

## IV-184 データ領域の特性の予測誤差への影響について

北海道大学工学部 正員 ○山形 耕一

### 1. はじめに

交通計画において予測は、計画手段の配列の結果としての予測年次の交通諸量を出力し、意志決定のための情報提供する機能をもつが、現象の観測、モデル設定、予測のプロセスを通じて諸々の原因のため予測誤差がもたらされる。本研究では、これらの予測誤差を発生させる要因について分析し、このうち従来議論されることの少なかれデータ領域の拡張り、データ数および予測領域との関係が予測精度に及ぼす影響を考察する。ここで、データを多次元空間における点の集合と理解し、点の分布する領域をデータ領域と呼ぶことにする。

### 2. 交通モデルにおける予測誤差

交通予測における誤差の混入は、データ収集、モデル設定、予測の3段階に分けて考えることができよう。データ収集段階としては、交通現象に影響する要因が多岐に亘り、その全てを調査し切れないとや価値観、効用の如く計測の困難・不能な要因を含むことはモデル設定を制約する。また実質上および標本抽出による調査誤差は、モデルのパラメータ誤差を通じて予測誤差をもたらす。これらに加え、現実に起きた事象しかデータとして採取できない制約がある。個々の要因のデータ領域における値の分布域（値域と呼ぶ）は調査時点以前に現実に起きた事象における値の範囲に限られる。また、要因同志は互に相関をもつため、データは空間全体に散在するのではなく、相関関係に従って局在する傾向がある。この限られたデータ領域を対象に現象把握を行わなければならないため、データ領域外への予測性の制約、モデルの不安定性、さらには値域の小さい要因の影響の分析の困難さ等の問題が生じる。

モデル設定段階では上記のデータ収集上の制約に基づく問題点に加えて、複雑な現象を操作可能なモデルで記述する困難さを伴う。この困難さの一因としては、交通現象が混合的な性質を持ち、これらとに考えられる。交通は個人の、さらに個々の時点での交通行動の集合であり、個人による価値観の差異や個々の時点の行動環境の違いによる価値意識構造の差異が内在している。このため、交通行動に対する要因の作用性が異なる現象をマクロ的に記述することに要求されているのである。従って、モデルは平均的な現象を記述することとなり、モデルで説明し切れない残差を生じる。すなわち、モデルにはデータの誤差や残差等の搅乱項が含まれるが、搅乱項を含むデータから関係式を安定的に抽出するために必要なデータ数は明らかにされているとは言い難い。従来のサンプル理論は集團における要素の数標識の推定に関する理論であり、関係式の安定的な把握のためのサンプル理論が今後開発されることが必要であろう。

予測段階で誤差をもたらす原因としては、先ず予測のフレーム設定の問題である。予測モデルでは、予測年次における説明変数の将来値が必要であるが、これを定めるためのフレーム設定や説明変数値の将来値のもつ予測誤差は被説明変数の予測に大きな誤差をもたらす。第二には、データ領域と予測時点や予測対象値の存在する領域（予測領域と呼ぶ）の食い違いによる予測誤差がある。モデルはデータ領域外においては検証されない場合が多く、予測領域においてもモデルの構造やパラメータの値が保持されているかは保証されない。また、モデルの有効性が予測領域においても仮定できる場合でも、予測領域がデータ領域から遠く離れる程大きくなると、説明変数の相関関係の変化もまた予測誤差をもたらす。第三には、交通システムを取り巻く社会経済システム全体の構造変化があり、モデルに顕在的に取り込まれていない価値観等の要因の変化や予測時点においてその作用性が顕在化した要因の影響等が挙げられる。モデルに取り込まれなかった要因を含めて予測の空間を考えると、モデルに取り込まれていない要因の値域や相関関係の変化は予測領域の空間における位置を変化させている。さらに、予測モデルの説明変数は、説明変数と相関をもつところのモデルに取り込まれていない要因の影響をも代表

して表わしている。それゆえ、このような要因と説明変数の相関性の変化は、説明変数の作用性を変える予測誤差をもたらすと考えられる。したがって、予測領域の変化や相関性の変化に対するモデルの有効性の検証やデータ領域の拡大を意識モデル等を用いて行なうことは極めて重要と考えられる。

### 3. 重回帰モデルにおけるデータ領域の予測誤差への影響

重回帰モデルは、発生交通量予測等の交通計画のみならず多くの分野で用いられていくモデルであり、これを対象にデータの領域およびデータ数に関する検討を行う。対象とする現象は構造式

$$Y(X) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \varepsilon \quad (1) \quad \varepsilon \text{ は擾乱項で独立に } N(0, \sigma^2) \text{ に従うものとする。}$$

で記述されるものとする。説明変数  $X_1, \dots, X_p$ 、被説明変数  $Y$  に関する  $n$  組のデータ  $(y_{ij}, x_{1j}, \dots, x_{pj})$   $j=1, \dots, p$  から定められた (1) の回帰モデルを

$$\hat{Y}(X) = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots + \hat{\beta}_p X_p \quad (2) \quad S_{kj}^{ij}: \text{ 偏差平方和積和行列 } S_{kj} \text{ の逆行列の } k \text{ 行 } j \text{ 列の要素}$$

$$\hat{\beta}_k = \frac{1}{n} \sum_{j=1}^n S_{kj}^{ij} S_{jj}^{-1} \quad (3) \quad S_{jj}^{-1} = \sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ij} - \bar{x}_j)^T, \quad S_{kj} = \sum_{i=1}^n (x_{ij} - \bar{x}_j)(y_{ij} - \bar{y})$$

とする。推定パラメータ  $\hat{\beta}_k$  は  $\beta_k$  の不偏推定量であり、その分散  $V(\hat{\beta}_k)$  は、

$$V(\hat{\beta}_k) = S_{kk}^{-1} \sigma^2 = \sigma^2 / S_{kk} \quad (4)$$

となる。ここで  $R_k$  は  $X_k$  をそれ以外の説明変数  $X_1, \dots, X_{k-1}, X_{k+1} \dots X_p$  で説明した場合の重相関係数であり、説明変数間に相関がある場合のモデルの不安定性が  $(1-R_k^2)$  に表されている。パラメータ  $\hat{\beta}_k$  の推定誤差は  $S_{kk}$  が大きいときに小さくなる。 $S_{kk}$  は明らかに  $X_k$  の各々の値が平均  $\bar{x}_k$  から離れているとき大きくなるので、 $\hat{\beta}_k$  の安定性を増すためには、 $X_k$  の分布する領域を大きくとることおよび  $\bar{x}_k$  から離れたデータを収集することが有効となる。従って、 $Y$  に対する  $X_k$  の影響の線形性が経験的に知られ残差の分析を要しないときには、 $X_k$  のデータを上述の如く、計画的に配置して収集することが有効な方法となる。また、 $S_{kk}$  はデータ数が増すとき大きくなる。データに無作為に抽出された  $X_k$  が確率変数とみなせる場合には、 $X_k$  の分布が正規分布  $N(\mu_k, \sigma_k^2)$  で近似できるとすれば、 $S_{kk}$  の期待値が  $n-1$  となることから、 $\hat{\beta}_k$  の分散は  $1/(n-1)$  に比例して小さくなることが期待される。

説明変数の将来値  $\tilde{X}_0 = (x_{01}, \dots, x_{0p})$  に対する構造式上の値  $Y(\tilde{X}_0)$  と予測値  $\hat{Y}(\tilde{X}_0)$  の差の分散は、

$$V[Y(\tilde{X}_0) - \hat{Y}(\tilde{X}_0)] = \left\{ 1 + \frac{1}{n} + D_0^2 / (n-1) \right\} \sigma^2, \quad D_0^2 = \sum_{k=1, k \neq k_0}^p (x_{0k} - \bar{x}_k)(x_{0k} - \bar{x}_k)^T S_{kk}^{-1} / (n-1) \quad (5)$$

となる。ここに、 $D_0^2$  は  $P$  次元空間におけるデータの重心と  $\tilde{X}_0$  の間のマハラノビスの距離であり、 $D_0^2$  とデータ数が予測精度に関与している。すなわち、データの分布する領域と予測対象となる  $\tilde{X}_0$  の分布する領域との食い違いによる予測誤差は兩空間の間のマハラノビスの距離を用いて表現される。 $D_0^2$  はデータの分布する領域の外に出ると急激に大きくなるし、また、説明変数間の相関関係を反映しており、 $\tilde{X}_0$  がこの相関関係を満たしていない場合には著しく大きくなるためには、予測べ外挿となるときには  $D_0^2$  の増大に対応してデータ数を増加する必要がある。 $\hat{Y}(\tilde{X}_0) \pm d$  が確率  $1-\alpha$  をもつ予測対象値  $Y(\tilde{X}_0)$  を覆うようにして予測精度の一定化を図るために、

$$1 + \frac{1}{n} + D_0^2 / (n-1) \leq (d / S_{yy})^2 \quad (6) \quad d: \text{ 自由度 } n-p-1 \text{ の } \chi^2 \text{ 分布の } \frac{\alpha}{2} \text{ 点}, \quad \sigma^2 = S_{yy} (1-R^2) / (n-p-1)$$

なる関係を用いることが考えられる。必要データ数を事前に定めるためには、モデルの安定性および予測精度の場合共、未知の分散  $\sigma^2$  の大きさを仮定する必要があるが、 $S_{yy} (1-R^2) / \sigma^2$  が自由度  $n-p-1$  の  $\chi^2$  分布するときを用いて、分析で期待する重相関係数の大きさをもとに想定する方法が考えられる。

### 4. おわりに

本研究では、データ領域すなわち限られた次元およびデータの値域の観測に基づき予測を行うことの予測誤差への影響について考察した。そして、予測誤差がデータ領域の拡がり、データ数および予測領域とのマハラノビスの距離との関連について理解されることが分かった。関係式を抽出するためのサンプリングの考え方、意識モデルの併用による予測性の向上等が今後の課題として残されていく。