

管理工学における統計解析

鈴木 秀男, 松浦 峻, 山田 秀

筆者らは、慶應義塾大学理工学部管理工学科において統計分野の教育研究を行っている。管理工学科における統計なので、理論、応用の両面をバランスよく取り入れ、また、幅広い分野を対象とするように心がけている。このような多様性の紹介を目指すべく、(1) 野球、サッカーなどのスポーツデータの解析という応用に関する研究、(2) 多変量推測統計の理論研究の例として、多次元確率分布の principal points, seemingly unrelated regression モデルにおける最良共変推定量についての研究、(3) よりよいデータの収集方法という面で、実験計画の構成に関する研究を紹介する。

キーワード：スポーツデータ解析、多変量推測統計、実験計画法

1. はじめに

慶應義塾大学理工学部管理工学科では、鈴木秀男、松浦峻、山田秀の3名が統計分野での教育、研究に携わっている。管理工学科における統計なので、理論、応用の両面をバランスよく取り入れ、また、幅広い分野を対象とするように心がけている。このような多様性の紹介を目指すべく、まず2節では、鈴木秀男が野球、サッカーなどのスポーツデータの解析という応用に関する研究を紹介する。また3節では、松浦峻が多変量推測統計の理論研究の例として、多次元確率分布の principal points, seemingly unrelated regression モデルにおける最良共変推定量について紹介する。さらに4節では、山田秀がよりよいデータの収集方法という面で、実験計画の構成に関する研究を紹介する。

2. スポーツデータ競技解析の研究事例

2.1 基本的な考え方

応用統計解析は「事実に基づく管理」のための学問であるという意識のもとで、教育、研究活動を進めている。組織活動を効果的・効率的に進めていくために、過去の経験や勘だけに頼らず、「事実に基づく管理」が重要とされている。常に「事実に基づくデータがあるかどうか」が問われて、「データに基づく客観的な判断と行動」を行うべきとされる。この考え方は、モノづくりの現場での品質管理活動において、よくいわれており、非常に重視されている。問題解決にあたっては、現場で現物を見て、現実を知るという「三現主義」に従い、経験や勘だけに頼るのではなく、データや観察結果に基づいて

PDCA (Plan, Do, Check, Act) を回すことが重要である [1]。また、ビッグデータ、データサイエンスというキーワードをよく耳にするようになり、統計学の関連本がベストセラーになり、深層学習 (Deep Learning) も世間で話題になるなど、データの活用に注目が集まっている。情報通信技術 (ICT) の進化により、ビジネス、行政、医療、スポーツなどさまざまな分野で、大量で多様なデータが取得できるようになり、それらのデータや高度な解析技術を活用して価値のある情報を抽出し、アクションにつなげようということである。管理工学では、基礎理論を重視し、さまざまな研究分野を有機的に融合しながら、経営・社会の課題の発見・解決を目指している。データサイエンスは、管理工学とかなり親和性が高いと考えている。そのような状況から、産業界からもさまざまな課題・問題に関する依頼・相談が増えている。

2.2 野球の球種予測モデルの構築

野球におけるデータ解析の一例として、投手の配球予測システムというものがある。これは1球ごとのデータに対して機械学習を活用して投手の次の配球を予測するというものもある。本研究では、一般社団法人日本統計学会および日本統計学会 スポーツ統計分科会主催の第8回スポーツデータ解析コンペティションで提供された、プロ野球の2016年と2017年の全投球データに基づき、深層学習の一種である LSTM (Long Short Term Memory) (図1) を用いて、投手の次の球種を予測するモデルを構築した [2]。特徴として、投手によらずに一つの共通モデルを作成することで、汎用性の高いモデルになっている。また、データ入力に対してダイナミックな構造を取り入れたことで、予測精度を向上させた。具体的な検証結果について、九つある球種の中で次に投げられる球種を予測するときの精度は50.2%となり、ロジスティック回帰などのほかの手法に比べ予測

すずき ひでお, まつうら しゅん, やまだ しゅう
慶應義塾大学理工学部管理工学科
〒223-8522 神奈川県横浜市港北区日吉 3-14-1

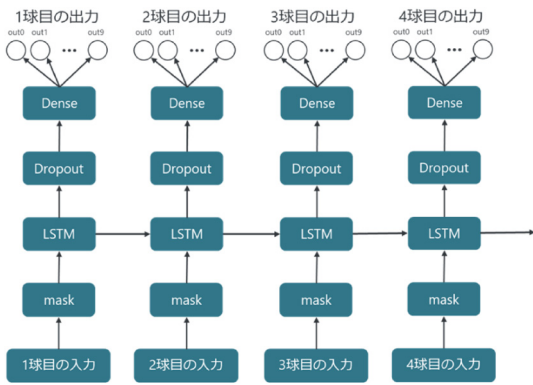


図1 球種予測のための LSTM のモデル図

の精度が高かった。本研究のモデルは、投球予測モデルとして一つのコンテンツになるだけでなく、実際の野球の現場での予測、野球ゲームなどの球種選択のアルゴリズムの発展にも応用できると考えられる。本成果は、当該コンペティションにおいて最優秀賞を受賞している。

2.3 サッカーの守備構造の評価

サッカーにおいて、被シュートや被ゴールといった指標のみでは守備の評価は難しく、たとえば、いくら相手に守備を崩されても相手が最後のシュートを成功させなければ失点とはならない。また、相手にチャンスをつくらなくても、一度のフリーキックやスーパーゴールから失点してしまうこともある。本研究では、そのような課題を解決するために、第8回スポーツデータ解析コンペティションで提供された J1 リーグ 2017 年 最終5 節 (= 45 試合) のデータを利用して、守備側選手の位置データからスペースがどの位置にどのくらいの大きさ(危険度)で存在するかを、ポロノイ図と空円(図2)という考え方をういて定量化した [3]。具体的には、次の二つの分析を行った。最初に、守備側選手を母点とするポロノイ図を作図し、ポロノイ点を中心とする空円を求めることで幾何学的にスペースの位置および半径を計算した。その後、ゴール、ボール、攻撃側選手の座標を用いて各空円の半径を再計算し、指標 AR(Adjusted Radius) を提案した。この指標は守備側選手の間スペースがどれほど危険であるかを定量化するものである。守備側選手の位置のみを使うと過大評価してしまうスペース(相手陣奥など)を適切に評価することができた。また、AR 平均値の時間変化を見ることで、守備構造が崩れていく様子を観察することができた。この手法はサッカーのみならず、バスケットボール、ラグビー、アメリカンフットボールなどの他スポーツに簡単に応用することができる。本成果は「第8回スポーツデータ解析コ

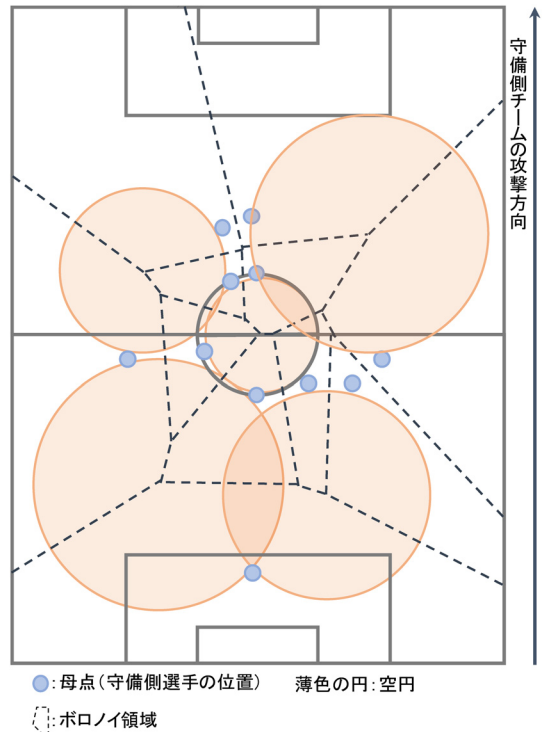


図2 守備側選手を母点とするポロノイ図と空円の例

ンペティション優秀賞」を受賞している。2.2 節および 2.3 節の研究事例では、LSTM やポロノイ図といった応用統計解析ではあまり用いられないアプローチを採用している。これも管理工学ならではの研究と考えている。

2.4 プロ野球チーム満足度調査

次に、プロ野球のサービスを対象にして、サービス品質と顧客満足度の因果モデルを構築し、因果関係の検証や数値化を行っている研究を紹介する。サービスには無形性という特徴があり、一般的に、サービス品質に関する評価の数値化が困難とされる。本研究では、顧客(ファン)の知覚品質(ファンがどのように品質を認識したか)でサービス品質を評価してもらい、統計解析を通して、プロ野球チームのサービス水準や満足度を数値化する。具体的には、プロ野球の顧客満足度指数化モデル [4-6] に基づき、2009 年から 2020 年の 1 月下旬において継続的に実施した 12 年間のプロ野球サービスの調査データから、経年変化に着目した分析を行う。毎回の調査において、各チームのファンに対してアンケート調査を実施し、チーム成績、チーム・選手の魅力、ファンサービス・地域貢献、ユニホーム・ロゴ、総合満足度、応援ロイヤルティ(チーム応援の意向など)、観戦ロイヤルティ(球場での観戦意向など)に関する項目、その他の関連項目を評価してもらう。回

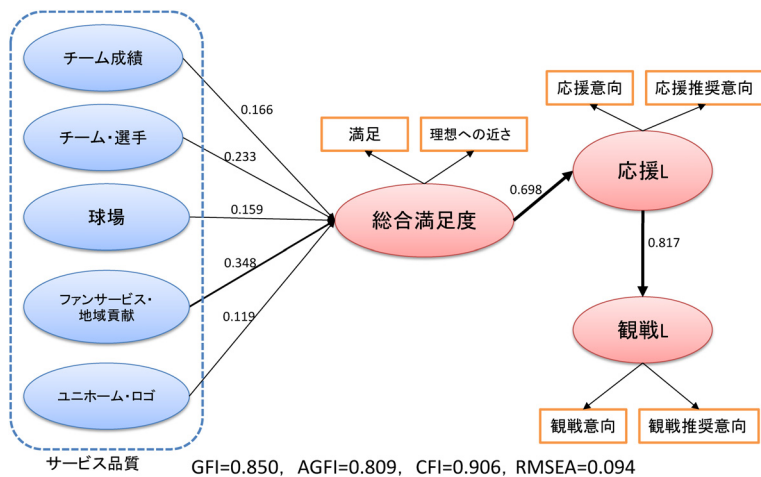


図3 プロ野球チームの顧客満足度指数化モデルと推定結果（2020年1月下旬調査）
注：矢線近くの数値は標準化係数を表す。

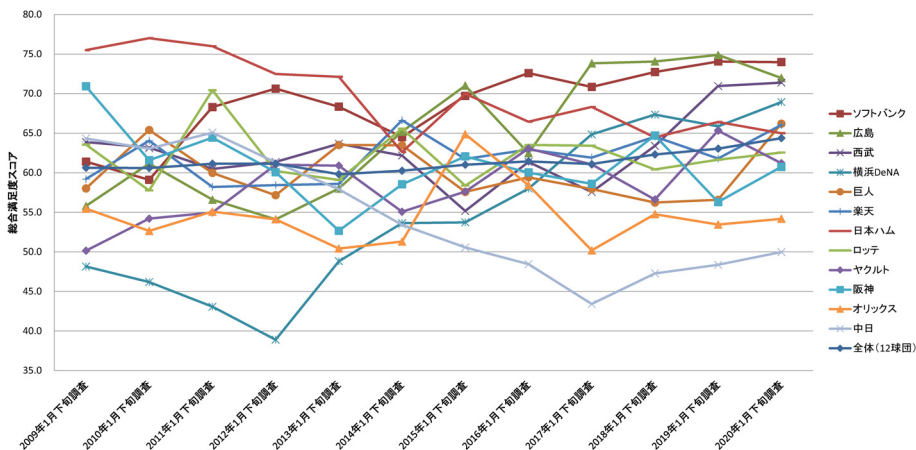


図4 各チームの総合満足度スコアの経年変化

答者の条件は、調査時点から1年以内に1回以上、応援するチームのホーム球場で試合観戦をしている人である。毎年の各チームの回答者数は100から120を確保、12年間の合計数は17,759である。なお、項目や調査などの詳細は、プロ野球のサービスの満足度調査HP [7]を参照されたい。図3は、共分散構造分析という手法を用いて構築した、サービス品質→総合満足度→応援ロイヤリティ→観戦ロイヤリティの因果モデルを示している。矢線近くの数値は標準化係数を表している。

この結果（標準化係数の大きさ）から、「ファンサービス・地域貢献」が総合満足度にもっとも影響を与えていることがわかる。すなわち、総合満足度を向上させるためには、ファンサービス・地域貢献活動が有効であることを示唆している。図4は、2020年1月下旬までの過去12年間の各チームの総合満足度スコア（潜

変数のスコアを100点満点に規準化したスコア）の経年変化を示している。直近の2020年1月下旬調査において、注目すべきチームの評価について考察する。

- ソフトバンク：2020年1月下旬調査では総合満足度1位（74.0）となった。2019年シーズンでもリーグ優勝は逃したものの、CS（クライマックスシリーズ）を勝ち進み、3年連続で日本一となった。チームの強さだけでなく、柳田悠岐選手、松田宣浩選手、今宮健太選手など魅力ある選手が健在である。球場でのファンサービス、地域貢献活動の取り組みも良いとされ、その結果が、チーム・選手の魅力、ファンサービス・地域貢献などが高い評価を得ており、総合力で高水準を維持している。
- 広島：2020年1月下旬調査では総合満足度2位（72.0）となり、前年までの2年連続の1位から順

位を落とした。2019年シーズンでは、セ・リーグでの成績が4位であったが、逆に言えば、総合満足度の高水準は維持されている。丸佳浩選手のFA移籍による影響が心配されたが、鈴木誠也選手、菊池涼介選手らの人気・主力選手は健在であること、球場の要素を含めたファンサービスに対する高評価が、高い満足度の維持の主要因であると考えられる。

- 西武：2020年1月下旬調査では総合満足度3位(71.4)となった。2年連続でリーグ優勝を成し遂げ、森友哉選手、山川穂高選手、源田壮亮選手らの活躍が光った。ただし、2年連続で日本シリーズ進出を逃していることが、総合満足度が伸びていない要因と考えられる。一方、球場の改修の一部が完成し、ファンサービスの高評価につながっている。改修はまだ継続中であり、今後のさらなるファンサービスの向上が期待される。
- 横浜DeNA：2020年1月下旬調査では4位(68.9)となった。近年のファンサービスの向上、地域貢献活動の活性化、球場改修の成果が出ており、回答者の声としても、「横浜市民の誇り」、「選手とファンとの一体感」などのコメントが目を引く。山崎康晃選手、宮崎敏郎選手、今永昇太選手らの魅力ある選手も存在している。さらに、優勝すれば、総合満足度1位となることが期待される。

以上のように、プロ野球チームの満足度指数化モデル、サービス品質や顧客満足度の数値化から、チームの取り組みに関するさまざまな「気づき」が得られる。良い成果をあげているチームは、ファンサービス・地域貢献、チーム成績、チーム・選手のバランスがよく、ファンとチームとの一体感がとれていることが読み取れる。このような数値化は、客観的な判断とさまざまなアクション、すなわち「事実に基づく管理」へとつながる。

3. 多変量推測統計の理論

3.1 多次元確率分布の principal points の統計的推定に関する研究

ある p 次元確率分布の k -principal points とは、その確率分布に従う確率変数ベクトルを \mathbf{X} とおくと、

$$E\left[\min_{j=1,\dots,k} \|\mathbf{X} - \boldsymbol{\gamma}_j\|^2\right] \quad (1)$$

を最小にする k 個の p 次元ベクトル $\boldsymbol{\gamma}_1^*, \dots, \boldsymbol{\gamma}_k^*$ として定義される (e.g., Tarpey et al. [8]). 一般に、1-principal point は分布の平均 $E[\mathbf{X}]$ になる。 k -

principal points は分布の平均の概念を1点から k 個の点に拡張したものと見ることができる。式(1)の表現自体は、オペレーションズ・リサーチ、品質管理、情報理論などのさまざまな分野で見受けられるものであり、その理論や応用に関して古くから多くの議論がなされている。本節では、特に統計学的な側面での研究に焦点を当て、確率分布の母数が未知な場合の k -principal points を推定する、という立場での近年の成果を紹介する。以下に述べる内容は Matsuura et al. [9] および Matsuura and Tarpey [10] の結果の一部である。

p 次元確率変数ベクトル \mathbf{X} の分布が、以下のような位置尺度分布族の確率密度関数 f をもつとする。 $\boldsymbol{\mu} \in R^p$ は位置母数ベクトル、 $\sigma \in (0, \infty)$ は尺度母数である。

$$f(\mathbf{x}; \boldsymbol{\mu}, \sigma) = \frac{1}{\sigma^p} g\left(\frac{\mathbf{x} - \boldsymbol{\mu}}{\sigma}\right) \quad (2)$$

for some function $g: R^p \rightarrow [0, \infty)$

このとき、 \mathbf{X} の分布の k -principal points $\boldsymbol{\gamma}_1^*, \dots, \boldsymbol{\gamma}_k^*$ は、関数 g (の定数倍) を確率密度関数としてもつ分布の k -principal points を $\boldsymbol{\delta}_1, \dots, \boldsymbol{\delta}_k$ とおくと、

$$\boldsymbol{\gamma}_j^* = \boldsymbol{\mu} + \sigma \boldsymbol{\delta}_j, \quad j = 1, \dots, k$$

と表される。ここで、母数 $\boldsymbol{\mu}, \sigma$ が未知であり、 \mathbf{X} の分布から大きさ n の無作為標本 $\mathbf{X}_1, \dots, \mathbf{X}_n$ が得られている場合を想定する。 $\boldsymbol{\mu}, \sigma$ の最尤推定量をそれぞれ $\hat{\boldsymbol{\mu}} = \hat{\boldsymbol{\mu}}(\mathbf{X}_1, \dots, \mathbf{X}_n)$ および $\hat{\sigma} = \hat{\sigma}(\mathbf{X}_1, \dots, \mathbf{X}_n)$ とおくと、

$$\hat{\boldsymbol{\gamma}}_j = \hat{\boldsymbol{\mu}} + \hat{\sigma} \boldsymbol{\delta}_j, \quad j = 1, \dots, k \quad (3)$$

は \mathbf{X} の分布の k -principal points の最尤推定量となるが、この最尤推定量は

$$E\left[\min_{j=1,\dots,k} \|\mathbf{X} - \hat{\boldsymbol{\gamma}}_j\|^2\right] \quad (4)$$

を最小にするとは限らない。そこで、式(3)を一般化して、 k -principal points の推定量を

$$\hat{\boldsymbol{\gamma}}_j = \hat{\boldsymbol{\mu}} + \hat{\sigma} \mathbf{a}_j, \quad j = 1, \dots, k$$

とおき、式(4)を最小にする $\mathbf{a}_1^*, \dots, \mathbf{a}_k^*$ を導出することを考える。その結果は、具体的には以下のように与えられる。

定理 無作為標本 $\mathbf{X}_1, \dots, \mathbf{X}_n$ とは独立に確率密度関数(2)をもつ分布に従う確率変数ベクトルを \mathbf{Y} とおく。 $\mathbf{U} = \frac{\mathbf{Y} - \hat{\boldsymbol{\mu}}}{\hat{\sigma}}$ 、 $W = \frac{\hat{\sigma}}{\sigma}$ とおき、 W 条件付きの \mathbf{U} の確率密度関数を $\psi(\mathbf{u}|W)$ と表記する。このとき、式(4)を最小にする $\mathbf{a}_1^*, \dots, \mathbf{a}_k^*$ は、以下の確率密度関数

$$h(z) = \frac{E[W^{p+2}\psi(zW|W)]}{E[W^2]}$$

をもつ分布の k -principal points で与えられる。

例 1 p 次元正規分布 $N_p(\boldsymbol{\mu}, \sigma^2 \boldsymbol{\Psi})$ から大きさ n の無作為標本 $\mathbf{X}_1, \dots, \mathbf{X}_n$ が観測されているとする。 $\boldsymbol{\mu}, \sigma$ は未知母数であるとし、 $\hat{\boldsymbol{\mu}} = \bar{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$ および $\hat{\sigma} = \sqrt{\frac{\sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}})^\top \boldsymbol{\Psi}^{-1} (\mathbf{X}_i - \bar{\mathbf{X}})}{np}}$ とおく。このとき、上の定理より、 k -principal points の推定量を $\hat{\gamma}_j = \hat{\boldsymbol{\mu}} + \hat{\sigma} \mathbf{a}_j$, $j = 1, \dots, k$ とおくと、式 (4) を最小にする $\mathbf{a}_1^*, \dots, \mathbf{a}_k^*$ は、 $\sqrt{\frac{(n+1)p}{(n-1)p+2}} t_{(n-1)p+2}(\mathbf{0}, \boldsymbol{\Psi})$ (位置母数ベクトル $\mathbf{0}$, 尺度母数行列 $\boldsymbol{\Psi}$ をもつ自由度 $(n-1)p+2$ の多変量 t 分布の $\sqrt{\frac{(n+1)p}{(n-1)p+2}}$ 倍の分布) の k -principal points であることが示される。

例 2 以下のような確率密度関数

$$f(x; \mu, \sigma) = \begin{cases} \frac{1}{\sigma} e^{-\frac{x-\mu}{\sigma}}, & x \geq \mu \\ 0, & x < \mu \end{cases}$$

をもつ(一次元)確率分布(位置母数 μ , 尺度母数 σ をもつ指数分布) から大きさ n の無作為標本 X_1, \dots, X_n が観測されているとする。 μ, σ は未知母数であるとし、 $\hat{\mu} = \min\{X_1, \dots, X_n\}$ および $\hat{\sigma} = \bar{X} - \hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i - \hat{\mu}$ とおく。このとき、上の定理より、 k -principal points の推定量を $\hat{\gamma}_j = \hat{\mu} + \hat{\sigma} \mathbf{a}_j$, $j = 1, \dots, k$ とおくと、式 (4) を最小にする $\mathbf{a}_1^*, \dots, \mathbf{a}_k^*$ は、以下の確率密度関数

$$h(z) = \begin{cases} (1 + \frac{z}{n})^{-(n+2)}, & z \geq 0 \\ (1 - z)^{-(n+2)}, & z < 0 \end{cases}$$

をもつ分布の k -principal points であることが示される。

3.2 Seemingly unrelated regression モデルにおける最良共変推定量に関する研究

以下の定式化で表されるような互いに相関がある複数個 (p 個) の重回帰モデルを想定する。

$$\mathbf{y}_i = \mathbf{X}_i \boldsymbol{\beta}_i + \boldsymbol{\varepsilon}_i, \quad i = 1, \dots, p$$

$$\text{with } \begin{cases} E[\boldsymbol{\varepsilon}_i] = \mathbf{0}, V[\boldsymbol{\varepsilon}_i] = \sigma_i^2 \mathbf{I}_m, & i = 1, \dots, p, \\ \text{Cov}[\boldsymbol{\varepsilon}_i, \boldsymbol{\varepsilon}_j] = \sigma_i \sigma_j \rho_{ij} \mathbf{I}_m, & i \neq j. \end{cases}$$

ただし、 $\mathbf{y}_i, \mathbf{X}_i, \boldsymbol{\beta}_i, \boldsymbol{\varepsilon}_i$ はそれぞれ i 番目の重回帰モデルにおける目的変数ベクトル (m 次元ベクトル)、説明変数行列 ($m \times k_i$ 行列)、偏回帰係数ベクトル (k_i 次元ベクトル)、誤差ベクトル (m 次元ベクトル) であり、 σ_i^2 は i 番目の重回帰モデルにおける誤差分散、 ρ_{ij}

は i 番目と j 番目の重回帰モデルにおける誤差間の相関係数である。また、 \mathbf{I}_m は m 次元単位行列を表す。

上記のモデルは Seemingly Unrelated Regression モデルと呼ばれ(以下、SUR モデル)、Zellner [11] をはじめとして多くの議論がなされている。このモデルは、 $n = mp, k = \sum_{i=1}^p k_i$ とおき、さらに

$$\begin{aligned} \mathbf{y} &= (\mathbf{y}_1^\top, \dots, \mathbf{y}_p^\top)^\top : n \text{ 次元ベクトル,} \\ \mathbf{X} &= \text{diag}\{\mathbf{X}_1, \dots, \mathbf{X}_p\} : n \times k \text{ ブロック対角行列,} \\ \boldsymbol{\beta} &= (\boldsymbol{\beta}_1^\top, \dots, \boldsymbol{\beta}_p^\top)^\top : k \text{ 次元ベクトル,} \\ \boldsymbol{\varepsilon} &= (\boldsymbol{\varepsilon}_1^\top, \dots, \boldsymbol{\varepsilon}_p^\top)^\top : n \text{ 次元ベクトル,} \\ \boldsymbol{\Sigma}_d &= \text{diag}\{\sigma_1, \dots, \sigma_p\} : p \times p \text{ 対角行列,} \\ \boldsymbol{\Lambda} &= (\rho_{ij})_{1 \leq i, j \leq p} : p \times p \text{ 行列 (ただし } \rho_{ii} = 1) \end{aligned}$$

とおくと、以下のように書き直すことができる。

$$\mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon} \text{ with } E[\boldsymbol{\varepsilon}] = \mathbf{0}, V[\boldsymbol{\varepsilon}] = \boldsymbol{\Sigma}_d \boldsymbol{\Lambda} \boldsymbol{\Sigma}_d \otimes \mathbf{I}_m. \quad (5)$$

ただし、 \otimes はクロネッカー積を表す。

この SUR モデルにおける偏回帰係数ベクトル $\boldsymbol{\beta}$ の推定について、さまざまなものが提案されている。本節では、近年、共変推定の観点で議論した Kurata and Matsuura [12] および Matsuura and Kurata [13] の結果の一部を紹介する。

SUR モデル (5) において、誤差ベクトル $\boldsymbol{\varepsilon}$ が以下のような楕円対称分布の確率密度関数をもつとする。

$$f(\boldsymbol{\varepsilon}) = |\boldsymbol{\Sigma}_d \boldsymbol{\Lambda} \boldsymbol{\Sigma}_d|^{-\frac{m}{2}} g(\boldsymbol{\varepsilon}^\top (\boldsymbol{\Sigma}_d \boldsymbol{\Lambda} \boldsymbol{\Sigma}_d \otimes \mathbf{I}_m)^{-1} \boldsymbol{\varepsilon})$$

for some function $g : [0, \infty) \rightarrow [0, \infty)$.

また、相関係数行列 $\boldsymbol{\Lambda}$ は既知であるとする。このとき、次のような共変性をもつ $\boldsymbol{\beta}$ の推定量のクラスを想定する。

$$\begin{cases} \mathbf{y}_i \rightarrow a_i \mathbf{y}_i + \mathbf{X}_i \mathbf{c}_i \\ \hat{\boldsymbol{\beta}}_i \rightarrow a_i \hat{\boldsymbol{\beta}}_i + \mathbf{c}_i \end{cases} \quad (6)$$

for any $a_i \in (0, \infty)$, $\mathbf{c}_i \in R^{k_i}$, $i = 1, \dots, p$.

この共変推定量のクラスにおいて、以下の損失関数 $L(\hat{\boldsymbol{\beta}})$ の期待値を最小にする推定量(最良共変推定量)を導出することを考える。

$$L(\hat{\boldsymbol{\beta}}) = (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \mathbf{X}^\top (\boldsymbol{\Sigma}_d \boldsymbol{\Lambda} \boldsymbol{\Sigma}_d \otimes \mathbf{I}_m)^{-1} \mathbf{X} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}).$$

定理 準備として、 $i = 1, \dots, p$ について

$$\begin{aligned} \mathbf{e}_i &= (\mathbf{I}_m - \mathbf{X}_i (\mathbf{X}_i^\top \mathbf{X}_i)^{-1} \mathbf{X}_i^\top) \mathbf{y}_i, \\ \mathbf{u}_i &= \frac{\mathbf{e}_i}{\|\mathbf{e}_i\|} \end{aligned}$$

とおき, $i, j = 1, \dots, p$ について

$$t_{ij} = E_{(\beta, \Sigma_d) = (\mathbf{0}, \mathbf{I}_p)}[\|e_i\| \|e_j\| | (\mathbf{u}_1, \dots, \mathbf{u}_p)]$$

とおく. ただし, $E_{(\beta, \Sigma_d) = (\mathbf{0}, \mathbf{I}_p)}[\cdot]$ は, $\beta = \mathbf{0}$, $\Sigma_d = \mathbf{I}_p$ の下で期待値を取ることを表す. さらに,

$$\begin{aligned} \mathbf{S} &= \text{diag}\{\|e_1\|, \dots, \|e_p\|\} : p \times p \text{ 対角行列,} \\ \mathbf{T} &= (t_{ij})_{1 \leq i, j \leq p} : p \times p \text{ 行列} \end{aligned}$$

とおくと, 共変性 (6) を満たす β の推定量のクラスにおいて, $E[L(\hat{\beta})]$ を最小にする最良共変推定量は以下のように表される.

$$\hat{\beta} = \left[\mathbf{X}^\top \{ \mathbf{S}^{-1} (\mathbf{T} \circ \Lambda^{-1}) \mathbf{S}^{-1} \otimes \mathbf{I}_m \} \mathbf{X} \right]^{-1} \mathbf{X}^\top \{ \mathbf{S}^{-1} (\mathbf{T} \circ \Lambda^{-1}) \mathbf{S}^{-1} \otimes \mathbf{I}_m \} \mathbf{y}.$$

ただし, \circ はアダマール積 ($\mathbf{A} = (a_{ij})$, $\mathbf{B} = (b_{ij})$ のとき $\mathbf{A} \circ \mathbf{B} = (a_{ij} b_{ij})$) を表す.

このほか, Matsuura and Kurata [13] では, 誤差ベクトル ε の分散共分散行列に対応する $\Sigma_d \Lambda \Sigma_d$ に関する最良共変推定量についても議論している.

4. 実験計画の構成

4.1 よい計画とは

これまでの節は, 収集されているデータの解析やそのための理論という視点での研究紹介である. この節では, データの計画的な収集という視点で, 実験計画法 (Design of Experiments) に関する研究を紹介する.

収集されたデータは, 対象とする系に介入せずに観察することにより得る観察データと, 対象とする系に介入し意図的に条件など管理して得る実験データに大別できる. 観察データは応答の値を相関のある変数から予測するという目的には適するが, 応答に影響を与える要因の分析や, 応答と要因の因果構造の定量化には一般に役立たない. 一方実験データは, 適切に実験の場を管理してデータを収集し, それを丹念に解析することにより, 予測のみならず, 要因分析, 因果構造の定量化に役立つ.

実際に実験を計画する段階では, 実験の意図, 実現可能性, 推定精度の確保などさまざまな点を考慮する. 本稿ではオペレーションズ・リサーチ誌という点に鑑み, 筆者の浅学さを顧みず, 実験計画の構成を最適化問題という視点で記述することを試みる.

応答 y とその因子 x_1, \dots, x_p の関係を, 実験データにより近似的に推定する. そのために, 因子 x_1, \dots, x_p の水準を意図的に変化させ, それ以外で応答 y に影響を与えらると思われるものはできる限り一定の水準と

し, さらに実験順序の無作為化を行う. 応答 y と因子 x_1, \dots, x_p に関する n 組のデータをこのような実験により収集し, 応答と因子の真の応答関数を推定する. その際, 応答 y が x_1, \dots, x_p の線形結合と誤差との和で表されると仮定する状況が最も基本的である. 応答 y について, n 次元ベクトル $\mathbf{y} = (y_1, \dots, y_n)^\top$ の真の応答関数が, 一般平均と因子の効果からなるベクトル β により

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon \quad (7)$$

で表されるとする. ただし, 一般平均を μ , x_1, \dots, x_p の効果からなるベクトルを $\beta_1 = (\beta_1, \dots, \beta_p)^\top$ とすると $\beta = (\mu, \beta_1)^\top$ である. また, \mathbf{x}_j を因子 x_j の n 組の実験水準からなる n 次元ベクトル, $\mathbf{X}_1 = (\mathbf{x}_1, \dots, \mathbf{x}_p)$ を $(n \times p)$ 計画行列とし, $\mathbf{1}$ を要素がすべて 1 の n 次元ベクトルとすると, $\mathbf{X} = (\mathbf{1}, \mathbf{X}_1)$ である. さらに, ε は n 次元誤差ベクトルであり, 独立性, 不偏性, 等分散性を仮定する.

一般平均と因子の効果からなる β の最小 2 乗推定量は,

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (8)$$

で与えられる. この推定に, 何らかの意味でのよさをもたらし計画がよい計画となる. 推定のよさを記述したうで, それを好ましくする計画を求める問題が, 実験計画の構成問題である.

4.2 2 水準一部実施要因計画

因子 x_1, \dots, x_p の中から効果の大きな因子を絞り込むために, 少数回の実験で式 (8) の推定を行うことを考える. このように因子 x_1, \dots, x_p の中から, 効果のある因子を絞り込む目的で実施される実験をスクリーニング実験と呼ぶ. この場合には, 実験回数が少ないことに加え, 因子の交互作用 (組合せ効果) があつたとしてもその影響を大きく受けない計画がよい.

これには, 因子を 2 水準とした一部実施要因計画をよく用いる. 因子 x_1, \dots, x_p が $-1, 1$ という 2 水準で表されるとすると, p 個の因子の水準の組合せは 2^p となる. このような水準の組合せをすべて実施する, すなわち実験回数 $n = 2^p$ の計画が要因計画である. また, その一部分だけを実施するのが一部実施要因計画である.

計画行列 $\mathbf{X}_1 = (\mathbf{x}_1, \dots, \mathbf{x}_p)$ が主効果の推定のための 2 水準直交計画であり, それぞれの列 \mathbf{x}_i ($i = 1, \dots, p$) に $-1, 1$ が $n/2$ ずつ含まれるものとする. 応答変数 y について, 真の応答関数が, 一般平均 μ , p 因子の主効果からなるベクトル β_1 , すべての交

相互作用の効果からなる $(p(p-1)/2)$ 次元ベクトル β_2 により

$$\mathbf{y} = \mu\mathbf{1} + \mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2 + \varepsilon \quad (9)$$

で表されるとする。ベクトルの要素ごとの積であるアダマール積を \circ とすると、 \mathbf{x}_i と \mathbf{x}_j の交互作用列は $\mathbf{x}_i \circ \mathbf{x}_j$ となる。したがって、交互作用からなる \mathbf{X}_2 は、

$$\mathbf{X}_2 = (\mathbf{x}_1 \circ \mathbf{x}_2, \mathbf{x}_1 \circ \mathbf{x}_3, \dots, \mathbf{x}_{p-1} \circ \mathbf{x}_p)$$

となる。

解析のモデルには主効果 β_1 を含めているものの、交互作用 β_2 は含めていないものとする。主効果の推定量 $\hat{\beta}_1$ の期待値は、式 (8) に式 (9) を代入して期待値を考えると、別名行列と呼ぶ $(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_2$ を用いて

$$E(\hat{\beta}_1) = \beta_1 + (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_2 \beta_2 \quad (10)$$

となる。因子 i ($i = 1, \dots, p$) の主効果を表す列 \mathbf{x}_i と、因子 j, k の交互作用の列 $\mathbf{x}_j \circ \mathbf{x}_k$ との相関係数を $r(i, j \circ k)$ とすると、 $(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_2$ は、行に因子 i の主効果を、列に交互作用 $\mathbf{x}_j \circ \mathbf{x}_k$ 対応させた相関係数 $r(i, j \circ k)$ からなる $(p \times (p(p-1)/2))$ の行列となる。

この主効果と交互作用の相関構造は、(i) 実験数 n が 2 のべき乗数の場合と、(ii) 実験数 n が 2 のべき乗数ではなく 4 の倍数の場合で異なる。実験数 n が 2 のべき乗数の場合には、主効果列 \mathbf{x}_i と交互作用列 $\mathbf{x}_j \circ \mathbf{x}_k$ が直交するか、完全に交絡するかのいずれかになる。すなわち、 $r(i, j \circ k)$ は 0 か 1 かのいずれかになる。

これに対して、実験数 n が 2 のべき乗数ではなく 4 の倍数の場合には、一般に、0 か 1 かではなく部分的に交絡する。たとえば $n = 12$ の場合には、 $|r(i, j \circ k)| = 1/3$ ($i \neq j \neq k \neq i$)、 $|r(i, j \circ k)| = 0$ ($i = j$ or $i = k$) となることはよく知られている。また、 $n = 20$ の場合には、 $|r(i, j \circ k)|$ が 0, 1/5, 3/5 のいずれかになる。このうち、3/5 となる組合せは、交互作用による影響が大きく表れるために避けるのがよい。このために、 x_1, \dots, x_{19} を頂点、交互作用を辺とするグラフを考え、このグラフの構造をもとに 3/5 となる組合せを避ける選択方法がある [14]。

4.3 過飽和実験計画の最適性

4.3.1 2水準過飽和計画とその評価基準

前述の一部実施要因計画では、因子 x_1, \dots, x_p の列ベクトルが直交するよう計画を構成しているので、 n 回

の実験で最大で $n-1$ の因子の効果も推定できる。この直交性の制約を緩め、 n 次元空間内で n 以上の 2 水準ベクトルをできる限り直交に近い形で列挙する計画が過飽和実験計画である。その意味で過飽和実験計画には、割付けられる因子数が実験回数よりも多いというよさがある。そのため、効果の推定の前段階として、多数の因子の中から少数の重要な因子を選別するために用いる。要素が -1 または 1 からなる n 次元ベクトルを \mathbf{x}_i とするとき、 $n \times p$ 行列 $\mathbf{X}_1 = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p)$ について、 $p \geq n$ の場合には 2 水準過飽和実験計画と呼ぶ。

過飽和実験計画は、Satterthwaite [15] によって基本的考え方が提示され、Booth and Cox [16] によって系統的な構成方法が提示された。過飽和実験計画が広く研究されるようになったのは、Lin [17] が示した構成方法である。この構成方法は、単純、かつ、いくつかの意味での正当性がある。

過飽和実験計画の場合には、列数が行数よりも多くなるためにすべての列間での直交性を確保できない。一方で、因子の効果をより精密に推定する立場からは、列間に直交性に近い性質が成り立つことが好ましい。そこでまず、列間にどの程度直交性が成り立っているのかの評価尺度として、二つの列間の内積の 2 乗

$$s_{ij} = \mathbf{x}_i^\top \mathbf{x}_j \quad (11)$$

を導入する。また計画全体の評価は、2 列間の内積の 2 乗値についてすべての列の対に関する平均

$$E(s^2) = \frac{1}{p(p-1)/2} \sum_{1 \leq i < j \leq p} s_{ij}^2 \quad (12)$$

を用いる。この基準は、初期の段階において Booth and Cox [16] がランダムサーチにより過飽和実験計画を構成することを考察したなどから、内積の 2 乗の期待値 $E(s^2)$ として略すことが多い。

この $E(s^2)$ について、実験数 n が 4 の倍数で、 -1 と 1 が $n/2$ ずつ含まれ互いに直交する列ベクトル \mathbf{x}_i ($i = 1, \dots, n-1$) からなる $n \times (n-1)$ の行列を $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n-1})$ とすると、 -1 と 1 が $n/2$ ずつ含まれる任意の n 次元ベクトル \mathbf{x} について、

$$\sum_{i=1}^{n-1} (\mathbf{x}^\top \mathbf{x}_i)^2 = \mathbf{x}^\top \mathbf{X} \mathbf{X}^\top \mathbf{x} = n^2 \quad (13)$$

が成り立つ。また、 $E(s^2)$ の下界

$$E(s^2) \geq \frac{n^2(p-n+1)}{(n-1)(p-1)} \quad (14)$$

を Nguyen [18], Tang and Wu [19] は独立に導いている。これは、実験数 n 、列数 p の過飽和実験計画において、 $E(s^2)$ が式 (14) の右辺に一致していれば、その過飽和実験計画は $E(s^2)$ が最小である。この $E(s^2)$ の下界について、Tang and Wu [19] は因子 \mathbf{x}_i の 4 次モーメントからなる $\mathbf{X}_1^\top \mathbf{X}_1 \mathbf{X}_1^\top \mathbf{X}_1$ の固有値の 2 乗和に関する最小化問題に帰着させている。さらに、Jones and Majumdar [20] では、 -1 と 1 が $n/2$ ずつという制約を外した場合の $E(s^2)$ の下界について示している。その誘導には、2 部グラフに関する Turan の定理を用いている。

4.3.2 多水準過飽和計画

2 水準過飽和計画の構成に関する性質は、多水準過飽和計画、混合水準過飽和計画に展開できる。以下では、簡単のために 3 水準列について説明する。たとえば、 $n = 9$ で四つの 3 水準列からなる計画

$$\mathbf{X} = (\mathbf{x}_1^3, \mathbf{x}_2^3, \mathbf{x}_3^3, \mathbf{x}_4^3) = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 2 & 2 \\ 1 & 3 & 3 & 3 \\ 2 & 1 & 2 & 3 \\ 2 & 2 & 3 & 1 \\ 2 & 3 & 1 & 2 \\ 3 & 1 & 3 & 2 \\ 3 & 2 & 1 & 3 \\ 3 & 3 & 2 & 1 \end{pmatrix} \quad (15)$$

を考える。列間の内積の 2 乗 s_{ij}^2 、計画全体の平均 $E(s^2)$ の多水準計画への拡張として、 $(\mathbf{x}_i^3, \mathbf{x}_j^3)$ の直交性をピアソンの χ^2 統計量

$$\begin{aligned} \chi^2(\mathbf{x}_i^3, \mathbf{x}_j^3) &= \sum_{k=1}^3 \sum_{l=1}^3 \frac{(\text{実測度数} - \text{期待度数})^2}{\text{期待度数}} \\ &= \sum_{k=1}^3 \sum_{l=1}^3 \frac{(n(k, l) - n/9)^2}{n/9} \end{aligned} \quad (16)$$

により評価する。ただし、 $n(k, l)$ は \mathbf{x}_i と \mathbf{x}_j において、 (k, l) となる組合せの個数である。前述の例では、すべての列が互いに直交しているので、どの 2 列の組合せを求めても $n(k, l) = 1$ ($k, l = 1, 2, 3$) となる。

1, 2, 3 がそれぞれ 3 ずつ含まれる $n = 9$ の任意のベクトル \mathbf{x}^3 と $(\mathbf{x}_1^3, \mathbf{x}_2^3, \mathbf{x}_3^3, \mathbf{x}_4^3)$ について $\chi^2(\mathbf{x}_i^3, \mathbf{x}_j^3)$ の和を考えると、

$$\sum_{i=1}^4 \chi^2(\mathbf{x}_i^3, \mathbf{x}_j^3) = 18 \quad (17)$$

が成立する。一般に、 l 水準の $(n \times (n-1)/l)$ 計画行列 \mathbf{X} が、互いに直交するベクトル $(\mathbf{x}_1^l, \dots, \mathbf{x}_{(n-1)/2}^l)$ からなるとき、 $1, \dots, l$ がそれぞれ同数ずつ含まれる

任意のベクトル \mathbf{x}^3 について、

$$\sum_{i=1}^{(n-1)/2} \chi^2(\mathbf{x}_i^l, \mathbf{x}^l) = 2n \quad (18)$$

が成立する。これは、直交列との二乗和が一定という意味で式 (13) の多水準への拡張になっている。

加えて、式 (12) による $E(s^2)$ に関する下界と同様に、 χ^2 値の合計に関する下界がある。任意の過飽和計画 \mathbf{X} に対して、その χ^2 値の合計 $\chi^2(\mathbf{X})$ について、

$$\chi^2(\mathbf{X}) \geq \frac{1}{2}v(v-1)n(n-1) \quad (19)$$

が成立する [21]。ただし、 v は過飽和度であり、列数 \times (水準数 $- 1$) を $n-1$ で除したものである。このように、多水準、混合水準の過飽和計画において、 χ^2 を列間の直交性評価に用いることにより、内積の 2 乗 s_{ij}^2 で成り立つ性質と同様の性質が導ける。これらの誘導では、計画行列 \mathbf{X} に関する制約付き最適化問題に帰着させる接近法が多数用いられている。

4.4 最適計画

式 (7) と等しい解析モデルを用いる場合には、 $\hat{\beta}$ の期待値は母数 β に等しく、 $E(\hat{\beta}) = \beta$ という不偏性が成り立つ。そのため、推定量 $\hat{\beta}$ の分散共分散行列

$$V(\hat{\beta}) = (\mathbf{X}^\top \mathbf{X})^{-1} \sigma^2 \quad (20)$$

を好ましくする計画がよい計画となる。一般に、多変量の推定量の分散共分散行列について、スカラー関数での評価指標として分散共分散行列の行列式による一般化分散が用いられる。この値が小さいほど推定が好ましい。また、この一般化分散 $\det(V(\hat{\beta}))$ において、 σ^2 は実験を適用する場によって決まり、計画のよしあしで解消できるものではない。したがって、推定の分散を好ましくするという点では、 $\det(V(\hat{\beta})) = \det(\mathbf{X}^\top \mathbf{X}^{-1} \sigma^2)$ を小さくするために、

$$\det(\mathbf{X}^\top \mathbf{X}) \quad (21)$$

が大きい計画がよい計画となる。これは行列式 (determinant) を用いているという点で、 D 最適基準と呼ぶ。計画行列 \mathbf{X} を構成するもとなる実験水準 x_{ij} については、実験可能な領域 $R(\mathbf{x})$ 内に選ぶ必要がある。たとえば、因子ごとに上限と下限がある場合が典型である。したがって、計画 \mathbf{X} の構成問題は、

$$x_{ij} \in R(\mathbf{x}) \quad (22)$$

の制約のもとに、 $\det(\mathbf{X}^\top \mathbf{X})$ を最大化する \mathbf{X} を求める問題となる。

実験可能な領域 $R(\mathbf{x})$ が、それぞれの因子に対して上限と下限が規定されている場合において、応答 y が因子 x_1, \dots, x_p の線形結合と誤差との和で表現されるならば、 $\det(\mathbf{X}^T \mathbf{X})$ 最適な計画の構成は簡単である。これには、因子の上限と下限に水準をとり、それぞれの因子からなるベクトルが直交するように選べばよい。

一方、 $R(\mathbf{x})$ がそれぞれの因子に対して上限と下限により規定されている場合でも、応答 y が因子 x_1, \dots, x_p の線形結合と誤差との和ではなく、たとえば、 x_j^2 や $x_j x_k$ を含む 2 次モデルの場合には、計画行列 \mathbf{X} に x_{i1}, \dots, x_{ip} に加え $x_{ij}^2, x_{ij} x_{ik}$ などが含まれるので、その計算が煩雑になり、 D 最適な計画の構成が困難になる。さらに、 $R(\mathbf{x})$ が凸集合にならない場合にも、状況は困難になる。加えて、触媒の種類を取り上げたような場合には因子が名義尺度となるため、その計算が一般に煩雑になる。

このような最適計画の構成問題は、Kiefer and Wolfowitz [22] から議論が始まり、最適性の基準、最適計画の構成方法などが論じられている。提案された当初は、数理的な問題として取り扱われていた。2000 年ころから統計ソフトウェアへの実装が始まり、現在では多くのソフトウェアでも導入されつつある。

5. おわりに

真の意味でのデータサイエンスの実践には、対象となる場の理解とモデル化、情報機器の活用によるデータの収集、統計的データ解析という三つの融合が欠かせない。管理工学科の教育体系は、これらの三つの領域をよいバランスで提供している。たとえば、経営、経済、オペレーションズ・リサーチ、情報、統計などである。今後も統計について、隣接分野との融合を考えつつ理論と応用について教育研究を進めていきたい。

参考文献

[1] 鈴木秀男, 『この一冊で合格! QC 検定 3 級集中テキスト&問題集』, ナツメ社, 2015.
 [2] 柴田頼仁, 鈴木秀男, “LSTM を用いた球種子測モデルの構築,” 統計数理研究所共同研究レポート 423, **6**, pp. 1–8, 2019.
 [3] 白戸豪大, 鈴木秀男, “ポロノイ図の空円性を用いたサッカーの守備構造評価,” 統計数理研究所共同研究レポート 423, **6**, pp. 83–90, 2019.
 [4] 鈴木秀男, 『顧客満足度向上のための手法—サービス品質の獲得—』, 日科技連出版社, 2010.
 [5] 鈴木秀男, 『サービス品質の構造を探る—プロ野球の事例

から学ぶ—』, 日本規格協会, 2011.

[6] 鈴木秀男, “「事実に基づく管理」のすゝめ—「プロ野球チームの満足度調査」から学ぶ—,” *Direct Marketing Review*, **17**, pp. 6–14, 2018.
 [7] 鈴木秀男, 「プロ野球のサービスの満足度調査: 慶應義塾大学理工学部 管理工学科 鈴木研究室 HP」, <http://lab.ae.keio.ac.jp/~hsuzuki/baseball0901/index.html> (2020 年 10 月 20 日閲覧)
 [8] T. Tarpey, L. Li and B. Flury, “Principal points and self-consistent points of elliptical distributions,” *Annals of Statistics*, **23**, pp. 103–112, 1995.
 [9] S. Matsuura, H. Kurata and T. Tarpey, “Optimal estimators of principal points for minimizing expected mean squared distance,” *Journal of Statistical Planning and Inference*, **167**, pp. 102–122, 2015.
 [10] S. Matsuura and T. Tarpey, “Optimal principal points estimators of multivariate distributions of location-scale and location-scale-rotation families,” *Statistical Papers*, **61**, pp. 1629–1643, 2020.
 [11] A. Zellner, “An efficient method of estimating seemingly unrelated regressions and tests for aggregation bias,” *Journal of the American Statistical Association*, **57**, pp. 348–368, 1962.
 [12] H. Kurata and S. Matsuura, “Best equivariant estimator of regression coefficients in a seemingly unrelated regression model with known correlation matrix,” *Annals of the Institute of Statistical Mathematics*, **68**, pp. 705–723, 2016.
 [13] S. Matsuura and H. Kurata, “Covariance matrix estimation in a seemingly unrelated regression model under Stein’s loss,” *Statistical Methods and Applications*, **29**, pp. 79–99, 2020.
 [14] Y. Oishi and S. Yamada, “Evaluation on alias relation of interactions in Plackett Burman design and its application to guide assignments,” *Total Quality Science*, **5**, pp. 81–91, 2020.
 [15] F. E. Satterthwaite, “Random balance experimentation,” *Technometrics*, **1**, pp. 111–137, 1959.
 [16] K. H. V. Booth and D. R. Cox, “Some systematic supersaturated designs,” *Technometrics*, **4**, pp. 489–495, 1962.
 [17] D. K. J. Lin, “A new class of supersaturated designs,” *Technometrics*, **35**, pp. 28–31, 1993.
 [18] N. K. Nguyen, “An algorithmic approach to constructing supersaturated designs,” *Technometrics*, **38**, pp. 67–73, 1996.
 [19] B. Tang and C. F. J. Wu, “A method for constructing supersaturated designs and its E_s^2 optimality,” *Canadian Journal of Statistics*, **25**, pp. 191–201, 1997.
 [20] B. Jones and D. Majumdar, “Optimal supersaturated designs,” *Journal of American Statistical Association*, **109**, pp. 1592–1600, 2014.
 [21] S. Yamada and T. Matsui, “Optimality of mixed-level supersaturated designs,” *Journal of Statistical Planning and Inference*, **104**, pp. 459–468, 2002.
 [22] J. Kiefer and J. Wolfowitz, “Optimum Designs in Regression Problems,” *Annals of Mathematical Statistics*, **30**, pp. 271–294, 1959.