

第1章 統計的学習の考え方

1.1 導入

データサイエンスにおいて、データの視覚化や整理¹は重要であるが、主な課題はデータの数理的な解析である。データの不確実性(ランダム性)のモデル化や定量化を目的とするとき、これらのことを扱う分野を統計的推測/統計的学習といい、大規模データを用いた予測に重点をおく分野を機械学習またはデータマイニングという。何に重きをおくかの違いで、ふたつの文化²が存在する。

データのモデリングにはふたつの目的がある。

- (1) ある変量の将来の値を正確に予測すること。
- (2) データの中にある未知の面白いパターンを発見すること。

この目的を達成するために、数理科学の3つの重要な支柱からの知識を利用する。

- (1) 関数の近似: データに対して数理モデルを作ることは、ある変数が他の変数にどのように依存/関係しているかを理解することである。変数間の関連を表現する最も自然なやり方は、関数もしくは写像を用いることである。この関数ないしは写像が解析者³には完全には特定できないとき、データに基づきこの関数の近似を試みる。
- (2) 最適化: 統計的モデル⁴が与えられたとき、この族の中で最もよい確率分布をみつきたい。このために、効率的な探索と最適化の手続きが必要となる。最適化のステップは観測データへの関数のあてはめや校正と考えることができる。このステップでは、最適化アルゴリズムや効率的なプログラミングが必要となる。

¹いわゆる記述統計学の重要な内容。

²Breiman (2001): Statist Sci 16(3), pp. 199-231 を参照のこと。

³統計家、データサイエンティストなどデータを分析する人々の総称として解析者とよぶことにする。

⁴統計的推測では、確率分布の族を統計的モデルという。「モデル」という言葉が何をさしているかを文脈から正確に理解することが大切である。

- (3) 確率論と統計的推測: モデルをあてはめるデータはランダムな過程の実現値とみなす. ここでは, 分布や確率が重要になる. 将来の予測を行うときに内包する不確実性を定量化するために, 確率論や統計的推測論の知識が肝要である.

多次元のデータを扱う場合には, これらの知識に加えて, 線型代数の知識が重要となる.

1.2 教師あり学習と教師なし学習

入力ベクトル x が与えられたとき, 出力 y の予測をすることが機械学習の主な目的である. たとえば, x をデジタル化された署名データとし, y を 2 値の出力で, 署名が「本人のもの」か「本人のものではない」のいずれかを表すとする. 別の例では, 入力 x は妊婦の体重と喫煙嗜好とし, 出力 y を生まれてくる新生子の体重とする. 入力 x による出力 y の予測値を関数 g を用いて, $g(x)$ と表す. g の定義域は, 入力 x の取りうる値の集合を含むものであり, 終域は y の取りうる値の集合である. 関数 g のことを予測関数とよぶことにする. 関数 g は, データがもつランダム性を除いた入力 x と出力 y の関連性についての情報を内包するものである. 出力 y が実数値をとるとき, 出力から入力を予測することを回帰問題といい, 出力 y が 2 値もしくは有限集合に値をとるとき, 判別問題⁵という.

入力 x による出力 y の予測値を $\hat{y} := g(x)$ と書いたとき, この予測値の精度を損失関数 $\text{Loss}(y, \hat{y})$ で測る. 損失関数は非負値関数で, 精度がわるいほどおおきな値をとるものである. 2 値の判別問題では, 損失関数を 0-1 損失とすることが基本的である. すなわち

$$\text{Loss}(y, \hat{y}) = \mathbb{1}\{y \neq \hat{y}\} = \begin{cases} 1 & (y \neq \hat{y}) \\ 0 & (y = \hat{y}) \end{cases}$$

である. 後で 2 値の判別問題でも別の損失関数を考えることになる. 回帰問題では, 2 乗損失

$$\text{Loss}(y, \hat{y}) = (y - \hat{y})^2$$

が最もよく使われる.

すべての入力と出力の組 (x, y) に対して正確な予測をすることができる予測関数を見つけることは一般に不可能である. これはデータのもつランダム性に由来する. さきほどの例において同じ値の入力に対して異

⁵分類問題ともいう.

なる出力をもつ組が存在することを考えればよい。すなわち各 (x, y) はある同時分布に従う確率ベクトル (X, Y) の実現値と考えることにする。このことで予測関数 g の精度を期待損失 (誤差/リスク) で定量化する。すなわち

$$\text{Err}(g) := E[\text{Loss}(Y, g(X))]$$

である。ただし期待値 $E[\cdot]$ は (X, Y) の同時分布に関するものである。0-1 損失を用いた判別問題では、誤差は誤差判別の確率となる。すなわち

$$\text{Err}(g) = \Pr[Y \neq g(X)]$$

である。この文脈では予測関数のことを判別器とよぶ。

確率ベクトル (X, Y) の同時分布と損失関数を与えられると、最良の判別器を求めることができる。 $c \in \mathbb{N}$, $c \geq 2$ とし、 Y は $\{0, 1, \dots, c-1\}$ に値を取るとする。このとき

$$g^* \in \arg \min_g E[\text{Loss}(Y, g(\mathbf{X}))]$$

を誤差を最小とする判別器とする。すなわち、任意の判別器 g に対して

$$E[\text{Loss}(Y, g(\mathbf{X}))] \geq E[\text{Loss}(Y, g^*(\mathbf{X}))]$$

が成立している。

$X = \mathbf{x}$ が与えられたとき、 $g^*(\mathbf{x})$ を具体的に求めてみる。そのために、 $X = \mathbf{x}$ が与えられたときの Y の条件付き p.d.f. (または p.m.f.) を $p^*(y|\mathbf{x})$ と書く。すなわち

$$p^*(y|\mathbf{x}) = \Pr(Y = y | \mathbf{X} = \mathbf{x})$$

である。このとき

$$E[\text{Loss}(Y, g(\mathbf{X}))] = \Pr[Y \neq g(\mathbf{X})] = 1 - \Pr[Y = g(\mathbf{X})]$$

なので

$$\begin{aligned} g^*(\mathbf{x}) &= \arg \min_{y \in \{0, 1, \dots, c-1\}} E[\text{Loss}(Y, y) | \mathbf{X} = \mathbf{x}] \\ &= \arg \min_{y \in \{0, 1, \dots, c-1\}} \{1 - \Pr(Y = y | \mathbf{X} = \mathbf{x})\} \\ &= \arg \max_{y \in \{0, 1, \dots, c-1\}} \Pr(Y = y | \mathbf{X} = \mathbf{x}) \\ &= \arg \max_{y \in \{0, 1, \dots, c-1\}} p^*(y|\mathbf{x}) \end{aligned}$$

となる。上式の 3 番目の等号は $1 - \Pr(Y = y | \mathbf{X} = \mathbf{x})$ を最小にする y の値は $\Pr(Y = y | \mathbf{X} = \mathbf{x})$ を最大にする y の値に一致することからわかる。回帰問題で 2 乗損失を用いたとき最適の予測関数 g^* を回帰関数とよぶ。

定理 1.1. 回帰問題において 2 乗損失 $\text{Loss}(y, \hat{y}) = (y - \hat{y})^2$ を用いる. $E[Y^2] < \infty$ を仮定する. このとき最適予測関数 g^* は次で与えられる.

$$g^*(\mathbf{x}) = E[Y | \mathbf{X} = \mathbf{x}].$$

すなわち $\mathbf{X} = \mathbf{x}$ を与えたときの Y の条件付き期待値である.

Proof. $g^*(\mathbf{x}) = E[Y | \mathbf{X} = \mathbf{x}]$ とする. 任意の (2 乗可積分な) 関数 g を考える. すると

$$\begin{aligned} E[\{Y - g(\mathbf{X})\}^2] &= E[\{Y - g^*(\mathbf{X}) + g^*(\mathbf{X}) - g(\mathbf{X})\}^2] \\ &= E[\{Y - g^*(\mathbf{X})\}^2] + \underbrace{E[\{g^*(\mathbf{X}) - g(\mathbf{X})\}^2]}_{\geq 0} \\ &\quad + 2E[\{Y - g^*(\mathbf{X})\}\{g^*(\mathbf{X}) - g(\mathbf{X})\}] \\ &\geq E[\{Y - g^*(\mathbf{X})\}^2] + 2E[\{Y - g^*(\mathbf{X})\}\{g^*(\mathbf{X}) - g(\mathbf{X})\}] \end{aligned}$$

となる. 上の不等式の等号は $g(\mathbf{x}) = g^*(\mathbf{x})$ のとき成立する. しかし

$$\begin{aligned} E[\{Y - g^*(\mathbf{X})\}\{g^*(\mathbf{X}) - g(\mathbf{X})\}] &= E[E[\{Y - g^*(\mathbf{X})\}\{g^*(\mathbf{X}) - g(\mathbf{X})\} | \mathbf{X}]] \\ &= E[\{g^*(\mathbf{X}) - g(\mathbf{X})\}E[Y - g^*(\mathbf{X}) | \mathbf{X}]] \\ &= E[\{g^*(\mathbf{X}) - g(\mathbf{X})\}\underbrace{\{E[Y | \mathbf{X}] - g^*(\mathbf{X})\}}_{=0}] \\ &= 0 \end{aligned}$$

となる. よって

$$E[\{(Y - g(\mathbf{X}))\}^2] \geq E[\{Y - g^*(\mathbf{X})\}^2]$$

がわかる. □

定理 1.1 により $\mathbf{X} = \mathbf{x}$ を与えたときの出力 Y を次のように表現できる.

$$Y = g^*(\mathbf{x}) + \epsilon(\mathbf{x}).$$

すると

$$E[\epsilon(\mathbf{x}) | \mathbf{X} = \mathbf{x}] = E[Y - g^*(\mathbf{x}) | \mathbf{X} = \mathbf{x}] = E[Y | \mathbf{X} = \mathbf{x}] - g^*(\mathbf{x}) = 0.$$

さらに

$$\text{Var}[\epsilon(\mathbf{X}) | \mathbf{X} = \mathbf{x}] =: \nu^2(\mathbf{x})$$

と書ける. $\nu(\mathbf{x})$ は未知であり $\epsilon(\mathbf{x})$ の分布は平均が 0 であること以外は未知である.

一般に、最適予測関数 g^* は (X, Y) の同時分布に依存し、その分布は未知である。したがって g^* を実際の問題では用いることができない。 (X, Y) の同時分布 $p_{(X,Y)}(x, y)$ に従う n 個の確率変数列を

$$\mathfrak{T}_n = \{(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)\}$$

と書き \mathfrak{T}_n を訓練集合という。この訓練集合の実現値を

$$t_n = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

と書く。

目標は訓練集合 \mathfrak{T}_n が与えられたとき、未知の最適予測関数 g^* を推定/学習することである。訓練集合 \mathfrak{T} に基づく g^* の近似/推定量を $g_{\mathfrak{T}}$ と書くこと⁶にする。ここで $g_{\mathfrak{T}}$ はランダムな関数であることに注意せよ。実現値 t_n により計算されたものが実現値になり、 g_t と書く。訓練集合 \mathfrak{T}_n により未知の関数

$$g^* : \mathbf{x} \mapsto \hat{y}$$

を学習する学習器が $g_{\mathfrak{T}}$ である。教師が真の関係に基づき n 個の出力と入力 t_n の対 \mathfrak{T}_n を見本として与え、その見本で学習器 $g_{\mathfrak{T}}$ を訓練する。新しい入力 X に対して教師によって与えられていない出力 Y を $g_{\mathfrak{T}}(X)$ で予測することになる。この設定を教師あり学習といい、入力のことを説明変数、出力のことを応答変数ということもある。

一方、教師なし学習では変数に説明変数と応答変数の区別はない。データの同時分布 $p^*(x, y)$ を学習することになる。この変数を (X, Y) と書くことにする。 (X, Y) の分布 p^* の推定量 p に対して誤差は

$$\text{Err}(p) = E[\text{Loss}(p^*(X, Y), p(X, Y))]$$

となる。教師なし学習の例としては、コンビニエンスストアの客の購買行動の解析を考える。100 個の商品があるとする。ある客が、 i ($i = 1, 2, \dots, 100$) 番目の商品を買うか買わないの結果を $0, 1$ に対応させる。客の購買パターンは $\mathbf{x} \in \{0, 1\}^{100}$ となる。訓練集合の実現値 $t_{100} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{100}\}$ に基づいて、客の購買パターンを見つけたい。教師なし学習の手法としては、クラスタリング、主成分分析、カーネル密度推定、決定木等がある。

以下の節では、教師あり学習についてみていくことにする。

⁶本来ならば、 $g_{\mathfrak{T}_n}$ と書くべきだろうが、煩雑になるので、簡単に $g_{\mathfrak{T}}$ と書くことにした。

1.3 教師あり学習における訓練誤差と汎化誤差

任意の予測関数 (以後は学習器とよぶことにする.) g が与えられたとき予測誤差を

$$\text{Err}(g) := E[\text{Loss}(Y, g(\mathbf{X}))]$$

で定義する. これを求めることは解析者は一般にはできない. なぜならば統計家にとって (\mathbf{X}, Y) の同時分布は未知だからである. 学習器 g の予測誤差 $\text{Err}(g)$ を近似/推定するために訓練集合 \mathfrak{T}_n に基づく訓練誤差

$$\widehat{\text{Err}} := \frac{1}{n} \sum_{j=1}^n \text{Loss}(Y_j, g(\mathbf{X}_j))$$

を用いる. 訓練集合 \mathfrak{T}_n は (\mathbf{X}, Y) の同時分布に独立同一に従う確率変数列なので

$$\text{Err}(g) = E[\widehat{\text{Err}}(g)]$$

が成立していることが簡単にわかる. すなわち $\widehat{\text{Err}}(g)$ は $\text{Err}(g)$ の不偏推定量である.

最適予測関数を g^* を近似/学習/推定するために, 学習器の候補の集まりである関数族を \mathcal{G} を選択する. すなわち, 訓練誤差を最小にする学習器を見つけることを目指す. 新しいデータ \mathbf{X} の学習器の予測精度を測る尺度を導入する. 訓練集合の実現値 \mathfrak{t}_n を固定する. すなわち, $(\mathbf{X}_j, Y_j) = (\mathbf{x}_j, y_j) (j = 1, 2, \dots, n)$ とする.

訓練集合の実現値 \mathfrak{t}_n の基づく学習器 $g_{\mathfrak{t}}^{\mathcal{G}}$ を

$$g_{\mathfrak{t}}^{\mathcal{G}} \in \arg \min_{g \in \mathcal{G}} \frac{1}{n} \sum_{j=1}^n \text{Loss}(y_j, g(\mathbf{x}_j))$$

で定義する. (\mathbf{X}, Y) を訓練集合 \mathfrak{T}_n とは独立に同時分布 $p_{(\mathbf{X}, Y)}(\mathbf{x}, y)$ に従う確率ベクトルとしたとき学習器 $g_{\mathfrak{t}}^{\mathcal{G}}$ の汎化誤差 $\text{Err}(g_{\mathfrak{t}}^{\mathcal{G}})$ を

$$\text{Err}(g_{\mathfrak{t}}^{\mathcal{G}}) := E[\text{Loss}(Y, g_{\mathfrak{t}}^{\mathcal{G}}(\mathbf{X}))]$$

で定義する. ただし期待値 $E[\cdot]$ は (\mathbf{X}, Y) の同時分布 $p_{(\mathbf{X}, Y)}(\mathbf{x}, y)$ に関して取ったものである.

ランダムな訓練集合 \mathfrak{T}_n に基づく学習器 $g_{\mathfrak{x}}^{\mathcal{G}}$ の汎化誤差 $\text{Err}(g_{\mathfrak{x}}^{\mathcal{G}})$ はランダムになるので, \mathfrak{T}_n の同時分布に関して期待値をとったものを期待汎化誤差という.

$$E[\text{Err}(g_{\mathfrak{x}}^{\mathcal{G}})] = E[\text{Loss}(Y, g_{\mathfrak{x}}^{\mathcal{G}}(\mathbf{X}))].$$

右辺の $E[\cdot]$ は \mathfrak{T}_n の同時分布に関する期待値であり, 左辺の $E[\cdot]$ は (\mathbf{X}, Y) と \mathfrak{T}_n の同時分布に関する期待値である.

注意 1.2. (多項式回帰) (U, Y) の同時分布を以下のように定める. $U \sim \text{Unif}(0, 1)$ とし $U = u$ ($0 < u < 1$) が与えられたときの Y の条件付き分布は $N(\theta, \sigma^2)$ とする. ただし,

$$\theta = 10 - 140u + 400u^2 - 250u^3, \quad \sigma^2 = 25$$

である. 2 乗損失を用いたとき最良予測関数は

$$g^*(u) = E[Y|U = u] = 10 - 140u + 400u^2 - 250u^3$$

となる. これは以下の議論からわかる.

$$\begin{aligned} \text{Err}(g) &= E[(Y - g(U))^2] = E[\{Y - E[Y|U] + E[Y|U] + g(U)\}^2] \\ &= E[\{E[Y|U] + g(U)\}^2] + E[\{E[Y|U] + g(U)\}^2] \\ &\quad + 2E[\{Y - E[Y|U]\}\{E[Y|U] + g(U)\}] \\ &= E[\{E[Y|U] + g(U)\}^2] + E[\{E[Y|U] + g(U)\}^2] \\ &\quad + 2E\left[E\left[\{Y - E[Y|U]\}\{E[Y|U] + g(U)\} \middle| U\right]\right] \\ &= E[\{E[Y|U] + g(U)\}^2] + E[\{E[Y|U] + g(U)\}^2] \\ &\quad + 2E\left[\{E[Y|U] + g(U)\} \underbrace{E\left[Y - E[Y|U] \middle| U\right]}_{=0}\right] \\ &= E[\{Y - E[Y|U]\}^2] + E[\{E[Y|U] - g(U)\}^2] \\ &\geq E[\{E[Y|U] - g(U)\}^2]. \end{aligned}$$

$d \in \mathbb{N}$, $d \geq 2$ とする. $d - 1$ 次元の u の多項式がなす関数族を考える.

$$\mathcal{G}_d := \{g(u) = \beta_0 + \beta_1 u + \beta_2 u^2 + \cdots + \beta_{d-1} u^{d-1}; \beta_0, \beta_1, \dots, \beta_{d-1} \in \mathbb{R}\}.$$

$$\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_{d-1})^\top, \quad \boldsymbol{x} = (1, u, u^2, \dots, u^{d-1})^\top$$

とおくと

$$\mathcal{G}_d = \{g(u) = \boldsymbol{x}^\top \boldsymbol{\beta}; \boldsymbol{\beta} \in \mathbb{R}^d\}$$

と書ける. $d \geq 4$ とする.

$$\boldsymbol{\beta}^* = (10, -140, 400, -250, 0, \dots, 0)^\top$$

としたとき $g^*(u) = \boldsymbol{x}^\top \boldsymbol{\beta}^*$ が最良予測関数となる. 訓練データを

$$(u_1, y_1), (u_2, y_2), \dots, (u_n, y_n)$$

とし

$$\boldsymbol{x}_j = (1, u_j, u_j^2, \dots, u_j^{d-1})^\top$$

と書く. これらを行列とベクトルで表記する.

$$\underbrace{\mathbf{X}}_{n \times d \text{ 行列}} = \begin{pmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_n^\top \end{pmatrix} = \begin{pmatrix} 1 & u_1 & u_1^2 & \cdots & u_1^{d-1} \\ 1 & u_2 & u_2^2 & \cdots & u_2^{d-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & u_n & u_n^2 & \cdots & u_n^{d-1} \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

とおく. すると学習器 $g(u) = \mathbf{x}^\top \boldsymbol{\beta}$ の訓練誤差は

$$\widehat{\text{Err}}(g) = \frac{1}{n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_{2,n}^2$$

となる. ただし $\|\cdot\|_{2,n}$ は \mathbb{R}^n の Euclid ノルムである. すると

$$\widehat{\boldsymbol{\beta}} \in \arg \min_{\boldsymbol{\beta} \in \mathcal{G}_d} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_{2,n}^2$$

とすれば $\widehat{\boldsymbol{\beta}}$ は最小 2 乗推定量となる. \mathbf{X} の列で張られた線型部分空間を $\text{span}(\mathbf{X})$ と書く. $\text{span}(\mathbf{X})$ への射影行列を \mathbf{P} とすれば $\widehat{\boldsymbol{\beta}}$ は

$$\mathbf{X}\widehat{\boldsymbol{\beta}} = \mathbf{P}\mathbf{y}$$

をみたすことがわかる. すると $d \times n$ の行列 \mathbf{X}^\dagger が存在⁷して

$$\mathbf{P} = \mathbf{X}\mathbf{X}^\dagger$$

と書けることが知られている. したがって

$$\widehat{\boldsymbol{\beta}} = \mathbf{X}^\dagger \mathbf{y}$$

と書けることがわかる. よって

$$g_t^{\mathcal{G}_d}(u) = \mathbf{x}^\top \widehat{\boldsymbol{\beta}}, \quad \widehat{\boldsymbol{\beta}} = \mathbf{X}^\dagger \mathbf{y}$$

となる. 特に $\mathbf{X}^\top \mathbf{X}$ が正則のとき

$$\widehat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

となる. □

⁷ \mathbf{X}^\dagger は \mathbf{X} の Moore-Penrose の一般逆行列と呼ばれるものである. 補遺の B 章を参照のこと.

1.4 教師あり学習における兼ね合い (trade-off)

教師ありの機械学習において汎化誤差

$$\text{Err}(g_t^{\mathcal{G}}) = \mathbb{E}[\text{Loss}(Y, g_t^{\mathcal{G}}(\mathbf{X}))]$$

もしくは

$$\mathbb{E}[\text{Err}(g_t^{\mathcal{G}})] = \mathbb{E}[\mathbb{E}[\text{Loss}(Y, g_t^{\mathcal{G}}(\mathbf{X}))]]$$

の最小化問題を考える. この問題を解くために関数族 \mathcal{G} を適切に選択することが肝要である. 選択の最には以下の点を考慮する必要がある.

- 関数族の複雑さ (最適予測関数を近似するために十分な豊かさを持つことが望ましい. 最適予測関数を含む族であればなおさらよい).
- 最適問題

$$g_{\mathcal{X}}^{\mathcal{G}} \in \arg \min_{g \in \mathcal{G}} \frac{1}{n} \sum_{j=1}^n \text{Loss}(Y_j, g(\mathbf{X}_j))$$

を解く学習器の訓練が容易である.

- 訓練誤差

$$\widehat{\text{Err}}(g) = \frac{1}{n} \sum_{j=1}^n \text{Loss}(Y_j, g(\mathbf{X}_j))$$

が予測誤差

$$\text{Err}(g) = \mathbb{E}[\text{Loss}(Y, g(\mathbf{X}))]$$

の推定量としてよい精度をもっているか?

- 特徴量のタイプは何か?

適切な関数族 \mathcal{G} の選択には衝突する因子間の兼ね合いが伴う. 簡単な関数族からの学習器の訓練は容易であるが, 最良予測関数 g^* をうまく近似/推定していないかもしれない. 逆に g^* を含むかもしれない豊かな関数族 \mathcal{G} からの学習器の訓練には計算コストがかかる. 関数族の複雑さ, 計算コスト, 推定精度の関係を理解するために汎化誤差を複数の因子に分解する. 分解には二通りある. ひとつは近似と推定精度の兼ね合いを説明するもので, 他方はバイアスと分散の兼ね合いを説明する分解である.

まず近似と推定誤差の兼ね合いを理解するための分解を考える.

$$\text{Err}(g_t^{\mathcal{G}}) = \text{Err}^* + \underbrace{\text{Err}(g^{\mathcal{G}}) - \text{Err}^*}_{\text{近似の精度}} + \underbrace{\text{Err}(g_t^{\mathcal{G}}) - \text{Err}(g^{\mathcal{G}})}_{\text{推定の精度}}. \quad (1.1)$$

ただし

$$\text{Err}^* = \text{Err}(g^*), \quad g^* \in \arg \min_g \text{Err}(g), \quad g^{\mathcal{G}} \in \arg \min_{g \in \mathcal{G}} \text{Err}(g).$$

(1.1) の右辺の第 2 項目を近似の精度といい、誤差の下限と関数族 \mathcal{G} 中の最適な予測関数の誤差との差である。関数族 \mathcal{G} の選択と \mathcal{G} 上の $\text{Err}(g)$ の最小化問題は関数解析と数値解析の問題である。訓練データ \mathbf{t} は関わらない。近似精度をよくするためには、関数族 \mathcal{G} を大きくとればよい。第 3 項は統計的誤差になる。これは訓練データにかかわるものである。 $g_t^{\mathcal{G}}$ がよりうまく $g^{\mathcal{G}}$ を推定するかどうかである。標本数 n が無限大に近づくとこの精度は 0 に収束する。近似と推定の精度をよくするためには互いに相反することが要求される。推定精度をよくするためには \mathcal{G} を小さくすればよい。一方、近似精度をよくするためには関数族 \mathcal{G} を大きく取ればよい。したがって兼ね合いが肝要となる。

2 乗損失を考える。この場合には汎化誤差は

$$\text{Err}(g_t^{\mathcal{G}}) = \mathbb{E}[\{Y - g_t^{\mathcal{G}}(\mathbf{X})\}^2].$$

最良予測関数は

$$g^*(\mathbf{X}) = \mathbb{E}[Y | \mathbf{X} = \mathbf{x}]$$

であることを思い出そう。すると

- (1) 誤差の下限: $\text{Err}^* = \mathbb{E}[(Y - g^*(\mathbf{X}))^2]$.
- (2) 第 2 因子の近似精度: $\text{Err}(g^{\mathcal{G}}) - \text{Err}(g^*) = \mathbb{E}[(g^{\mathcal{G}}(\mathbf{X}) - g^*(\mathbf{X}))^2]$.
- (3) 第 3 因子の推定精度: $\mathcal{G} = \{g(\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\beta}; \boldsymbol{\beta} \in \mathbb{R}^d\}$ としたとき、

$$\text{Err}(g_t^{\mathcal{G}}) - \text{Err}(g^{\mathcal{G}}) = \mathbb{E}[(g_t^{\mathcal{G}}(\mathbf{X}) - g^{\mathcal{G}}(\mathbf{X}))^2].$$

よって線型関数族 \mathcal{G} に対して

$$\begin{aligned} \text{Err}(g_t^{\mathcal{G}}) &= \mathbb{E}[(g_t^{\mathcal{G}}(\mathbf{X}) - Y)^2] \\ &= \text{Err}^* + \mathbb{E}[(g^{\mathcal{G}}(\mathbf{X}) - g^*(\mathbf{X}))^2] + \mathbb{E}[(g_t^{\mathcal{G}}(\mathbf{X}) - g^{\mathcal{G}}(\mathbf{X}))^2] \end{aligned}$$

となる。

次に 2 乗損失の場合についてバイアスと分散の兼ね合いについて説明する。まず

$$\begin{aligned} \text{Err}(g_t^{\mathcal{G}}) &= \mathbb{E}[(g_t^{\mathcal{G}}(\mathbf{X}) - Y)^2] \\ &= \mathbb{E}[(g_t^{\mathcal{G}}(\mathbf{X}) - g^*(\mathbf{X}) + g^*(\mathbf{X}) - Y)^2] \\ &= \mathbb{E}[(Y - g^*(\mathbf{X}))^2] + \mathbb{E}[(g_t^{\mathcal{G}}(\mathbf{X}) - g^*(\mathbf{X}))^2] \\ &= \text{Err}^* + \mathbb{E}[D^2(\mathbf{X}, \mathbf{t})]. \end{aligned}$$

ただし

$$D(\mathbf{X}, \mathbf{t}) = g_t^{\mathcal{G}}(\mathbf{X}) - g^*(\mathbf{X})$$

の \mathfrak{t} に \mathfrak{T} を代入すると

$$\begin{aligned} E[(g_{\mathfrak{T}}^G(\mathbf{x}) - g^*(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}] &= E[D^2(\mathbf{x}, \mathfrak{T}) | \mathbf{X} = \mathbf{x}] \\ &= \{E[D(\mathbf{x}, \mathfrak{T}) | \mathbf{X} = \mathbf{x}]\}^2 + \text{Var}[D(\mathbf{x}, \mathfrak{T}) | \mathbf{X} = \mathbf{x}] \\ &\quad (\because \text{注意 A.47}) \\ &= \{E[g_{\mathfrak{T}}^G(\mathbf{x}) - g^*(\mathbf{x}) | \mathbf{X} = \mathbf{x}]\}^2 + \text{Var}[g_{\mathfrak{T}}^G | \mathbf{X} = \mathbf{x}]. \end{aligned}$$

よって

$$E[\text{Err}(g_{\mathfrak{T}}^G)] = \text{Err}^* + \underbrace{\{E[g_{\mathfrak{T}}^G(\mathbf{X}) | \mathbf{X}] - g^*(\mathbf{X})\}^2}_{\text{バイアス}} + \underbrace{E[\text{Var}[g_{\mathfrak{T}}^G(\mathbf{X}) | \mathbf{X}]]}_{\text{分散}}$$

と書ける.

1.5 教師なし学習における兼ね合い (trade-off)

応答変数 Y と説明変数ベクトル X の区別がある教師あり学習とは対照的に教師なし学習ではこのような変数の区別がない.

$n \in \mathbb{N}$ を標本 (観測データ) 数とする. 各標本 (観測データ) の次元を $k \in \mathbb{N}$ とし, 縦ベクトルで表すことにする. したがって, n 個の標本を \mathbf{x}_j ($j = 1, 2, \dots, n$) と書いたとき, $k \times n$ の観測データ行列 $[\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$ から有用情報やパターンを抽出することが教師なし学習の主目的になる.

標本 (観測データ) $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ はある未知の確率分布からの n の独立同一標本の実現値と考え, データ行列からこの未知の確率分布を学習/推測することが教師なし学習は本質的に目標とする.

観測データはある未知の確率分布からの n 個のランダム標本の実現値であるときデータの経験分布は未知の確率分布に関する重要な情報を含んでいる. 経験分布の概念とカーネル型密度推定については後ほど説明する.

教師なし学習における損失と誤差を導入する. 訓練データを

$$\mathfrak{T}_n := \{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n\}$$

とする. ここで \mathbf{X}_j ($j = 1, 2, \dots, n$) は確率ベクトル $\mathbf{X} \in \mathbb{R}^k$ の独立複製とする. また \mathbf{X} の未知の同時 p.d.f. を p^* と書くことにする. この p を真の分布ということにする. さらに, 訓練データ \mathfrak{T}_n の実現値 $\mathfrak{t}_n = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ と書くことにする.

教師なし学習では未知の同時 p.d.f. (真の分布) p^* をうまく近似する p.d.f. p をみつけることが目標である. すなわち, $\text{Loss}(p^*, p)$ を損失関

数としたときその汎化誤差

$$\begin{aligned} \text{Err}(p) &:= E[\text{Loss}(p^*(\mathbf{X}), p(\mathbf{X}))] \\ &= \int \text{Loss}(p^*(\mathbf{x}), p(\mathbf{x})) p^*(\mathbf{x}) d\mathbf{x} \end{aligned}$$

を小さくする p.d.f. p をみつけることである. ただし $E[\cdot]$ は未知の真の確率分布 p^* に関する期待値である. 損失関数を

$$\text{Loss}(p^*, p) = \log \frac{p^*}{p} = \begin{cases} \log \frac{p^*(\mathbf{x})}{p(\mathbf{x})} & (p^*(\mathbf{x})p(\mathbf{x}) > 0), \\ 0 & (p^*(\mathbf{x})p(\mathbf{x}) = 0) \end{cases}$$

と取れば, 汎化誤差は Kullback-Leibler 情報量 (K-L 情報量とかくことにする) と呼ばれるものになる.

以下では K-L 情報量のもとで考える. K-L 情報量を KL と書くことにする. すなわち, K-L 情報量 KL を

$$\text{KL}(p) = E\left[\log \frac{p^*(X)}{p(X)}\right] = E[\log p^*(X)] - E[\log p(X)]$$

で定める. 上の式の最右辺の第 1 項目は p に関係ないので $\text{KL}(p)$ の最小化は

$$-E[\log p(X)] = -\int p^*(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x}$$

の最小化となる. ここで

$$\text{CE}(p) := -E[\log p(\mathbf{X})] \quad (1.2)$$

と書くことにする. これをクロス・エントロピー誤差と呼ぶことにする.

次に探索の対象となる確率密度関数の集まりを \mathcal{P} となくことにする. この確率分布の集まりを統計的モデルという. 統計的モデル \mathcal{P} が真の分布 p を含んでいる ($p^* \in \mathcal{P}$ が成立する) とき, クロス・エントロピー誤差 $\text{CE}(p)$ を最小とするものは p^* となる. すなわち, 任意の $p \in \mathcal{P}$ に対して

$$\text{CE}(p) \geq \text{CE}(p^*)$$

が成立する. しかし (1.2) の最小化問題は一般に実行不可能である. なぜならばこの誤差は未知の分布 p^* に依存するからである. $\text{CE}(p)$ の代わりに訓練クロス・エントロピー誤差

$$\widehat{\text{CE}}_t := -\frac{1}{n} \sum_{j=1}^n \log p(\mathbf{x}_j)$$

の $p \in \mathcal{P}$ に関する最小化問題を考える.

ここで重要なステップは探索の対象となる統計的モデル \mathcal{P} の選択である. まず \mathcal{P} の元 p を添え字 θ で添え字付ることにする. 添え字 θ の全体の成す集合 (この集合を添え字集合) を Θ と書くことにする. すなわち集合 \mathcal{P} を集合 Θ により母数化する. 添え字の θ を母数といい, Θ を母数空間と呼ぶことにする. すると \mathcal{P} の選択は Θ の選択となる.

いま, $d \in \mathbb{N}$ とし, $\Theta \subset \mathbb{R}^d$ とする. すなわち

$$\mathcal{P} := \{p(\cdot | \theta); \theta \in \Theta \subset \mathbb{R}^d\}. \quad (1.3)$$

(1.3) のような形で与えられる統計的モデルを母数モデルと統計学では伝統的に呼んでいる. すなわち, 母数モデルとは有限次元の母数空間によって添え字づけられている分布の集まり (統計的モデル) である.

以降では, (1.3) の中から未知の真の分布 p^* をよく近似する確率分布を学習することにする.

まず, データから近似モデルを探するときの情報な概念を導入する. データの実現値 $X = x$ が与えられたときに関数

$$\theta \mapsto p(x | \theta) \quad (1.4)$$

を尤度関数という. この関数の $\theta = (\theta_1, \theta_2, \dots, \theta_d)^\top$ に関する勾配をスコア関数といい, $S(x | \theta)$ と書くことにする.

$$S(x | \theta) := \frac{1}{p(x | \theta)} \frac{\partial p(x | \theta)}{\partial \theta} := \frac{1}{p(x | \theta)} \left(\frac{\partial p(x | \theta)}{\partial \theta_1}, \frac{\partial p(x | \theta)}{\partial \theta_2}, \dots, \frac{\partial p(x | \theta)}{\partial \theta_d} \right)^\top.$$

$S(x | \theta)$ の x に X を代入しランダムにしたものを $S(X | \theta)$ と書くことにする. すると

$$\begin{aligned} E_\theta [S(X | \theta)] &= \int \frac{1}{p(x | \theta)} \left(\frac{\partial p(x | \theta)}{\partial \theta} \right) p(x | \theta) dx = \int \frac{\partial p(x | \theta)}{\partial \theta} dx \\ &= \frac{\partial}{\partial \theta} \int p(x | \theta) dx = 0 \end{aligned} \quad (1.5)$$

となる. ただし最後の等号は, 微分記号と積分記号の交換が可能であることを仮定⁸している. すなわち $p(\cdot | \theta)$ の関する $S(X | \theta)$ の期待値は零ベクトルとなる. さらに $p(\cdot | \theta)$ の関する $S(X | \theta)$ の共分散行列を考える.

⁸無条件でこの交換が可能なのわけではない. 通常は, 交換が可能になるような正則条件を母数モデル \mathcal{P} や真の分布 p^* に課す. ここの議論では, 厳密性に欠くが, 正則条件を気にせずに議論をする.

これを Fisher 情報量行列といい, $\mathbf{I}(\theta)$ と記すことにする.

$$\begin{aligned} \mathbf{I}(\theta) &:= E_{\theta}[\mathbf{S}(\mathbf{X}|\theta)\mathbf{S}^{\top}(\mathbf{X}|\theta)] \\ &= E_{\theta} \left[\begin{array}{cccc} \frac{\partial \log p(\mathbf{X}|\theta)}{\partial \theta_1} \frac{\partial \log p(\mathbf{X}|\theta)}{\partial \theta_1} & \frac{\partial \log p(\mathbf{X}|\theta)}{\partial \theta_1} \frac{\partial \log p(\mathbf{X}|\theta)}{\partial \theta_2} & \cdots & \frac{\partial \log p(\mathbf{X}|\theta)}{\partial \theta_1} \frac{\partial \log p(\mathbf{X}|\theta)}{\partial \theta_d} \\ \frac{\partial \log p(\mathbf{X}|\theta)}{\partial \theta_2} \frac{\partial \log p(\mathbf{X}|\theta)}{\partial \theta_1} & \frac{\partial \log p(\mathbf{X}|\theta)}{\partial \theta_2} \frac{\partial \log p(\mathbf{X}|\theta)}{\partial \theta_2} & \cdots & \frac{\partial \log p(\mathbf{X}|\theta)}{\partial \theta_2} \frac{\partial \log p(\mathbf{X}|\theta)}{\partial \theta_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \log p(\mathbf{X}|\theta)}{\partial \theta_d} \frac{\partial \log p(\mathbf{X}|\theta)}{\partial \theta_1} & \frac{\partial \log p(\mathbf{X}|\theta)}{\partial \theta_d} \frac{\partial \log p(\mathbf{X}|\theta)}{\partial \theta_2} & \cdots & \frac{\partial \log p(\mathbf{X}|\theta)}{\partial \theta_d} \frac{\partial \log p(\mathbf{X}|\theta)}{\partial \theta_d} \end{array} \right] \end{aligned}$$

で定義する. さらに $-\log p(\mathbf{X}|\theta)$ の Hesse 行列の分布 p に関する期待値を $\mathbf{J}(\theta)$ を

$$\begin{aligned} \mathbf{J}(\theta) &:= -E_{\theta} \left[\frac{\partial \mathbf{S}(\mathbf{X}|\theta)}{\partial \theta} \right] \\ &:= -E_{\theta} \left[\begin{array}{cccc} \frac{\partial^2 \log p(\mathbf{X}|\theta)}{\partial \theta_1 \partial \theta_1} & \frac{\partial^2 \log p(\mathbf{X}|\theta)}{\partial \theta_1 \partial \theta_2} & \cdots & \frac{\partial^2 \log p(\mathbf{X}|\theta)}{\partial \theta_1 \partial \theta_d} \\ \frac{\partial^2 \log p(\mathbf{X}|\theta)}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2 \log p(\mathbf{X}|\theta)}{\partial \theta_2 \partial \theta_2} & \cdots & \frac{\partial^2 \log p(\mathbf{X}|\theta)}{\partial \theta_2 \partial \theta_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 \log p(\mathbf{X}|\theta)}{\partial \theta_d \partial \theta_1} & \frac{\partial^2 \log p(\mathbf{X}|\theta)}{\partial \theta_d \partial \theta_2} & \cdots & \frac{\partial^2 \log p(\mathbf{X}|\theta)}{\partial \theta_d \partial \theta_d} \end{array} \right] \end{aligned}$$

で定義する. ある $\theta^* \in \Theta$ が存在して $p^* = p(\cdot|\theta^*)$ となるとき

$$\mathbf{J}(\theta^*) = \mathbf{I}(\theta^*)$$

となること⁹がわかる.

⁹ $d = 1$ として確認する. $d \geq 2$ の場合も本質的には同じである.

$$\begin{aligned} J(\theta) &= -E \left[\frac{\partial^2 \log \mathbf{p}(X|\theta)}{\partial \theta^2} \right] = -E \left[\frac{\partial}{\partial \theta} \left(\frac{1}{\mathbf{p}(X|\theta)} \frac{\partial \mathbf{p}(X|\theta)}{\partial \theta} \right) \right] \\ &= E \left[\frac{1}{\mathbf{p}^2(X|\theta)} \left(\frac{\partial \mathbf{p}(X|\theta)}{\partial \theta} \right)^2 - \frac{1}{\mathbf{p}(X|\theta)} \frac{\partial^2 \mathbf{p}(X|\theta)}{\partial \theta^2} \right] \\ &= E \left[\left(\frac{\partial \log \mathbf{p}(X|\theta)}{\partial \theta} \right)^2 \right] - E \left[\frac{1}{\mathbf{p}(X|\theta)} \frac{\partial^2 \mathbf{p}(X|\theta)}{\partial \theta^2} \right] \\ &= I(\theta) - E \left[\frac{1}{\mathbf{p}(X|\theta)} \frac{\partial^2 \mathbf{p}(X|\theta)}{\partial \theta^2} \right]. \end{aligned}$$

しかし

$$\begin{aligned} E \left[\frac{1}{\mathbf{p}(X|\theta)} \frac{\partial^2 \mathbf{p}(X|\theta)}{\partial \theta^2} \right] &= \int \frac{1}{\mathbf{p}(x|\theta)} \frac{\partial^2 \mathbf{p}(x|\theta)}{\partial \theta^2} \mathbf{p}(x|\theta) dx \\ &= \int \frac{\partial^2 g(x|\theta)}{\partial \theta^2} dx = \frac{1}{\partial \theta^2} \int g(x|\theta) dx = 0. \end{aligned}$$

n が十分大きいとき行列 $I(\theta)$, $J(\theta)$ はクロス・エントロピー誤差の近似において重要な役割を担う. 設定を記述するために, $p^{\mathcal{P}} = p(\cdot | \theta^*)$ をクロス・エントロピー誤差

$$CE(\theta) := -E[p(X | \theta)]$$

を最小にする点とする. すなわち任意の $\theta \in \Theta$ に対して

$$CE(\theta) \geq CE(\theta^*)$$

である. CE は θ の関数として「なめらか」な関数と仮定する. 特に, これらの行列が θ^* の近傍で狭義凸で 2 回連続微分可能¹⁰なときこの仮定をみたしていることがわかる. θ^* の定義より

$$0 = \frac{\partial}{\partial \theta} CE(\theta) \Big|_{\theta=\theta^*} = -\frac{\partial}{\partial \theta} E[\log p(X | \theta)] \Big|_{\theta=\theta^*} = -E \left[\frac{\partial}{\partial \theta} \log p(X | \theta) \right] \Big|_{\theta=\theta^*}$$

となる. ただし微分記号と積分記号の交換が可能であることを仮定している. 同じように $J(\theta)$ は Err の Hesse 行列であることがわかる. $p(\cdot | \hat{\theta}_n)$ を訓練誤差

$$CE_{\mathfrak{X}_n}(\theta) := -\frac{1}{n} \sum_{j=1}^n \log p(X_j | \theta)$$

を最小にする点とする. ただし $\mathfrak{X}_n = \{X_1, X_2, \dots, X_n\}$ はランダム標本の集合である. CE^* をすべての確率分布に関するクロス・エントロピーの最小値とする. 明らかに $X \sim p^*$ のとき最小となる. すなわち $CE^* = -E[\log p^*(X)]$ である. 教師あり学習の場合と同様に汎化誤差 $CE(\hat{\theta}_n)$ を分解する.

$$\begin{aligned} CE(\hat{\theta}_n) &= CE^* + \underbrace{CE(\theta^*) - CE^*}_{\text{近似精度}} + \underbrace{CE(\hat{\theta}_n) - CE(\theta^*)}_{\text{推定精度}} \\ &= CE(\theta^*) + E \left[\frac{p(X | \theta^*)}{p(X | \hat{\theta}_n)} \right]. \end{aligned}$$

上の式の最左辺と最右辺の第 2 項目は \mathfrak{X}_n に依存しているのでランダムな量であることに注意せよ.

1.6 最尤法

$d \in \mathbb{N}$ とする. 母数モデル $\mathcal{P} = \{p(x, \theta) : \theta \in \Theta \subset \mathbb{R}^d\}$ を考える. ただし $p(x | \theta)$ は p.d.f. もしくは p.m.f. とする. 真の母数 θ^* は母数空間 Θ に含まれるとする. すなわち, $\theta^* \in \Theta$ である.

¹⁰たとえば g が正規分布のときにはこの条件をみたしている.

いま, $n \in \mathbb{N}$ とし,

$$\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n \sim \text{i.i.d. } p(\mathbf{x} | \boldsymbol{\theta}^*)$$

とする.

定義 1.3. $\mathbf{X}_1 = \mathbf{x}_1, \mathbf{X}_2 = \mathbf{x}_2, \dots, \mathbf{X}_n = \mathbf{x}_n$ を観測したとき $\boldsymbol{\theta}$ の尤度関数 $\text{lik}_n(\boldsymbol{\theta})$ を

$$\text{lik}_n(\boldsymbol{\theta}) = \prod_{j=1}^n p(\mathbf{x}_j | \boldsymbol{\theta})$$

で定義し, 対数尤度 $\ell_n(\boldsymbol{\theta})$ を

$$\ell_n(\boldsymbol{\theta}) = \log \text{lik}_n(\boldsymbol{\theta})$$

で定義する.

注意 1.4. 尤度関数は

$$\text{lik}_n : \Theta \ni \boldsymbol{\theta} \mapsto \text{lik}_n(\boldsymbol{\theta}) = \prod_{i=1}^n p(\mathbf{x}_i | \boldsymbol{\theta}) \in [0, \infty)$$

であることに注意をせよ. 一方, $\mathbb{X} = \{\mathbf{x} \in \mathbb{R}^k; p(\mathbf{x} | \boldsymbol{\theta}^*) > 0\}$ とおいたとき, $(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n)$ の同時 p.d.f. または p.m.f. $\prod_{i=1}^n p(\mathbf{x}_i | \boldsymbol{\theta}^*)$ は

$$\mathbb{X} \times \mathbb{X} \times \dots \times \mathbb{X} \ni (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) \mapsto \prod_{i=1}^n p(\mathbf{x}_i | \boldsymbol{\theta}^*) \in (0, \infty)$$

である. □

ここで, $\arg \max$ の記号を導入する. 関数 $g(x)$ の最大値を取る点を表す集合を

$$\arg \max_{x \in \mathbb{R}} g(x)$$

と書く. たとえば $g(x) = -(x-1)^2$ のとき

$$\arg \max_{0 \leq x \leq 4\pi} g(x) = \{1\}$$

となる. $g(x) = \sin x$ のとき

$$\arg \max_{0 \leq x \leq 4\pi} g(x) = \{\pi/2, 5\pi/2\}$$

となる.

定義 1.5. $X_1 = x_1, X_2 = x_2, \dots, X_n = x_n$ を観測したとき θ^* の最尤推定値 (maximum likelihood estimate) を $\text{lik}_n(\theta)$ を最大にする値 $\hat{\theta}_n(x_1, x_2, \dots, x_n)$ で定義する. すなわち

$$\hat{\theta}_n(x_1, x_2, \dots, x_n) \in \arg \max_{\theta \in \Theta} \text{lik}_n(\theta).$$

(x_1, x_2, \dots, x_n) に (X_1, X_2, \dots, X_n) を代入したものの $\hat{\theta}_n(X_1, X_2, \dots, X_n)$ を θ の最尤推定量 (maximum likelihood estimator=m.l.e.) という.

注意 1.6. $X_1, X_2, \dots, X_n \sim \text{i.i.d. Ber}(\theta^*)$ とする. ただし $(0, 1) =: \Theta \ni \theta^*$ は未知の真の母数とする. すなわち

$$p(x|\theta^*) = (\theta^*)^x (1 - \theta^*)^{1-x}, \quad (x = 0, 1)$$

である. このとき

$$\text{lik}_n(\theta) = \prod_{j=1}^n p(x_j|\theta) = \prod_{j=1}^n \theta^{x_j} (1 - \theta)^{1-x_j} = \theta^{t_n} (1 - \theta)^{n-t_n}$$

となる. ただし $t_n = \sum_{j=1}^n x_j$ である. よって対数尤度は

$$\ell_n(\theta) = t_n \log \theta + (n - t_n) \log(1 - \theta) \quad (0 < \theta < 1)$$

となる. このことから, $0 < t_n < n$ のとき,

$$\frac{t_n}{n} \in \arg \max_{\theta \in (0, 1)} \ell_n(\theta)$$

がわかる. したがって, $0 < t_n < n$ のとき, θ の最尤推定量は $\hat{\theta}_n = \frac{\sum_{j=1}^n X_j}{n}$ となる. 一方, $t_n = 0, n$ のとき, θ^* の最尤推定値は存在しない. \square

注意 1.7. $\theta^* = (\mu^*, \sigma^*) \in \mathbb{R} \times (0, \infty)$ を真の母数とし,

$$X_1, X_2, \dots, X_n \sim \text{i.i.d. N}(\mu^*, (\sigma^*)^2)$$

とする. $X_j = x_j (j = 1, 2, \dots, n)$ を観測したとき, 尤度関数は

$$\begin{aligned} \text{lik}_n(\mu, \sigma) &= \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x_j - \mu)^2}{\sigma^2}\right\} \\ &= \left(\frac{1}{\sqrt{2\pi}}\right)^n \frac{1}{\sigma^n} \exp\left\{-\frac{1}{2\sigma^2} \sum_{j=1}^n (x_j - \mu)^2\right\} \\ &= \left(\frac{1}{\sqrt{2\pi}}\right)^n \frac{1}{\sigma^n} \exp\left\{-\frac{ns_n^2}{2\sigma^2}\right\} \exp\left\{-\frac{n(\bar{x}_n - \mu)^2}{2\sigma^2}\right\}. \end{aligned}$$

ただし

$$\bar{x}_n = \frac{1}{n} \sum_{j=1}^n x_j, \quad s_n^2 = \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x}_n)^2.$$

最後の等号は

$$\sum_{j=1}^n (x_j - \mu)^2 = ns_n^2 + n(\bar{x}_n - \mu)^2 \quad (1.6)$$

からわかる. 対数尤度は

$$\ell_n(\mu, \sigma) = -n \log \sigma - \frac{ns_n^2}{2\sigma^2} - \frac{n(\bar{x}_n - \mu)^2}{2\sigma^2}.$$

よって

$$\begin{cases} \frac{\partial \ell_n(\mu, \sigma)}{\partial \mu}(\mu, \sigma) = -\frac{n(\bar{x}_n - \mu)}{2\sigma^2} = 0 \\ \frac{\partial \ell_n(\mu, \sigma)}{\partial \sigma}(\mu, \sigma) = -\frac{n}{\sigma} + \frac{ns_n^2}{\sigma^3} + \frac{n(\bar{x}_n - \mu)^2}{\sigma^3} = 0 \end{cases}$$

を解くと

$$\mu = \bar{x}_n, \quad \sigma = \sqrt{\frac{1}{n} \sum_{j=1}^n (x_j - \bar{x}_n)^2}$$

となる. $\ell_n(\mu, \sigma)$ の Hessian を求める.

$$\mathbf{H} := \begin{pmatrix} \frac{\partial^2 \ell_n}{\partial \mu^2}(\bar{x}_n, s_n) & \frac{\partial^2 \ell_n}{\partial \mu \partial \sigma}(\bar{x}_n, s_n) \\ \frac{\partial^2 \ell_n}{\partial \sigma \partial \mu}(\bar{x}_n, s_n) & \frac{\partial^2 \ell_n}{\partial \sigma^2}(\bar{x}_n, s_n) \end{pmatrix} = \begin{pmatrix} -\frac{1}{2\sigma^2} & 0 \\ 0 & -\frac{2n}{s_n^2} \end{pmatrix} \quad (1.7)$$

より, $-\mathbf{H}$ は正定値行列となるので, 関数

$$\Theta \ni (\mu, \sigma) \mapsto \ell_n(\mu, \sigma)$$

は $(\mu, \sigma) = (\bar{x}_n, s_n)$ で最大となる. よって (μ, σ) の最尤推定量は

$$\hat{\mu}_n = \frac{1}{n} \sum_{j=1}^n X_j =: \bar{X}_n, \quad \hat{s}_n = \sqrt{\frac{1}{n} \sum_{j=1}^n (X_j - \bar{X}_n)^2}$$

となる. □

問 1.1. (1.6) と (1.7) を確認せよ.