する必要がある。 ここで bs = "gp" はス プライン関数にガウスモデルを適用させる方法であり、



図 6 時間効果を用いた GAM (m4) の結果 黒破線:回帰曲線 赤色領域:95%信頼区間 コンター図:x 軸及び y 軸の説明変数の関係が見られる場合の 0x 濃度の平均濃度 [ppb]

```
リスト8時間効果及びその他説明変数を用いた GAM (m5)の構築
```

```
(m5 \leq gam(formula = 0x \sim s(time, bs = "gp") + s(date, bs = "gp") + s(solar_radiation) + s(t)
emperature) + s(wind_speed),
          data = Kyushu)) |> summary()
Family: gaussian
Link function: identity
Formula:
0x \sim s(time, bs = "gp") + s(date, bs = "gp") + s(solar_radiation) +
   s(temperature) + s(wind_speed)
Parametric coefficients:
           Estimate Std. Error t value Pr(>|t|)
                                693 <2e-16 ***
(Intercept) 57.56980 0.08307
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '' 1
Approximate significance of smooth terms:
                   edf Ref. df F p-value
s(time)
                  3.230 3.528 902.01 <2e-16 ***
s(date)
                 4.743 4.959 1939.84 <2e-16 ***
s(solar radiation) 3.921 4.831 55.97 <2e-16 ***
s(temperature) 5.179 6.363 651.66 <2e-16 ***
s(wind_speed)
                 8.552 8.943 87.63 <2e-16 ***
____
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '' 1
Rank: 37/50
R-sq. (adj) = 0.708 Deviance explained = 70.8%
GCV = 161.4 Scale est. = 161.22 n = 23361
```





コンター図:x 軸及び y 軸の説明変数の関係が見られる場合の 0x 濃度の平均濃度 [ppb]

このオプションをつけることで時間及び空間効果が考慮 されるようになる¹⁹⁾。

はじめに,説明変数に時間効果のみで構築したモデル (リスト7,図6:m4)について検討する。なお,図6には, これまで同様 plot()関数による説明変数(単独または複 数)毎のOx濃度の予測値に加え,vis.gam()関数による, 2つの説明変数間の関係をコンター図にして示している。

m4 の結果から,24時間のうちでOx 濃度が最も高くなる時間帯は15時付近,通算時間では3日~4日目付近 (2019/5/23~5/24)であり,図1に示した時系列図と一致する結果が得られた。

この m4 モデルの説明変数に日射量,気温及び風速を 追加したモデルの結果をリスト 8 及び図 7 に示す (m5)。 m5 では,日射量が 600 [W/m²] 付近で最低値となってい た。気温は 20 [°C] 以上で Ox 濃度との関係が最大になっ ていた。風速は 2~6 [m/s] 付近でプラスであり,それ以 上風速が早くなっても影響度はほぼ変わらなかった。 時 間情報と日射量,気温及び風速の関係では,日射量が多く, 気温も高くなりやすい 15 時~20 時の間で Ox 濃度が高く なる傾向を示していた。 なお,風速は 15 時~20 時の間や 5/3~5/4 付近で,6 [m/s]の風があるときに Ox 濃度が最大 値となっていた。

次に,日射量,気温及び風速を時間情報の因子とするモ デルを構築する。 この場合,説明変数を s(date, bs = "gp", by = temperature)のように記述して GAM を構築する²³⁾。 by = temperature 部分が気温のデータを時間情報の因子と して取り込む設定を表している。 なお,説明変数の冒頭 に 0 + を入れると,計算の過程で自動的に決定される切 片係数を明示的に初期化し,各説明変数(本報の場合,日 射量,気温及び風速)の値を切片係数に割り当てるように なる²³⁾。

日射量及び気温を日内変動時間(0~23時),風速を通 算時間の変動モデルとして構築した結果をリスト9及び 図8に示す(m6)。この結果,各説明変数の効果がm6 ではm5に比べて時間情報に適合しており,モデルが改 善されていることが分かる。また,m6の逸脱度の説明 性は95.7 [%]で,m5の70.8 [%]に比べて大幅に向上 していた。 リスト9時間効果を説明変数の因子として導入した GAM (m6)の構築

```
(m6 <- gam(formula = 0x \sim 0 +
            solar_radiation + s(time, bs = "gp", by = solar_radiation) +
            temperature + s(time, bs = "gp", by = temperature) +
            wind_speed + s(date, bs = "gp", by = wind_speed), data = Kyushu)) |> summary()
Family: gaussian
Link function: identity
Formula:
0x \sim 0 + solar_radiation + s(time, bs = "gp", by = solar_radiation) +
   temperature + s(time, bs = "gp", by = temperature) + wind_speed +
   s(date, bs = "gp", by = wind_speed)
Parametric coefficients:
               Estimate Std. Error t value Pr(>|t|)
solar_radiation -0.01999 0.02369 -0.844 0.39868
              0.83077
                          0.28834 2.881 0.00397 **
temperature
wind_speed
               25.72510
                          3.09354 8.316 < 2e-16 ***
____
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '' 1
Approximate significance of smooth terms:
                        edf Ref.df
                                    F p-value
s(time):solar_radiation 3.964 4.595 108 <2e-16 ***
                     4.496 5.081 539 <2e-16 ***
s(time):temperature
                     4.709 5.220 1755 <2e-16 ***
s(date):wind_speed
____
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '' 1
Rank: 26/39
R-sq. (adj) = 0.696 Deviance explained = 95.7%
GCV = 167.99 Scale est. = 167.88
                                 n = 23361
```

以上の結果から,時間変動を考慮したモデルでは m6 が最も良い結果となったため,以降では m6 を TVC の 代表として比較を行う。

補足だが,日内変動時間のように循環する数値情報を扱う場合,GAMのモデル式では bs = "cc" とすることで, スプラインが循環する情報を扱うことができるようになる²⁰⁾。このように設定するほうが良い結果となる場合も あるため,適用するスプラインや説明変数の関係を詳しく 見ながらモデルを構築する必要がある。

3. 空間効果を考慮した GAM

本節では、空間効果を取り入れた GAM (SVC)の構築を

行う。 位置情報の経度 (lon) と 緯度 (lat) を交互作用と
 し、スプラインはガウシアンモデルを設定して s(lon, lat, bs = "gp") のようにしてモデルを構築すると、空間効果を
 取り入れたモデルができる。

まず, リスト 10 及び図 9 に空間効果のみで構築した GAM の結果 (m7) を示し, Ox 濃度との関連度を確認す る。m7 の結果 (図9) では, X - Y 平面上に経度及び緯 度の効果が等高線で示されている。 これによると, 福岡 県の北九州市周辺や大分県の一部で Ox 濃度が高くなる傾 向にあり, 長崎県や鹿児島県では Ox 濃度が低い傾向にあ った。 なお, 海上で Ox の生成にプラスの影響が見られ ている領域は, もともとデータが存在しない範囲であるた





リスト10 空間効果を導入した GAM (m7) の構築

m7 <- gam(formula = 0x $^{\sim}$ s(lon, lat, bs = "gp "), data = Kyushu)

めこのようになっているが、本検討では海上の予測は行 わないため検討は省略する。



図 9 空間効果を導入した GAM (m7) の結果 黒破線:回帰曲線 赤色領域:95%信頼区間 コンター図:x 軸及び y 軸の説明変数の関係が見られる 場合の 0x 濃度の平均濃度 [ppb] リスト11 空間効果を説明変数に追加した GAM (m8)の構築

```
(m8 \leftarrow gam(formula = 0x \sim s(lon, lat, bs = "gp") + s(solar_radiation) + s(temperature) + s(w \sim s(temperature)) + s(w \sim s(temperature))
ind_speed), data = Kyushu)) |> summary()
Family: gaussian
Link function: identity
Formula:
0x \sim s(lon, lat, bs = "gp") + s(solar_radiation) + s(temperature) +
    s(wind_speed)
Parametric coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 57.5698 0.1017 565.9 <2e-16 ***
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '' 1
Approximate significance of smooth terms:
                      edf Ref.df
                                       F p-value
                   31.702 31.965 60.43 <2e-16 ***
s(lon,lat)
s(solar_radiation) 8.875 8.995 36.08 <2e-16 ***
                  8.833 8.992 1359.28 <2e-16 ***
s(temperature)
s(wind_speed)
                   5.626 6.823 82.29 <2e-16 ***
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '' 1
R-sq. (adj) = 0.562 Deviance explained = 56.3%
GCV = 242.39 Scale est. = 241.81 n = 23361
```

続いて, m7 に日射量, 気温及び風速を説明変数として 追加した結果をリスト 11 及び図 10 に示す (m8)。m8 で は,経度 131 度~131.5 度付近及び緯度 33.5 度以上での領 域で,日射の強さに関係なく Ox 濃度が高くなる傾向にあ った。気温及び風速効果は,位置情報に関係なく,気温 が高く風速が遅いほど Ox 濃度が高くなる傾向にあった。 また, m8 の逸脱度の説明性は 56.3 [%] であり, m6 と 比べて極端に低い値となっていた。

最後に、日射量、気温及び風速を空間情報の因子に設定 して構築したモデルの結果をリスト 12 及び図 11 に示す (m9) 。 m9 結果では、 m8 の結果に比べて日射量及び風 速の効果がより現実的なものとなっており、 m8 で見ら れた極端な地域性は見られなかった。 なお、 m9 の逸脱 度の説明性は 93.1 [%] となっており、 m7 や m8 と比べ て良い結果であった。

以上の結果から,空間変動を考慮したモデルでは m9 が最も良い結果となり,以降では m9 を SVC の代表と

して比較を行う。

4. 時空間効果を考慮した GAM

最後に、本節では時空間効果を考慮したモデル (STVC) を構築する。 時空間効果によるモデル化では、位置情報 lon 及び lat に加え、時間情報 date または time を交互 作用とし、 s(lon, lat, date, bs = "gp") のようにモデルを構 築する^{19) 23)}。

まずは時空間情報のみで構築したモデルを構築し (m10),その結果をプロットする(リスト13及び図12)。 なお、これまで各説明変数による Ox 濃度の予測値をプロ ットしていたが、時空間効果の情報量が多くなり、結果の 表示が見にくくなってしまうため、本節では vis.gam()関 数だけで結果を表示する。



図 10 空間効果を説明変数に追加した GAM (m8) の結果 黒破線:回帰曲線 赤色領域:95%信頼区間 コンター図:x 軸及び y 軸の説明変数の関係が見られる場合の 0x 濃度の平均濃度 [ppb]

リスト12 空間効果を説明変数の因子として導入した GAM (m9)の構築

```
(m9 \leq -gam(formula = 0x \sim 0 +
            solar_radiation + s(lon, lat, bs = "gp", by = solar_radiation) +
            temperature + s(lon, lat, bs = "gp", by = temperature) +
            wind_speed + s(lon, lat, bs = "gp", by = wind_speed), data = Kyushu)) |> summar
y ()
Family: gaussian
Link function: identity
Formula:
0x \sim 0 + solar_radiation + s(lon, lat, bs = "gp", by = solar_radiation) +
   temperature + s(lon, lat, bs = "gp", by = temperature) +
    wind_speed + s(lon, lat, bs = "gp", by = wind_speed)
Parametric coefficients:
                Estimate Std. Error t value Pr(>|t|)
solar_radiation 0.04438 0.09031 0.491
                                              0.623
temperature
               16. 20156 10. 78568 1. 502
                                              0.133
wind_speed
               -70.56022 74.90184 -0.942
                                              0.346
Approximate significance of smooth terms:
                           edf Ref.df
                                          F p-value
s(lon, lat):solar_radiation 26.60 29.02 8.614 <2e-16 ***
s(lon, lat):temperature 30.75 31.44 20.808 <2e-16 ***
s(lon,lat):wind_speed
                        30.42 31.14 16.827 <2e-16 ***
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '' 1
Rank: 99/102
R-sq. (adj) = 0.514 Deviance explained = 93.1%
GCV = 270.07 Scale est. = 269.04 n = 23361
```



図 11 空間効果を説明変数の因子として導入した GAM (m9) の結果 黒実線:回帰曲線

コンター図:x 軸及び y 軸の説明変数の関係が見られる場合の 0x 濃度の平均濃度 [ppb]

リスト13 時空間効果を導入した GAM (m10)の構築

m10 <- gam(formula = 0x $^{\sim}$	s(lon,	lat,	time,	bs =	″gp″)	+ s(lon,	lat,	date,	bs =	″gp″),	data
= Kyushu)											





m10 の結果では、時間及び空間変動効果は m6 (TVC) 及び m9 (SVC) と同様の結果となっていたが、空間変動 において m8 及び m9 で見られた地域差がほぼなくな り、北東から南西にかけて徐々に Ox 濃度が低くなる傾 向になっていた。

次に ml1 に日射量, 気温及び風速を説明変数として追加したモデルの結果をリスト 14 に示す。 ml1 では, 各説明変数の p 値が 0.01 以下で当てはまりは良かったが, モデル結果のプロット (vis.gam()) がエラーで出力されなかった。そのため, ml1 のモデルの評価は省略することとする。

最後に、日射量、気温及び風速を時空間効果の因子に設

定して構築したモデルの結果をリスト 15 及び図 13 に示 す (m12)。 なお, m12 ではプロットの数が多いため, 時間属性を上2段,空間属性を下2段に分けて表示してい る。m12 の時間変動は m6 とほぼ同じ結果となっていた。 また,空間変動は m9 とは異なり,空間依存性をほとんど 示さなっていた。 日射量のみ空間的な変動が確認された が,空間効果は小さかった。

なお, ml2 モデルの R² 値は 0.72 , 逸脱度の説明性 は 96.0 [%] であり,これまでに構築したモデルの中でも 良い結果となった。

このように、GAM では空間効果より時間効果の影響度 が強く現れる結果となった。

リスト14 時空間効果を説明変数に追加した GAM (m11)の構築

```
(m11 \leq gam(formula = 0x \sim s(lon, lat, date, bs = "gp") + s(lon, lat, time, bs = "gp") + s(s)
olar_radiation) + s(temperature) + s(wind_speed),
           data = Kyushu)) |> summary()
Family: gaussian
Link function: identity
Formula:
0x \sim s(1on, lat, date, bs = "gp") + s(1on, lat, time, bs = "gp") +
    s(solar_radiation) + s(temperature) + s(wind_speed)
Parametric coefficients:
           Estimate Std. Error t value Pr(>|t|)
(Intercept) 57.56980 0.08764 656.9 <2e-16 ***
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '' 1
Approximate significance of smooth terms:
                    edf Ref.df
                                    F p-value
s(lon, lat, date)
                  6.963 6.999 1534.95 <2e-16 ***
s(lon, lat, time)
                  2.000 2.000
                               64.51 <2e-16 ***
s(solar_radiation) 8.897 8.997 132.95 <2e-16 ***
                 6.849 7.980 1610.95 <2e-16 ***
s(temperature)
s(wind speed)
                  7.486 8.441 68.40 <2e-16 ***
____
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '' 1
Rank: 35/131
R-sq. (adj) = 0.675 Deviance explained = 67.6%
GCV = 179.69 Scale est. = 179.45 n = 23361
```

リスト15時空間効果を説明変数の因子として導入したGAM(m12)の構築

```
(m12 <- gam(formula = 0x ^{\sim} 0 +
            solar_radiation + s(lon, lat, time, bs = "gp", by = solar_radiation) +
            temperature + s(lon, lat, time, bs = "gp", by = temperature) +
            wind_speed + s(lon, lat, date, bs = "gp", by = wind_speed), data = Kyushu)) |>
summary()
Family: gaussian
Link function: identity
Formula:
0x \sim 0 + solar_radiation + s(lon, lat, time, bs = "gp", by = solar_radiation) +
   temperature + s(lon, lat, time, bs = "gp", by = temperature) +
   wind_speed + s(lon, lat, date, bs = "gp", by = wind_speed)
Parametric coefficients:
               Estimate Std. Error t value Pr(>|t|)
solar_radiation 0.8676
                          1.4072 0.617 0.5375
temperature
                 9.8369
                            4.1690 2.360 0.0183 *
                7.7110
                           0.7122 10.827 <2e-16 ***
wind_speed
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '' 1
Approximate significance of smooth terms:
                                 edf Ref.df
                                                 F p-value
s(lon, lat, time):solar_radiation 31.042 32.086 25.11 <2e-16 ***
s(lon, lat, time):temperature
                             13.848 14.336 180.45 <2e-16 ***
s(lon, lat, date):wind speed
                               7.917 7.924 1307.07 <2e-16 ***
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '' 1
Rank: 58/315
R-sq. (adj) = 0.72 Deviance explained =
                                           96%
GCV = 154.97 Scale est. = 154.61 n = 23361
```

以上の結果から,時空間変動を考慮したモデルでは m12 が最も良い結果となり,以降では m12 を TSVC の代表 として比較を行う。

5. 最適な回帰モデル式の選択

これまで構築した m1 (LM), m2 (BASE), m6 (TVC), m9 (SVC) 及び m12 (STVC) について比較検証を行う。 モデルの比較には, *R²* 値, *RMSE* 値, *AIC*, *BIC* 及 び逸脱度の説明性 (Dev.exp) を用いる。 また,予測精度 の比較には、 10-fold-cross-validation で求めた R^2 値及び *RMSE* 値を用いる。

各モデルの統計量を表4に示す。 この結果によると, *R*² 値, *RMSE*, *AIC* 及び *BIC* は, m6(TVC) 及び m12 (SVC) が他のモデルに比べて良い結果となっていた。 ま た,逸脱度の説明性については, m2(BASE) が最も低く, それ以外の m6(TVC), m9(SVC) 及び m12(STVC) で高 い値を示していた。

続いて、交差検証法による予測精度検証の結果を図14



図 13 時空間効果を説明変数の因子として導入した GAM (m12)の結果 黒破線:回帰曲線 赤色領域:95%信頼区間

コンター図:x 軸及び y 軸の説明変数の関係が見られる場合の 0x 濃度の平均濃度 [ppb]

表 4 GAM の予測	则精度の比較 R2 値	(R.sq), AIC,	BIC 及ひ逸脫度の) 祝明性 (Dev. exp)
	R. sq AIC		BIC	Dev. exp
m1 (LM)	0. 5016	197540.5	197580. 8	_
m2 (BASE)	0. 5265	196364.1	196565.3	0. 5269
m6 (TVC)	0. 6961	185996.5	186122.3	0.9566
m9 (SVC)	0. 5136	197087.9	197815.3	0. 9307
m12 (STVC)	0. 7201	184112.0	184554. 0	0. 9601

表46	iAMの予測精度の比較	R2	値	(R. sq)		AIC,	BIC及び逸脱度の説明性	(Dev. ex	p
-----	-------------	----	---	---------	--	------	--------------	----------	---





に示す。これらの散布図では, y 軸に観測値, x 軸に交差 検証法の結果を示しており, データの重なり部分の色を変 えて表示できるよう密度分布で示している。

図14によると、m6(TVC)及びm12(STVC)で精度が 高い結果となり、それ以外のm1(LM),m2(BASE)及び m9(SVC)とは R²値及びRMSE値に大きな差が見られ た。m6(TVC)及びm12(STVC)では、データの密度が 大きい領域が回帰直線付近に集まっており、かつ、線形に 伸びていることから、これらのモデルの予測精度が高いこ とが示された。なお、m6(TVC)とm12(STVC)との予 測精度の差は小さく、空間変動の効果は時間変動の効果に 比べて小さいことが示された。しかしながら、時間効果 だけでは説明ができなかった不確実性を空間効果が補っ てGAMが構築されていると思われるため、時間効果単独 よりも時空間効果を取り込んだモデルを構築することで モデルの説明性が高くなっていると思われる。 これらの結果から,時空間効果を組み込んだ ml2 (STVC)の精度及び回帰の説明性が,これまで構築したモ デルの中で最も良い結果であることが示された。

6. 時空間データの予測

m12(STVC)を使って空間濃度分布の予測を行う。表2 に示す条件で予測地点のデータを作成し, R の predict.gam() 関数で空間濃度分布の予測計算を行った (図 15)。結果は 2019-05-23 6h~2019-05-24 23h までの時 間帯について示している。 なお, 図 15 では m12(STVC) による Ox 空間濃度分布予測値を地図上にコンター図とし て表示し, 更に, 常時監視測定局の位置及び観測値をコン ター図に重ねて表した。

図 15 から, ml2(STVC)の予測値は観測値を概ね再現 できていることが示された。特に,Ox 濃度の予測値が低 濃度になりやすい夜間~早朝の時間帯の予測精度につい



図 15 m12 (STVC) モデルを用いた空間濃度分布の予測結果 ポイントデータ:常時監視測定局の位置及び観測値 コンター図:m12 モデルを用いた 0x の空間濃度分布

ては高い予測精度となっていた。 一方で,日中に Ox 観 測値が 120 [ppb] を超過するような極端に高い濃度を観測 に既に Ox 濃度が上昇している時間, 2019/5/24 0h 付近の

していた時間, 2019/5/23 6h 付近で日射量が少ない時間帯

ように日中に上昇した Ox 濃度が夕方~夜間でも低下しな かった時間帯について, ml2 (STVC) では十分な予測が できなかった。

予測が十分ではなかった理由として,以下の2つの要因 が考えられる。1つ目の要因は,ml2のモデルに組み込 んだ通算時間 (date)の変動が実際のOx 濃度の変動を捉 えきれていなかったことである。m6~m9において, date の効果は3~4日目の濃度が最も高くなるように,な だらかな曲線としてモデル化がされていた。 観測値であ れば日内変動が見られるように,Ox 濃度の増減が表され ているはずであるが,GAMを構築する際の設定でこのよ うな結果となった。 機械学習では,予測を行う場合に学 習データからの過学習が度々問題になっていることを考 えると,この程度の予測誤差は起こりうることと思われる。 一方で,日射量や気温の変動は毎日同じでは無いため,時 間効果を time として構築した GAM では,時間変動の情 報だけでは極端な観測値の予測は難しいと思われる。

2つ目の要因は、説明変数とし使用した日射量、気温及 び風速では説明できない事象が実際に起きていたことで ある。 Ox の生成に大きく影響している日射量と気温は, 日の出から日の入りの時間帯に大きく依存しており、日の 入り後の時間は日射の影響がないことから、Ox の生成量 が0になる。 更に, Ox の主成分である O3 は, NO によ る分解で減少することが知られている²⁴⁾。そのため、Ox の国内における生成量は日射量に大きく影響を受けてい ると言えるが、国外で生成された Ox、又は Ox の前駆物 質が国内に流入する現象である越境移流が起きていた場 合,日射量,気温,風速及びこれらの相互関係からだけで は説明することができない濃度の変動を示すこととなり, 通常であれば Ox 濃度が低下している時間帯に Ox 濃度が 上昇する可能性がある。 本報の計算対象期間では、夜間 に観測値で Ox 濃度が高い時間及び場所が確認でき、国外 からの越境移流があった可能性が高い。 これらの事象の 詳細な解析は本報の検討内容から逸脱するため追加検討 は行わないこととするが、 越境移流を説明することがで きる説明変数には、化学輸送モデルによるシミュレーショ ンの結果等,明らかに国外からの影響と識別することがで きるデータに限られている ⁶⁾。他にも、回帰の結果に地 球統計学手法である Kriging を組み合わせて、予測誤差を 補間する手法の検討が行われており 11, この手法を取り 入れて予測精度を高めることも改善策の一つと言える。

まとめ

GAM により構築される予測精度の高いモデルは、その 結果から多くの情報を得ることができ,説明変数及び予測 変数間の関係性を詳しく確認することができた。 構築し たモデルの中では,時空間効果を導入したモデルで最も高 い予測精度及び説明性を示していたが,説明性が高い反面, m12(STVC)による空間濃度分布の予測結果のように,空 間的な変動は観測値と異なる部分もあった。そのため, 計算結果をそのまま受け入れるのではなく,GAMの結果 を確認し,必要に応じてモデルを改善していく必要がある ことが示された。加えて,本報では時空間効果を定義する 際に,説明変数を因子として指定する方法が最も精度が高 い結果となっていたが,m5,m8及びm11のように " 説明変数 + 時空間効果の交互作用"とする方法が良い場 合もある。予測変数の種類でもモデルの構築方法は異な ってくるため,最低限,本報で解説した方法を試して欲し い。

文 献 (MS ゴシック 9pt)

- e-Gov ポータル. <u>https://elaws.e-gov.go.jp/</u> (2024年7月閲覧).
- 環境省 : 大気汚染防止法第22条の規定に基づく 大気の汚染の状況の常時監視に関する事務の処理基 準

 <u>https://www.env.go.jp/air/osen/law22_kijun.htm</u>
 (2024年7月閲覧).
- 環境省 : 令和4年度 大気汚染状況について. <u>https://www.env.go.jp/press/press_03287.html</u> (2024年7月閲覧).
- 菅田誠治,大原利眞,黒川純一,早崎将光:大気
 汚染予測システム(VENUS)の構築と検証, J. Jpn.
 Soc. Atmos. Environ., 46(1), 49-59 (2011).
- 山村由貴,新谷俊二,力寿雄,中川修平,王哲,鵜 野伊津志:夏季の太平洋高気圧条件下における高 濃度 PM2.5 に対する火山の寄与解析, J. Jpn. Soc. Atmos. Environ., 55(4), 169-180 (2020).
- 古澤尚英,板橋秀一,豊永悟史,村岡俊彦: 2014 年冬季の熊本県中心部における微小粒子状物質の発 生源感度解析, J. Jpn. Soc. Atmos. Environ., 53(5), 194-205 (2018).
- 7) 国立環境研究所:地方環境研究所等との共同研究. <u>https://www.nies.go.jp/kenkyu/chikanken/</u>(2024 年7月閲覧).
- 古澤尚英, 曽我稔, 豊永悟史, 荒木真, 菅田誠司: 大気環境学会九州支部第 24 回研究発表会講演要旨 集, p14, (2024).
- 古澤尚英,豊永悟史,荒木真,曽我稔,菅田誠司: 第65回大気環境学会年会講演要旨集,p306,(2024).
- 10) 辻本昌礼, 山本浩平, 亀田貴之 : 気象モデル推定 値を取り入れた Land Use Regression モデルによ

る国内大気汚染物質濃度分布推定, J. Jpn. Soc. Atmos. Environ., **57(1)**, 1-14 (2022).

- 小原大翼,豊永悟史,古澤尚英,荒木真,山本裕典, 矢野弘道,山崎文雅:地方自治体における PM2.5 常時監視ネットワークの効率化の検討(I) -Regression Kriging法による空間濃度分布予測-, J. Jpn. Soc. Atmos. Environ., 57(2), 53-65 (2022).
- 12) 杉山将,井手剛,神嶌敏弘,栗田多喜夫,前田 英作 監訳:統計的学習の基礎 データマイニング・推論・予測,p14,p337-348 (2014),(共立出版);
 {Trevor Hastire, Robert Tibshirani, Jerome Friedman: "The Elements of Statistical Learning -Data Mining, Inference, and Prediction Second Edition-", (2009), (Springer New York) }.
- 13) Araki, S., Shima, M., Yamamoto, K : Spatiotemporal land use random forest model for estimating metropolitan NO2 exposure in Japan, *Science of The Total Environment*, 634, 1269-1277 (2018).
- 14) 松田晃一 訳: 解釈可能な AI -機械学習モデルの解 釈手法を実践的に理解する-, p24-61 (2023), (株式 会社マイナビ出版); {Ajay Thampi: Interpretable AI -Building Explainable Machine Learning Systems-, (2022), (Manning Publications) }.
- 15) 辻谷将明,外山信夫: Rによる GAM 入門, 行動計量 学, 34(1), 111-131 (2007).
- 16) 大原利眞,若松伸司,鵜野伊津志,安藤保,泉川碩 雄:関東・関西地域における光化学オキシダントの

経年動向に関する解析, J. Jpn. Soc. Atmos. Environ., **30(2)**, 137-148 (1995).

- 17) 星純也,石井康一郎 : 関東地域における揮発性有 機化合物 (VOC) 排出量の変化と光化学オキシダン ト生成の関係について, J. Jpn. Soc. Atmos. Environ., 48(5), 215-222 (2013).
- 18) 熊本県 : 熊本県の大気汚染の状況.
 <u>https://kumamoto-taiki.jp/index.html</u> (2024年7 月閲覧).
- 村上大補: 実践 Data Science R ではじめる地理空 間データの統計解析入門, (2022), (講談社).
- 20) Simon Wood. : Mixed GAM Computation Vehicle with Automatic Smoothness Estimation (CRAN `mgcv` package manual.). <u>https://cran.r-</u> project.org/web/packages/mgcv/index.html (2024 年7月閲覧).
- Schwarz Gideon. : Estimating the dimension of a model, *The annals of statistics*, 461-464 (1978).
- 22) Simon N. Wood : Generalized additive models: an introduction with R, p40 (2017), (chapman and hall/CRC).
- 23) Lex Comber. : Modelling Space and Time with GAMS: spatially and temporally varying coefficient models. <u>https://bookdown.org/lexcomber/AGILE2024/</u> (2024年8月閲覧).
- 24) 秋元肇 : 大気反応化学, p167-168 (2014), (朝倉 書店).

リスト4前方ステップワイズ法による交互作用項を導入した GAM の構築

```
stepwise_gam <- function(target = character(), ex = character(), int_lev = NULL, FUN = BIC, da
ta) {
 FUN <- substitute(FUN)
  message_gam <- function(model) {</pre>
    message("======="")
    message("formula: ", deparse(formula(model)))
    message(deparse(FUN), ": ", eval(bquote(round(.(FUN)(model), digits = 2))))
    message("R.sq: ", round(summary(model)$r.sq, digits = 4))
    message("Deviance explained: ", round(summary(model)$dev.expl * 100, digits = 2), "%")
    }
 create_formula <- function(target, ex) {</pre>
    parse(text = paste(target, "\sim", paste(paste0("s(", ex, ")"), collapse = " + ")))[[1]] |> a
s.formula()
 }
  ## calc normal gam
 ms <- lapply(ex, function(x) {</pre>
    mgcv::gam(create_formula(target, x), data = data)
 })
  num <- eval(bquote(which.min(Map(. (FUN), ms))))</pre>
 m <- ms[[num]]</pre>
  message_gam(m)
  for (i in setdiff(seq(1, length(ex), by = 1), num)) {
    m2 <- mgcv::gam(create_formula(target, ex[c(num, i)]),
                    data = data)
    message_gam(m2)
    if (eval(bquote(.(FUN)(m2))) <= eval(bquote(.(FUN)(m)))) {
      num \langle -c(num, i) \rangle
      m <- m2
    }
    rm(m2)
 }
  ## calc interactions
  if (!is.null(int lev)) {
    pat <- combn(ex, int_lev)</pre>
    for (i in seq(1, ncol(pat), by = 1)) {
      formula <- parse(text = paste(paste(deparse(formula(m)), collapse = ""), paste0(" + s(",</pre>
paste(pat[, i], collapse = ", "), ")")))[[1]] > as.formula()
      m2 <- mgcv::gam(formula,
                      data = data)
      message_gam(m2)
      if (eval(bquote(.(FUN)(m2))) <= eval(bquote(.(FUN)(m)))) {
```

```
num <- c(num, i)
       m <- m2
     }
     rm(m2)
     rm(formula)
   }
 }
 message("Completed.")
 return(m)
}
stepwise_gam(target = "0x", ex = c("solar_radiation", "temperature", "wind_speed"), int_lev =
2, data = uto, FUN = BIC) |> summary()
_____
formula: Ox ^{\sim} s(temperature)
BIC: 1380.16
R. sq: 0.6928
Deviance explained: 70.44%
_____
_____
formula: 0x \sim s(temperature) + s(solar_radiation)
BIC: 1380.45
R. sq: 0.7103
Deviance explained: 72.56%
_____
_____
formula: 0x ^{\sim} s(temperature) + s(wind_speed)
BIC: 1372.93
R. sq: 0.7397
Deviance explained: 75.75%
_____
_____
formula: 0x \sim s(temperature) + s(wind_speed) + s(solar_radiation, temperature)
BIC: 1413.78
R. sq: 0.79
Deviance explained: 82.89%
_____
_____
formula: 0x \sim s(temperature) + s(wind_speed) + s(solar_radiation, wind_speed)
```

```
BIC: 1420.05
R. sq: 0.8161
Deviance explained: 85.82%
_____
_____
formula: 0x \sim s(\text{temperature}) + s(\text{wind}_{\text{speed}}) + s(\text{temperature}, \text{wind}_{\text{speed}})
BIC: 1381.82
R. sq: 0.7576
Deviance explained: 78.17%
_____
Completed.
_____
Family: gaussian
Link function: identity
Formula:
0x \sim s(temperature) + s(wind_speed)
Parametric coefficients:
         Estimate Std. Error t value Pr(>|t|)
(Intercept) 57.1617 0.9672 59.1 <2e-16 ***
____
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '' 1
Approximate significance of smooth terms:
              edf Ref.df F p-value
s(temperature) 7.468 8.429 38.812 < 2e-16 ***
s(wind_speed) 3.870 4.797 5.983 8.29e-05 ***
____
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '' 1
R-sq. (adj) = 0.74 Deviance explained = 75.8%
GCV = 168.67 Scale est. = 156.21 n = 167
```