

確率・統計特論 第 10 回

来嶋 秀治 (Shuji Kijima)

締切: 2011 年 6 月 29 日講義終了時回収

注意: 参照した文献等の情報を必ず記載すること.

今日の話題: 線形回帰, 多次元分布, AIC

1. 回帰分析

1-i. 線形回帰分析: 単回帰分析. 確率変数 X, Y は未知パラメータ α, β に対して, $Y = \alpha X + \beta$ を満たすと想定する. いま, k 個の標本 $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_k)$ を得たとする. 最小二乗推定量 (least square estimator) $(\hat{\alpha}, \hat{\beta})$ を

$$(\hat{\alpha}, \hat{\beta}) = \arg \min_{\alpha', \beta'} \sum_{i=1}^k (y_i - (\alpha' x_i + \beta'_i))^2$$

を満たすものとする. これを解析的に解くと, 具体的には

$$\hat{\alpha} = s_{x,y}/s_x^2, \quad \hat{\beta} = \bar{y} - \hat{\alpha}\bar{x}$$

となる. ただし

$$\bar{x} = \sum_{i=1}^n x_i/n, \quad \bar{y} = \sum_{i=1}^n y_i/n, \quad \overline{x^2} = \sum_{i=1}^n x_i^2/n, \quad \overline{xy} = \sum_{i=1}^n x_i y_i/n, \quad s_x^2 = \overline{x^2} - \bar{x}^2, \quad s_{x,y} = \overline{xy} - \bar{x}\bar{y}$$

とする. remember $\text{Cov}[XY] = E[(X - E[X])(Y - E[Y])] = E[XY] - E[X]E[Y]$.

1-ii. 線形回帰分析: 重回帰分析. 確率変数 X_1, \dots, X_n は未知パラメータ $\alpha_1, \dots, \alpha_n, \beta$ に対して, $Y = \sum_{i=1}^n \alpha_i X_i + \beta$ を満たすと想定する. いま, k 個の標本 x^1, x^2, \dots, x^k を得たとする. このとき, $\alpha_1, \dots, \alpha_n, \beta$ を推定する手法が重回帰分析である (\Rightarrow 演習 1).

1-iii. その他の回帰分析. 一般に, 想定する関係式が線形である必要はない. たとえば, 2 つの確率変数 X, Y が未知パラメータ α, C に対して, $Y = CX^\alpha$ を想定した方がよい場合もある (\Rightarrow 演習 2).

線形回帰分析以外の回帰としてロジスティック回帰などが挙げられる.

2. 多次元分布

(i) 多次元正規分布 (multivariate normal distribution) $N(\mu, \Sigma)$

($\mu = (\mu_1, \dots, \mu_n)^\top \in \mathbb{R}^n$, $\Sigma = (\sigma_{i,j}) \in \mathbb{R}^{n \times n}$: 分散共分散行列 (variance-covariance matrix).)

$$f(x) = \frac{1}{(2\pi)^{n/2} \sqrt{\det(\Sigma)}} \exp \left(-\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right) \quad (x \in \mathbb{R}^n)$$

ただし

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{1,2} & \cdots & \sigma_{1,n} \\ \sigma_{1,2} & \sigma_2^2 & \cdots & \sigma_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{1,n} & \sigma_{2,n} & \cdots & \sigma_n^2 \end{pmatrix}$$

(ii) Dirichlet 分布 (Dirichlet distribution) $\text{Dir}(\alpha)$ ($\alpha = (\alpha_1, \dots, \alpha_n) \in \mathbb{R}_{\geq 0}^n$)

標本空間 $\Omega = \{x = (x_1, \dots, x_n) \in \mathbb{R}_{\geq 0}^n \mid x_1 + \dots + x_n = 1\}$,

$$f(x) = \frac{1}{B(\alpha)} \prod_{i=1}^n x_i^{\alpha_i - 1} \quad (x \in \Omega)$$

ただし,

$$B(\alpha) = \int_{\Omega} \prod_{i=1}^n x_i^{\alpha_i-1} dx = \frac{\prod_{i=1}^n \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^n \alpha_i)}$$

(iii) 多項分布 (multinomial distribution) $M(\mathbf{p}, K)$ ($\mathbf{p} = (p_1, \dots, p_n) \in \mathbb{R}_{>0}^n, \sum_{i=1}^n p_i = 1$)
標本空間 $\Omega = \{\mathbf{z} \in \mathbb{Z}_{\geq 0}^n \mid z_1 + \dots + z_n = K\}$

$$f(\mathbf{z}) = \frac{K!}{z_1! z_2! \dots z_n!} p_1^{z_1} \dots p_n^{z_n} \quad (\mathbf{z} \in \Omega)$$

多次元正規分布, Dirichlet 分布は連続分布, 多項分布は離散分布であることに注意.

3. モデル選択 — オッカムの剃刀 (Occam's razor)

AIC (赤池情報量基準): 統計データ (k 個の標本) $\mathbf{x}^1, \dots, \mathbf{x}^k$ に対して, モデル (確率分布) $f(\mathbf{x}; \theta)$ の AIC を

$$\begin{aligned} \text{AIC} &:= -2 \sum_{i=1}^n \log L(\hat{\theta}_1, \dots, \hat{\theta}_q; \mathbf{x}^i) + 2(q+2) \\ &= -2(\text{最大対数尤度}) + 2(\text{パラメータ数}) \end{aligned}$$

とする. AIC の小さなモデルが, より良いモデルと推定される.

同様の概念として BIC (ベイズ情報量基準), MDL(最小記述長) などがある.

演習問題

演習 1. 正の確率変数 X と Y は未知パラメータ α と C に対して, $Y = CX^\alpha$ を満たすものと想定される. いま k 個の標本 $(x_1, y_1), \dots, (x_k, y_k)$ が得られたとする. このとき α と C を推定せよ.

ヒント: $Y = CX^\alpha$ の両辺の対数を考えよ.

演習 2. 確率変数 X, Y, Z は未知パラメータ a, b, c に対して, $Z = aX + bY + c$ を満たすものと想定される. いま k 個の標本 $(x_1, y_1, z_1), \dots, (x_k, y_k, z_k)$ が得られたとする. このとき a, b, c を推定せよ.

参考文献

樺島祥介, 北川源四郎, 甘利俊一, 赤池弘次, 下平英寿, 土谷隆 (編), 室田一雄 (編), 赤池情報量基準 AIC—モデリング・予測・知識発見, 共立出版 (2007).