

打ち切りデータ補完のための複数ターゲットトビットモデル Multi-target Tobit model for completion of censored data

高田 裕也*
Yuya Takada

加藤 毅†
Tsuayoshi Kato

1. はじめに

飲用水や環境水には感染症を引き起こす病原体が含まれており、水利用における微生物学的安全性が脅かされている [2]. 病原体には、サルモネラ菌、赤痢菌、ノロウイルス、ロタウイルスなど多岐にわたっており、その多くはヒトの糞便による汚染に由来している. 水利用の用途と水中微生物の生態に応じて適切な水質衛生基準値を設定するためには、複数の水中微生物濃度およびほかの水質データから、それらの関係性を解析する. 水中微生物濃度や水質データの関係性は、採水ごとに同時に測定された観測値を使って解析する. 水中微生物濃度の定量化に際して最も大きな障害は陽性率の低さであり、少なからぬ割合で濃度が定量限界を下回っている (図 1(a),(c)). 相関解析などの統計解析を可能にするためには、非検出値の真の値を高精度に推定することが求められている.

打ち切りデータの解析には、トビットモデル [1] という統計学的手法がある. トビットモデルは回帰分析の一種で、目標変量のみに打ち切りを許している. 病原体濃度の解析には、各病原体に対して目標変量にしてトビットモデルをあてはめることで非検出データを補完することができる. しかし、このアプローチを使うには、説明変量に使われる病原体濃度は打ち切られていても何らかの値で埋めてから適用する必要がある. そのため、トビットモデルは複数の病原体濃度を同時に推定することができない.

本研究では、トビットモデルの目標変量を複数化するように拡張を施すことで、打ち切られた複数の目的変量を同時に補完する新しい方法論、**複数ターゲットトビットモデル (MTTM)** を考案した. 本研究の主要な理論的結果は 2 個の発見からなっている.

第 1 の発見: トビットモデルで用いられている尤度関数の一部の項がカルバックライブラ距離 (KL 距離) の最小値で表される (定理 1).

目標変量の複数化は、その表現を使った尤度関数を複数の目標変量について重ね合わせることで実現した. この拡張によって得られた目的関数最大化のための最適化アルゴリズムに関して、次のことを発見をした.

第 2 の発見: 目的関数の最大化にブロック座標上昇法を用いたとき、各ステップが閉形式で表される (定理 3).

この発見により、ステップサイズのような調整を要するパラメータが不要となり、数値的に安定した最適化が可能となる. 水中微生物濃度の実データを用いた数

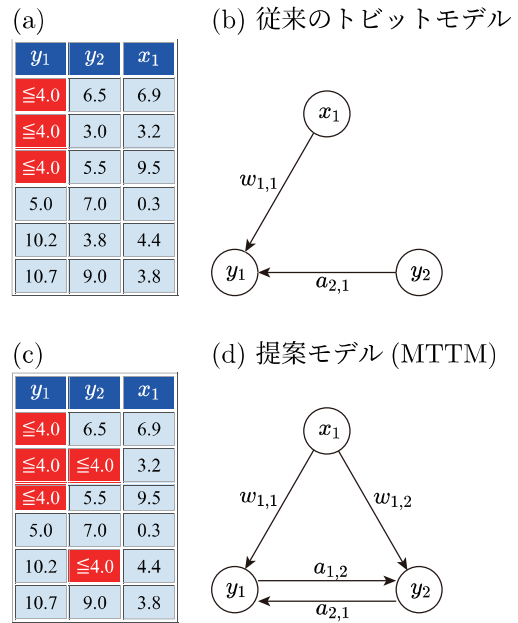


図 1: 従来トビットモデルと提案モデル. 従来のトビットモデルは、打ち切りがある列を含む表 ((a) 参照) から、その列を目標変数とした回帰モデル ((b) 参照) をデータにフィットさせる枠組みである. 提案モデルである MTTM は、複数の打ち切りを含む表 ((c) 参照) から、複数の回帰モデルを重ね合わせる ((d) 参照) ことで、同時にフィットさせることができる.

値実験によって検証したところ、提案する複数ターゲットトビットモデルは従来のトビットモデルを顕著に上回る精度で非検出値を補完できることを確認した.

2. 従来のトビットモデル (STTM)

トビットモデル [1] はデータセット $(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$ ($i \in [n] := \{1, \dots, n\}$) に対し、回帰モデル

$$y_i = \langle \mathbf{w}, \mathbf{x}_i \rangle + \epsilon_i$$

を仮定して、偏回帰係数 $\mathbf{w} \in \mathbb{R}^d$ の値を推定するモデルである. ϵ_i は雑音であり、 $\epsilon_i \sim \mathcal{N}(0, \beta^{-1})$ が仮定される. 定量限界を $u \in \mathbb{R}$ とおき、 n_v 個が定量限界を越えているとする. $i \in [n_v]$ に対し、 $y_i > u$, $i \in [n] \setminus [n_v]$ に対し、 $y_i \leq u$ となるよう、例題番号を入れ替える. トビットモデルが標準的な線形回帰分析と異なる点は、 y_i が定量限界 u を下回ったとき、 $y_i \leq u$ という情報は利用できるが、 y_i の値そのものは利用できないことである. トビットモデルのフィッティングには次の対数

*群馬大学理工学部
†群馬大学情報学部

尤度関数が用いられている：

$$\begin{aligned} \mathcal{L}_s(\boldsymbol{\theta}_s) := & \sum_{i=1}^{n_v} \log \mathcal{N}(y_i; \langle \mathbf{w}, \mathbf{x}_i \rangle, \beta^{-1}) \\ & + \sum_{i=n_v+1}^n \log \int_{-\infty}^u \mathcal{N}(t_i; \langle \mathbf{w}, \mathbf{x}_i \rangle, \beta^{-1}) dt_i, \end{aligned} \quad (1)$$

ただし、 $\boldsymbol{\theta}_s := (\mathbf{w}, \beta)$ はモデルパラメータを表す。 $\boldsymbol{\theta}_s$ の最適解はニュートン法などで見つけることができる [1].

筆者らは (1) の第 2 項の符号を逆転すると、KL 距離の最小値に等しいことを見つけた。次節では、この事実を利用するとトビットモデルの自然な拡張が得られることを示す。 \mathcal{P} を $\mathbb{R} \rightarrow \mathbb{R}$ なる確率密度関数の集合とし、次のような n 個の確率密度関数 $q_1, \dots, q_n \in \mathcal{P}$ を導入する。最初の n_v 個は Dirac デルタ関数 $\delta(\cdot)$ で定義する： $\forall i \in [n_v]$,

$$q_i(t) = \delta(t - y_i).$$

$i \in [n_v]$ に対して、 $Q_i := \{\delta(\cdot - y_i)\}$ とおく。残りの $n_h := n - n_v$ 個の確率密度関数 $q_{n_v+1}, \dots, q_n \in \mathcal{P}$ は以下を満たすとする： $\forall i \in [n] \setminus [n_v]$,

$$\forall t > u, \quad q_i(t) = 0. \quad (2)$$

$i \in [n] \setminus [n_v]$ に対して、 $Q_i := \{q_i \in \mathcal{P} \mid \forall t > u, q_i(t) = 0\}$ とおき、また、 $Q_s := Q_1 \times \dots \times Q_n$ とおく。(1) の第 2 項に関して、

$$\begin{aligned} -\log \int_{-\infty}^u \mathcal{N}(t_i; \langle \mathbf{w}, \mathbf{x}_i \rangle, \beta^{-1}) dt_i \\ = \min_{q_i \in Q_i} \text{KL}[q_i \parallel \mathcal{N}(\cdot; \langle \mathbf{w}, \mathbf{x}_i \rangle, \beta^{-1})] \end{aligned} \quad (3)$$

が成り立つことを示すことができる。ただし、 $\text{KL}(q \parallel p)$ は $p \in \mathcal{P}$ から $q \in \mathcal{P}$ への KL 距離を表す[†]。等式 (3) の導出は 5 節参照。この事実は直ちに次の定理を導く。

定理 1. $\beta > 0$ なる任意の $\boldsymbol{\theta}_s = (\mathbf{w}, \beta) \in \mathbb{R}^d \times \mathbb{R}$ に対して、

$$\mathcal{L}_s(\boldsymbol{\theta}_s) = \max_{q \in Q_s} \mathcal{F}_s(\boldsymbol{\theta}_s, q)$$

が成り立つ。ただし、 $\mathcal{F}_s(\boldsymbol{\theta}_s, q)$ は、 $q := (q_1, \dots, q_n) \in Q_s$ に対して、

$$\begin{aligned} \mathcal{F}_s(\boldsymbol{\theta}_s, q) := & \sum_{i=n_v+1}^n \mathcal{H}[q_i] + \sum_{i=1}^n \int q_i(t_i) \cdot \\ & \log \mathcal{N}(t_i; \langle \mathbf{w}, \mathbf{x}_i \rangle, \beta^{-1}) dt_i \end{aligned}$$

とおいた。 $\mathcal{H}[\cdot] : \mathcal{P} \rightarrow \mathbb{R}$ はエントロピーを表す[‡]。

[†]本稿では、KL 距離 $\text{KL}[q \parallel p]$ およびエントロピー $\mathcal{H}[q]$ の積分は $q(x) > 0$ なる区間 I のみとする。すなわち、 $\mathcal{H}[q] := -\int_{x \in I} q(x) \log q(x) dx$ および $\text{KL}[q \parallel p] := \int_{x \in I} q(x) \log(q(x)/p(x)) dx$ という定義を用いる。

3. 複数ターゲットトビットモデル (MTTM)

本節では、トビットモデルを複数のターゲット変数を扱えるようにするために、定理 1 から自然に導かれる拡張した MTTM を提案する。

目標変数の個数を m とおき、MTTM をフィットさせるデータセットを

$$(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n) \in \mathbb{R}^d \times \mathbb{R}^m$$

とする。各例題は説明変量 $\mathbf{x}_i \in \mathbb{R}^d$ と目標変量

$$\mathbf{y}_i := [y_{1,i}, \dots, y_{m,i}]^\top$$

からなる。ベクトル $\mathbf{y}_i \in \mathbb{R}^m$ のうち、第 k 要素を除いたベクトルを $\mathbf{y}_{\setminus k,i} := [y_{1,i}, \dots, y_{k-1,i}, y_{k+1,i}, \dots, y_{m,i}]^\top$ と書くことにする。MTTM の回帰モデルは m 個の回帰モデル

$$k \in [m], \quad y_{k,i} = \langle \mathbf{a}_k, \mathbf{y}_{\setminus k,i} \rangle + \langle \mathbf{w}_k, \mathbf{x}_i \rangle + \epsilon_{k,i}$$

で構成する。ただし、 $\mathbf{a}_k \in \mathbb{R}^{m-1}$ および $\mathbf{w}_k \in \mathbb{R}^d$ は第 k 回帰モデルの偏回帰係数である； $\epsilon_{k,i} \sim \mathcal{N}(0, \beta^{-1})$ は雑音である。各目標変量はほかの目標変量に対する説明変量になっていることに注意されたい。

MTTM では、目標変量の行列 $\mathbf{Y} := [\mathbf{y}_1, \dots, \mathbf{y}_n] \in \mathbb{R}^{m \times n}$ が不完全になっている状況を仮定する。目標変量 $y_{k,i} \in \mathbb{R}$ が定量限界 $u_k \in \mathbb{R}$ より大きければ、 $y_{k,i}$ の値を利用できる。しかし、 $y_{k,i} \leq u_k$ ならば、 $y_{k,i}$ の値を利用できない。欠損要素の添え字集合を \mathcal{E}_h 、観測要素の添え字集合を \mathcal{E}_v とおく；i.e.

$$\begin{aligned} \mathcal{E}_h &:= \{(k, i) \in [m] \times [n] \mid y_{k,i} \leq u_k\}, \\ \mathcal{E}_v &:= \{(k, i) \in [m] \times [n] \mid y_{k,i} > u_k\}. \end{aligned}$$

偏回帰係数の決定に用いる目的関数を定義するために、次のような mn 個の確率密度関数 $q_{k,i} \in \mathcal{P}$ を導入する：

$$\forall (k, i) \in \mathcal{E}_v, \quad q_{k,i}(t) = \delta(t - y_{k,i}) \quad (4)$$

$$\begin{aligned} \forall (k, i) \in \mathcal{E}_h, \quad q_{k,i}(t) &= 0 \quad \text{if } t > u_k, \\ q_{k,i}(t) &> 0 \quad \text{if } t \leq u_k. \end{aligned} \quad (5)$$

また、 $(k, i) \in \mathcal{E}_v$ に対して、 $Q_{k,i} := \{\delta(\cdot - y_{k,i})\}$ とおく。 $(k, i) \in \mathcal{E}_h$ に対して、 $Q_{k,i}$ を式 (5) を満たす確率密度関数の集合とおく。これらの直積集合を $Q_m := \prod_{(k,i) \in [m] \times [n]} Q_{k,i}$ とおく。 Q_m を使って、 \mathcal{F}_s を複数目標変量に対応できるように次のように拡張する：

$$\begin{aligned} \mathcal{F}_m(\boldsymbol{\theta}_m, q) := & \sum_{(k,i) \in \mathcal{E}_h} \mathcal{H}[q_{k,i}] + \sum_{i=1}^n \sum_{k=1}^m \int q_{1,i}(t_{1,i}) \cdot \\ & \dots \cdot q_{m,i}(t_{m,i}) \log \mathcal{N}(t_{k,i}; \langle \mathbf{a}_k, \mathbf{t}_i \rangle + \langle \mathbf{w}_k, \mathbf{x}_i \rangle, \beta^{-1}) dt_i \end{aligned}$$

ただし、 $\mathbf{t}_i := [t_{1,i}, \dots, t_{m,i}]^\top$ 、 $q := (q_{k,i})_{k,i} \in Q_m$ であり、 $\boldsymbol{\theta}_m := (\mathbf{a}_1, \dots, \mathbf{a}_m, \mathbf{w}_1, \dots, \mathbf{w}_m, \beta)$ と再定義した。MTTM のパラメータ $\boldsymbol{\theta}_m$ の値を決定するために、次の最適化問題を解く：

MTTM のフィッティング問題

$$\begin{aligned} \max \quad & \mathcal{L}_m(\boldsymbol{\theta}_m) \quad \text{wrt } \boldsymbol{\theta}_m, \\ \text{where } \quad & \mathcal{L}_m(\boldsymbol{\theta}_m) := \max_{q \in \mathcal{Q}_m} \mathcal{F}_m(\boldsymbol{\theta}_m, q). \end{aligned} \quad (6)$$

すなわち、 $\boldsymbol{\theta}_m$ と $q \in \mathcal{Q}_m$ の同時最適化によって偏回帰係数の値を決定する。

欠損値の補完. 最適解 $(\boldsymbol{\theta}_m^*, q^*)$ を得ることが出来れば、定量限界以下の値を期待値 $\mathbb{E}_{q_{k,i}(t_{k,i})}[t_{k,i}]$ を使って推定できる。次の定理によって、期待値計算が容易であることが分かる。(証明は 5 節参照。)

定理 2. 最大化問題 (6) の最適解 $(\boldsymbol{\theta}_m^*, q^*)$ において、 $q^* \in \mathcal{Q}_m$ の各要素 $q_{k,i}^* \in \mathcal{Q}_{k,i}$ は切断正規分布で表すことができる。すなわち、

$$\exists \mu_{k,i}^* \in \mathbb{R}, \exists \sigma_{k,i}^* \in \mathbb{R}, \quad q_{k,i}^* = f_{TN}(\cdot; \mu_{k,i}^*, \sigma_{k,i}^*, u_k).$$

ただし、 $f_{TN}(\cdot; \mu, \sigma, u)$ は切断正規分布の確率密度関数を表す：

$$f_{TN}(x; \mu, \sigma, u) := \begin{cases} \frac{\phi((x-\mu)/\sigma)}{\sigma\Phi((x-\mu)/\sigma)} & \text{for } x \leq u, \\ 0 & \text{for } x > u. \end{cases}$$

ただし、 ϕ および Φ はそれぞれ標準正規分布の確率密度関数および累積分布関数を表す。

定理 2 より、 $y_{k,i}$ の推定値は、 $\hat{u}_{k,i}^* := (u_k - \mu_{k,i}^*)/\sigma_{k,i}^*$ を使って

$$\mathbb{E}_{q_{k,i}(t_{k,i})}[t_{k,i}] = \mu_{k,i}^* - \sigma_{k,i}^* \frac{\phi(\hat{u}_{k,i}^*)}{\Phi(\hat{u}_{k,i}^*)}$$

と求められる。

4. MTTM のフィッティング

前節で、MTTM のデータへのフィッティングを $\mathcal{F}_m(\boldsymbol{\theta}_m, q)$ の最大化によって行うことを述べた。筆者らは、 $\mathcal{F}_m(\boldsymbol{\theta}_m, q)$ の最大化のためにブロック座標上昇法を採用した。Algorithm 1 に示すように、ブロック座標上昇法では、ある $(k, i) \in \mathcal{E}_h$ に対し、 $q_{k,i}$ 以外の q の構成要素および $\boldsymbol{\theta}_m$ を固定して、 $q_{k,i}$ に関して $\mathcal{F}_m(\boldsymbol{\theta}_m, q)$ を最大化するステップ (Line 5) と、 $q \in \mathcal{Q}_x$ を固定して $\boldsymbol{\theta}_m$ に関して $\mathcal{F}_m(\boldsymbol{\theta}_m, q)$ を最大化するステップ (Line 7) からなる。このアルゴリズムに関して、筆者らは次の定理を見つけた。(証明は 5 節参照。)

定理 3. Algorithm 1 の各ステップ (Line 5 および Line 7) の解は閉形式で表すことができる。

Algorithm 1: Block coordinate ascent algorithm for maximizing $\mathcal{F}_m(\boldsymbol{\theta}_m, q)$.

```

1 begin
2   Initialize  $\boldsymbol{\theta}_m$  and  $q \in \mathcal{Q}_x$ ;
3   for  $t := 1, 2, \dots$  do
4     for  $(k, i) \in \mathcal{E}_h$  do
5        $q_{k,i} := \operatorname{argmax}_{q_{k,i} \in \mathcal{Q}_{k,i}} \mathcal{F}_m(\boldsymbol{\theta}_m, q)$ ;
6     end
7      $\boldsymbol{\theta}_m := \operatorname{argmax}_{\boldsymbol{\theta}_m} \mathcal{F}_m(\boldsymbol{\theta}_m, q)$ ;
8   end
9 end

```

各ステップの具体的な更新則を以下に述べる。

4.1. 密度関数 q の更新

Algorithm 1 の Line 5 では、 $q_{k,i}$ を $\mathcal{F}_m(\boldsymbol{\theta}_m, q)$ を最大化する密度関数に更新する。そのような $q_{k,i}$ は次のような切断正規分布で表される：

$$q_{k,i}^{\text{new}}(t_{k,i}) = f_{TN}(t_{k,i}; \mu_{k,i}, \sigma_{k,i}, u_k). \quad (7)$$

$k = 1$ のとき、 $\mu_{k,i}$ および $\sigma_{k,i}$ はベクトル $\bar{\mathbf{y}}_r = [\bar{y}_{2,i}, \dots, \bar{y}_{m,i}]^\top$ および行列 $\mathbf{A}^{m \times m}$ を使って算出する。ただし、 $\bar{y}_{k',i} := \mathbb{E}_{q_{k',i}(t_{k',i})}[t_{k',i}]$ であり、行列 \mathbf{A} の第 k' 列には、 $k' \in [m]$ に対し、 $\mathbf{a}_{k'} \in \mathbb{R}^{m-1}$ の第 k' 要素に 0 を挿入した m 次元列ベクトルを入れる。これらを使って、次のように $\mu_{k,i}, \sigma_{k,i}$ を計算する：

$$\mu_{k,i} := \frac{\langle \mathbf{b}_1, \mathbf{W}^\top \mathbf{x} - \mathbf{B}_r^\top \bar{\mathbf{y}}_r \rangle}{\|\mathbf{b}_1\|^2}, \quad \sigma_{k,i} := \frac{1}{\sqrt{\beta} \|\mathbf{b}_1\|}.$$

ただし、 $\mathbf{b}_1 \in \mathbb{R}^m$ および $\mathbf{B}_r \in \mathbb{R}^{(m-1) \times m}$ は $[\mathbf{b}_1, \mathbf{B}_r^\top] = \mathbf{I} - \mathbf{A}^\top$ を満たすように算出する。 $k \neq 1$ のときは、添え字 k と 1 を入れ替えてこの手続きを実行する。

4.2. モデルパラメータ $\boldsymbol{\theta}_m$ の更新

Algorithm 1 の Line 7 における $\boldsymbol{\theta}_m$ の更新式も閉じた形で表される。分布 $q_{k,i}$ に関する期待値

$$\bar{y}_{k,i} := \mathbb{E}_{q_{k,i}(t_{k,i})}[t_{k,i}], \quad v_{k,i} := \mathbb{E}_{q_{k,i}(t_{k,i})}[t_{k,i}^2]$$

から構成されるベクトル

$$\bar{\mathbf{y}}_i := [\bar{y}_{1,i}, \dots, \bar{y}_{m,i}]^\top, \quad \mathbf{v}_i := [v_{1,i}, \dots, v_{m,i}]^\top$$

を考え、それぞれから第 k 要素を除いたベクトルを $\bar{\mathbf{y}}_{\setminus k,i}, \mathbf{v}_{\setminus k,i}$ とおく。また、 $\forall k \in [m], \forall i \in [n]$ について、

$$\tilde{\mathbf{w}}_k := \begin{bmatrix} \mathbf{a}_k \\ \mathbf{w}_k \end{bmatrix}, \quad \tilde{\mathbf{x}}_{k,i} := \begin{bmatrix} \bar{\mathbf{y}}_{\setminus k,i} \\ \mathbf{x}_i \end{bmatrix}, \quad (8)$$

$$\mathbf{G}_k := \operatorname{diag} \left(\begin{bmatrix} \sum_{i=1}^n \mathbf{v}_{\setminus k,i} \\ \mathbf{0}_d \end{bmatrix} \right)$$

とおく。すると、密度関数 q を固定して、目的関数 $\mathcal{F}_m(\boldsymbol{\theta}_m, q)$ を最大化する $(\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_m, \beta)$ は以下のよ

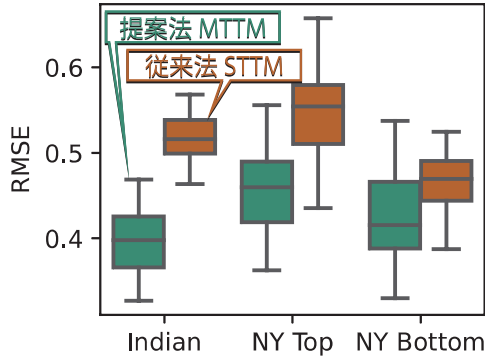


図 2: 性能評価.

うに表される:

$$\begin{aligned} \tilde{\mathbf{w}}_k^{\text{new}} &= \left(\mathbf{G}_k + \sum_{i=1}^n \tilde{\mathbf{x}}_{k,i} \tilde{\mathbf{x}}_{k,i}^\top \right)^{-1} \sum_{i=1}^n \tilde{y}_{k,i} \tilde{\mathbf{x}}_{k,i}, \\ \frac{1}{\beta^{\text{new}}} &= \frac{1}{mn} \sum_{k=1}^m \langle \tilde{\mathbf{w}}_k^{\text{new}}, \mathbf{G}_k \tilde{\mathbf{w}}_k^{\text{new}} \rangle \\ &\quad + \frac{1}{mn} \sum_{k,i} (\tilde{y}_{k,i} - \langle \tilde{\mathbf{w}}_k^{\text{new}}, \tilde{\mathbf{x}}_{k,i} \rangle)^2 + \frac{1}{mn} \sum_{k,i} v_{k,i}. \end{aligned}$$

5. 導出と証明

定理 1 を導く等式 (3) の導出, 定理 2, および定理 3 の証明において, 次の補助定理がカギとなる:

補助定理 4. 関数 $f: \mathbb{R} \rightarrow \mathbb{R}$ に対して, 次の最適化問題を考える:

$$\begin{aligned} \max \quad & F[q] \quad \text{wrt } q \in \mathcal{P} \\ \text{subject to} \quad & \forall x > u, \quad q(x) = 0, \\ \text{where} \quad & F[q] := \mathbb{E}_{q(x)} [\log f(x)] + \mathcal{H}[q], \\ & \forall x \leq u, \quad f(x) > 0. \end{aligned} \quad (9)$$

この問題の最適解は

$$q(x) := \begin{cases} \frac{f(x)}{\int_{-\infty}^u f(t) dt} & \text{for } x \leq u, \\ 0 & \text{for } x > u \end{cases} \quad (10)$$

で与えられる.

補助定理 4 において, 一部の区間で 0 という縛りが確率密度関数 $q \in \mathcal{P}$ になければ, 平均場近似で使われるトリックに帰着する. 補助定理 4 の証明は汎関数 $F[q]$ に対してラグランジュ乗数法を適用することで得ることができる.

式 (3) の導出 (スケッチ): 補助定理 4 における $F[q]$ は $f := \mathcal{N}(\cdot; \langle \mathbf{w}, \mathbf{x}_i \rangle, \beta^{-1})$ のとき, $\mathcal{F}[q] = -\text{KL}[q||f]$

が成り立つ. $Z := \int_{-\infty}^u f(t) dt$ とし, $\hat{q}(x) := \mathbf{1}[x \leq u] \cdot f(x)/Z$ とすると, 補助定理 4 より,

$$\begin{aligned} \text{RHS of (3)} &= -\min_{q \in \mathcal{Q}} \text{KL}[q||f] = \mathcal{F}[\hat{q}] \\ &= -\text{KL}[\hat{q}||f] = \int_{-\infty}^u \frac{f(t)}{Z} \log \frac{p(t)}{f(t)/Z} dt \\ &= \frac{\log Z}{Z} \int_{-\infty}^u f(t) dt = \log Z = \text{LHS of (3)}. \end{aligned}$$

定理 3 の証明 (スケッチ): ここでは, 本文中にすでに与えた閉形式解に対するそれぞれの導出のスケッチを示す. $f_{k,i}(t_{k,i}) := \mathcal{N}(t_{k,i}; \mu_{k,i}, \sigma_{k,i}^2)$ とおくと,

$$\mathcal{F}_m(\boldsymbol{\theta}, q) = \mathcal{H}[q_{k,i}] + \mathbb{E}_{q_{k,i}(t_{k,i})} [\log f_{k,i}(t_{k,i})] + \text{const} \quad (11)$$

のように整理できる. ただし, const は $q_{k,i}$ に依存しない項を表す. 補助定理 4 より, Algorithm 1 の Line 5 の解が導かれる.

Line 7 の解は $\frac{\partial \mathcal{F}_m(\boldsymbol{\theta}, q)}{\partial \boldsymbol{\theta}} = \mathbf{0}$ とおくことにより得られる. \square

定理 2 の証明 (スケッチ): 補助定理 4 を式 (11) に適用することにより, $\forall (k, i) \in \mathcal{E}_h$ に対して $q_{k,i}^*$ は切断正規分布の密度関数となることが導かれる. \square

6. 数値実験

MTTM の性能を評価するため, 3 個の水質データセット Indian, NY Top, NY Bottom を用いた. 総大腸菌群数 (TC) と糞便性大腸菌群数 (FC) を検出限界のある 2 種類の病原体濃度と見なし, 陰性率 10% として疑似的に打ち切りデータを用意した. データセット全体から無作為に選んだ $n = 100$ 個のサブセットを使った. これを 50 回繰り返して, それぞれで推定値の RMSE を算出した. 各データセットにおける RMSE の分布を図 2 にプロットした. 図中の STTM は, 説明変数に検出限界値を代入して STTM を適用したときの結果である. 3 個のデータセットいずれにおいても MTTM の RMSE が STTM の RMSE より顕著に低かった. よって, MTTM は打ち切られた複数の変数を補完するための有用な枠組みといえる.

謝辞 本研究の一部は, 環境研究総合推進費により実施され, JSPS 科学研究費 22K04372 の助成を受けたものである.

参考文献

- [1] Takeshi Amemiya. Tobit models: A survey. *Journal of Econometrics*, 24(1-2):3-61, January 1984.
- [2] T. Kato, A. Kobayashi, W. Oishi, S. S. Kadoya, S. Okabe, N. Ohta, M. Amarasiri, and D. Sano. Sign-constrained linear regression for prediction of microbe concentration based on water quality datasets. *J Water Health*, 17(3):404-415, Jun 2019.