

# 応用システム工学

## 第二回 確率統計の基礎

平成22年06月04日

正規分布の分散

多変数の確率分布

一般化線形モデル

# 指数型分布族

- 確率分布が指数関数で表される

$$f(y; \theta) = s(y)\tau(\theta)e^{a(y)b(\theta)}$$

- 二項分布 → サイコロの目の出る確率

$$P(x) = {}_n C_x p^x (1-p)^{n-x}$$

- ポアソン分布 → 一定の期間において、ごくまれに起こる事象の確率分布

$$P(x) = \frac{\mu^x e^{-\mu}}{x!}$$

- 正規分布 → 自然界で起こる現象の多くがその分布に当てはまる(特に期待値に関する分布)

- 二項分布, ポアソン分布の正規分布による近似もある

# 正規分布から導かれる分布

- 推定量, 検定統計量の標本分布の多くは正規分布と関係
  - 正規分布に従う確率変数から導かれる場合
  - 大標本に対する中心極限定理より漸近的に関係する場合
    - 正規分布 → 正規分布そのもの
    - カイ二乗分布(自由度 $n$ ) → 標準正規分布に従う $n$ 個の独立な確率変数の二乗和の分布
    - $t$ 分布(自由度 $n$ ) → 2つの独立な確率変数の比の分布。分子は標準正規分布に従う確率変数, 分母は中心カイ二乗分布に従う確率変数を自由度で除したものの平方根
    - $F$ 分布 → 2つの独立したカイ二乗確率変数(自由度 $n, m$ )を各々の自由度で除したものの比
      - 中心 $F$ 分布 → 中心カイ二乗分布のみを用いる
      - 非心 $F$ 分布 → 分子が非心カイ二乗分布, 分母が中心カイ二乗分布

# 期待値と分散

- 期待値の線形性

- 前回正規分布関数の期待値を求めた

$$E[A + B] = E[A] + E[B]$$

- 期待値の線形性を用いた分散  $\sigma^2$  の性質

$$\begin{aligned}\sigma^2 &= V[X] = E[(X - E[X])^2] \\ &= E[X^2 - 2XE[X] + E[X]^2] \\ &= E[X^2] - 2E[X]E[X] + E[X]^2 \\ &= E[X^2] - E[X]^2\end{aligned}$$

# 正規分布関数の分散1

- 正規分布の確率密度関数(確率変数X)

$$f_X(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right]$$

- 分散

$$V[X] = \int_{-\infty}^{\infty} (x-\mu)^2 f_X(x) dx$$

$$= \int_{-\infty}^{\infty} (x^2 - 2x\mu + \mu^2) \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] dx$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} (x^2 - 2x\mu + \mu^2) \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] dx$$

# 正規分布関数の分散2

$$\begin{aligned} V[X] = & \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} x^2 \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] dx \\ & - \frac{2\mu}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} x \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] dx \\ & + \frac{\mu^2}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] dx \end{aligned}$$

# 正規分布関数の分散3

- 第3項 正規分布関数の確率密度関数1参照

$$\begin{aligned}\frac{\mu^2}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] dx &= \mu^2 \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] dx \\ &= \mu^2\end{aligned}$$

- 第2項 正規分布関数の期待値2参照

$$\begin{aligned}\frac{-2\mu}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} y \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] dx &= \frac{-2\mu}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} y \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] dx \\ &= -2\mu^2\end{aligned}$$

# 正規分布関数の分散4

- 第1項

– 変数変換

$$z = \frac{x - \mu}{\sigma} \quad \Rightarrow \quad \frac{dz}{dx} = \frac{d}{dx} \frac{x - \mu}{\sigma} = \frac{1}{\sigma}$$
$$x = \sigma z + \mu \quad x: -\infty \leftrightarrow \infty \Leftrightarrow z: -\infty \leftrightarrow \infty$$

$$\frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} x^2 \exp\left[-\frac{1}{2}\left(\frac{x - \mu}{\sigma}\right)^2\right] dx = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} (z\sigma + \mu)^2 e^{-\frac{1}{2}z^2} \sigma dz$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} (z^2\sigma^2 + 2z\sigma\mu + \mu^2) e^{-\frac{1}{2}z^2} \sigma dz$$

$$= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} (z^2\sigma^2 + 2z\sigma\mu + \mu^2) e^{-\frac{1}{2}z^2} dz$$

# 正規分布関数の分散5

$$\int_{-\infty}^{\infty} z e^{-\frac{1}{2}z^2} dz = 0 \qquad \int_{-\infty}^{\infty} \exp\left[-\frac{1}{2}z^2\right] dz = \sqrt{2\pi}$$

$$\begin{aligned} \int_{-\infty}^{\infty} z^2 e^{-\frac{1}{2}z^2} dz &= \int_{-\infty}^{\infty} z z e^{-\frac{1}{2}z^2} dz = \left[ z \left( -e^{-\frac{1}{2}z^2} \right) \right]_{-\infty}^{\infty} - \int_{-\infty}^{\infty} -e^{-\frac{1}{2}z^2} dz \\ &= [0 - 0] + \sqrt{2\pi} = \sqrt{2\pi} \end{aligned}$$

$$\begin{aligned} \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \left( z^2 \sigma^2 + 2z\sigma\mu + \mu^2 \right) e^{-\frac{1}{2}z^2} dz &= \frac{1}{\sqrt{2\pi}} \left( \sigma^2 \sqrt{2\pi} + 2\sigma\mu \times 0 + \mu^2 \sqrt{2\pi} \right) \\ &= \sigma^2 + \mu^2 \end{aligned}$$

# 正規分布関数の分散6

$$\begin{aligned} V[X] &= \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} x^2 \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] dx \\ &\quad - \frac{2\mu}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} x \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] dx \\ &\quad + \frac{\mu^2}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] dx \\ &= \{\sigma^2 + \mu^2\} + \mu^2 - 2\mu^2 = \sigma^2 \end{aligned}$$

# 多変量の確率分布1

## 2変数

- 同時確率密度関数

- 二つの確率変数 $X, Y$

$$x \leq X \leq x + dx$$

$$y \leq Y \leq y + dy \quad \text{かつ}$$

となる確率が

$$f_{XY}(x, y) dx dy \quad \text{と書ける}$$

- $X$ と $Y$ の同時確率密度関数

$$f_{XY}(x, y)$$

# 多変数の確率分布2

## n変数

- 同時確率密度関数

- n個の確率変数

$$X_1, X_2, \dots, X_n$$

- 各々  $x_i \leq X_i \leq x_i + dx_i (1 \leq i \leq n)$  となる確率が

$$f_{X_1 X_2 \dots X_n}(x_1, x_2, \dots, x_n) dx_1 dx_2 \dots dx_n$$

- $X_1, X_2, \dots, X_n$  の同時確率密度関数

$$f_{X_1 X_2 \dots X_n}(x_1, x_2, \dots, x_n)$$

# 多変数の確率分布3

## 2変数

- 周辺密度関数
  - 二つの確率変数 $X, Y$
  - $Y$ を無視した $X$ のみに関する確率密度関数を指す

$f_X(x)$

- 同時確率密度関数  $f_{XY}(x, y)$   
と周辺密度関数  $f_X(x)$   
の関係

$$f_X(x) = \int_{-\infty}^{\infty} f_{XY}(x, y) dy$$

# 多変数の確率分布4

## n変数

- 周辺密度関数

- n個の確率変数

$$X_1, X_2, \dots, X_n$$

- m個のみに注目した同時確率密度

$$f_{X_1 X_2 \dots X_m}(x_1, x_2, \dots, x_m)$$

- $X_1, X_2, \dots, X_m$  についての周辺密度関数

$$f_{X_1 \dots X_m}(x_1, \dots, x_m) = \int_{-\infty}^{\infty} dx_{m+1} \cdots \int_{-\infty}^{\infty} dx_n f_{X_1 \dots X_n}(x_1, \dots, x_n)$$

# 確率変数の独立性1

- 二つの確率変数 $X, Y$
- 独立性の定義

$$\Leftrightarrow P(X \in A \text{かつ} Y \in B) = P(X \in A)P(Y \in B)$$

$$\Leftrightarrow f_{XY}(x, y) = f_X(x)f_Y(y)$$

# 確率変数の独立性2

- n個の確率変数

$$X_1, X_2, \dots, X_n$$

- 独立性の定義

$$\Leftrightarrow P(X_1 \in A_1 \text{かつ} X_2 \in A_2 \text{かつ} \dots \text{かつ} X_n \in A_n)$$

$$= P(X_1 \in A_1)P(X_2 \in A_2) \cdots P(X_n \in A_n)$$

$$\Leftrightarrow f_{X_1 X_2 \dots X_n}(x_1, x_2, \dots, x_n) = f_{X_1}(x_1) f_{X_2}(x_2) \cdots f_{X_n}(x_n)$$

- 自由度

– 独立に選べる確率変数の数

# 条件付き確率1

- 二つの確率変数 $X, Y$ が必ずしも独立でない

$$f_{XY}(x, y) = f(y|x)f_X(x)$$

–  $Y$ の $X$ に関する条件付き確率密度  $f(y|x)$

- $X$ が値 $x$ を取る条件下で $Y$ が値 $y$ を取る確率密度

–  $X$ と $Y$ が独立な場合  $f(y|x) = f_Y(y)$

- 同時確率密度 $f_{XY}(x, y)$ と条件付確率密度 $f(y|x)$ の関係

$$f(y|x) = \frac{f_{XY}(x, y)}{f_X(x)} = \frac{f_{XY}(x, y)}{\int_{-\infty}^{\infty} f_{XY}(x, y)dy}$$

# 条件付き確率2

- 条件  $X_1 = x_1, X_2 = x_2, \dots, X_m = x_m$   
に対する,  $X_{m+1} = x_{m+1}, X_{m+2} = x_{m+2}, \dots, X_n = x_n$   
となる確率密度関数

$$f(x_{m+1}, x_{m+2}, \dots, x_n \mid x_1, x_2, \dots, x_m)$$

$$\begin{aligned} & f(x_1, x_2, \dots, x_m, x_{m+1}, \dots, x_n) \\ &= f(x_{m+1}, x_{m+2}, \dots, x_n \mid x_1, x_2, \dots, x_m) \\ &\times f(x_1, x_2, \dots, x_m) \end{aligned}$$

# 確率変数の和の分布

- 二つの独立な確率変数 $X, Y$
- $Z=X+Y$ の分布

$X, Y$ が独立の時,  $Y=z-x$ かつ $X=x$ となる確率

– 事象  $Z = z$

$$\Leftrightarrow X + Y = z$$

$$\Leftrightarrow Y = z - x \text{ かつ } X = x$$

$$P_X(x)P_Y(z-x)$$

– 離散分布  $X \in M$

–  $Z=z$ となる確率

$$P(Z = z) = \sum_{x \in M} P_X(x)P_Y(z-x)$$

# 用語(最初にするべき)

- 測定値または観測値
  - － 説明変数
  - － 予測変数
  - － 独立変数
- 確率変数(上述の変数に応じて変動)
  - － 反応変数
  - － 結果変数
  - － 従属変数

# 説明変数の種類

- 一般化線形モデル

- 反応変数 $Y$ と説明変数 $x_1, x_2, \dots, x_m$ の関係

$$g[E(Y)] = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m$$

- 量的な説明変数 $x$

- パラメータ $\beta$ は、説明変数の変化に伴う反応変数の変化を表す

- 質的な説明変数

- 反応変数にパラメータが含まれるか否かを表す
    - ダミー変数,  $(0, 1)$ の場合は指示変数

# 相関係数

- 二つの確率変数の中の類似の度合いを示す統計学的な指標

– 二組の数値からなるデータ列

$$(x, y) = \{(x_i, y_i)\}, i = 1, \dots, n$$

– 相関係数  $\rho_{ij}$

$$\rho_{ij} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

$\bar{x}, \bar{y}$   
は各々の相加平均

# 共分散

- 共分散
  - 二組の対応するデータ間での、平均からの偏差の積の平均値

# 多変量正規分布

- 正規分布の確率変数  $Y_i \sim N(\mu_i, \sigma_i^2), i = 1, \dots, n$
- 確率変数  $Y_i, Y_j$  間の共分散

$$\text{cov}(Y_i, Y_j) = \rho_{ij} \sigma_i \sigma_j \quad \rho_{ij} : Y_i \text{ と } Y_j \text{ の相関係数}$$

–  $Y_1, Y_2, \dots, Y_n$  の同時分布は多変量正規分布となる

- 確率変数ベクトル  $y = [Y_1, \dots, Y_n]^t$
- 平均ベクトル  $\mu = [\mu_1, \dots, \mu_n]^t$
- 分散共分散行列  $V$ : 対角要素  $\sigma_i^2$ , 非対角要素  $\rho_{ij} \sigma_i \sigma_j$

$$y \sim N(\mu, V)$$

# 正規分布の線形結合

- 正規分布の確率変数  $Y_i \sim N(\mu_i, \sigma_i^2), i = 1, \dots, n$
- 確率変数の線形結合  $W$

$$W = a_1 Y_1 + a_2 Y_2 \cdots + a_n Y_n$$

– ただし  $a_1, a_2, \dots, a_n$  は定数

–  $W$  は正規分布となる

- 平均  $\sum_{i=1}^n a_i \mu_i$
- 分散  $\sum_{i=1}^n a_i^2 \sigma_i^2$

$$W = \sum_{i=1}^n a_i Y_i \sim N\left(\sum_{i=1}^n a_i \mu_i, \sum_{i=1}^n a_i^2 \sigma_i^2\right)$$

# カイ二乗分布1

- 自由度 $n$ の中心カイ二乗分布

- 独立な $n$ 個の標準正規分布の確率変数  $Z_1, \dots, Z_n$  の二乗和の分布

$$X^2 = \sum_{i=1}^n Z_i^2 = \mathbf{z}^T \mathbf{z} \sim \chi^2(n)$$

- ただし  $\mathbf{z} = [Z_1, \dots, Z_n]^t$        $\mathbf{z}^T \mathbf{z} = \sum_{i=1}^n Z_i^2$

- 期待値       $E(X^2) = n$

- 分散       $\text{var}(X^2) = 2n$

# カイ二乗分布

- 独立なn個の正規分布の確率変数  $Y_1, \dots, Y_n$  に対して  $Y_i \sim N(\mu_i, \sigma_i^2)$ 
  - 平均  $\mu_i$ , 分散  $\sigma_i^2$
  - 標準正規分布化  $Z_i = \frac{Y_i - \mu_i}{\sigma_i} = N(0,1)$
  - カイ二乗分布

$$X^2 = \sum_{i=1}^n \left( \frac{Y_i - \mu_i}{\sigma_i} \right)^2 \sim \chi^2(n)$$

# 非心カイ二乗分布

- 標準正規分布の確率変数  $Z_1, \dots, Z_n$
- 正規分布の確率変数化  $Y_i = Z_i + \mu_i$ 
  - 確率変数の二乗和

$$\sum_{i=1}^n Y_i^2 = \sum_{i=1}^n (Z_i + \mu_i)^2 = \sum_{i=1}^n Z_i^2 + 2 \sum_{i=1}^n Z_i \mu_i + \sum_{i=1}^n \mu_i^2 \sim \chi^2(n, \lambda)$$

- 自由度:  $n$
- 非心パラメータ  $\lambda = \sum_{i=1}^n \mu_i^2$
- 平均  $n + \lambda$
- 分散  $2n + 4\lambda$

# カイ二乗分布

- 多変量正規分布の確率変数ベクトル

$$y = [Y_1, \dots, Y_n]^t \quad y \sim N(\mu, V)$$

- 分散共分散行列 $V$ が正則, 逆行列を持つ時

$$X^2 = (y - \mu)^T V^{-1} (y - \mu) \sim x^2(n)$$

- 一般に多変量正規分布に対して

確率変数  $y^T V^{-1} y$  は非心カイ二乗分布となる

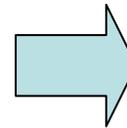
$$x^2(n, \lambda) \quad \lambda = \mu^T V^{-1} \mu$$

# カイ二乗分布の再生性

- 中心・非心カイ二乗分布に従う独立な確率変数  $X_i^2 \sim \chi_i^2(n_i, \lambda_i)$  の和のカイ二乗分布

– 自由度  $\sum_{i=1}^n n_i$   
– 非心パラメータ  $\sum_{i=1}^n \lambda_i$

$$\sum_{i=1}^m X_i^2 \sim \chi_i^2\left(\sum_{i=1}^n n_i, \sum_{i=1}^n \lambda_i\right)$$



カイ二乗分布の再生性

# カイ二乗分布の性質

- 多変量正規分布の確率変数(要素数n)

$$y \sim N(\mu, V)$$

- 分散共分散行列Vの階数 $k < n$ で特異の場合
  - 逆行列が定まらない
  - Vの一般逆行列 $V^-$ に対して
  - 確率変数  $y^T V^- y$  の性質 → 非心カイ二乗分布
    - 自由度k
    - 非心パラメータ  $\lambda = \mu^T V^- \mu$







# ポアソン分布2

- ポアソン分布に対する尤度関数  $L(\theta; y_1, \dots, y_n)$
- 最尤推定値  $\hat{\theta}$
- 対数尤度関数

$$\begin{aligned} l(\theta; y_1, \dots, y_n) &= \log L(\theta; y_1, \dots, y_n) \\ &= (\sum y_i) \log \theta - n\theta - \sum (\log y_i!) \end{aligned}$$

- 対数尤度関数の導関数

$$\begin{aligned} \frac{d}{d\theta} l(\theta; y_1, \dots, y_n) &= \frac{d}{d\theta} [(\sum y_i) \log \theta - n\theta - \sum (\log y_i!)] \\ &= \frac{1}{\theta} \sum y_i - n \end{aligned}$$

# ポアソン分布3

- 最尤推定値が極値を取る条件

$$\frac{d}{d\theta} l(\hat{\theta}; y_1, \dots, y_n) = \frac{1}{\hat{\theta}} \sum y_i - n = 0$$

$$\hat{\theta} = \frac{\sum y_i}{n} = \bar{y}$$

- 最大推定値の条件

$$\frac{d^2}{d\theta^2} l(\hat{\theta}; y_1, \dots, y_n) = -\frac{1}{\theta^2} \sum y_i < 0$$

$\theta = \hat{\theta}$  で最大となる。最尤推定値は  $\bar{y}$



# 正定値

- 二次形式  $y^T Ay$  と行列A
- 要素の全てが0でないy
- 正定値  $y^T Ay > 0$ 
  - 必要十分条件
    - 全ての行列式が正

$$|A_1| = a_{11} \quad |A_2| = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \quad |A_3| = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$

行列Aの階数は二次形式  $Q = y^T Ay$  の自由度  $|A_n| = \det A$

# Cochranの定理

- 正規分布  $N(0, \sigma^2)$  の独立な確率変数  $Y_1, \dots, Y_n$

- 二次形式  $Q = \sum_{i=1}^n Y_i^2$

– 和  $Q = Q_1 + \dots + Q_k$

– ただし  $Q_i$  の自由度  $m_i (i = 1, \dots, k)$   
 $m_1 + m_2 + \dots + m_k = n$  の時に限り

$Q_i, \dots, Q_k$  は独立な確率変数

$$\frac{Q_1}{\sigma^2} \sim x^2(m_1), \frac{Q_2}{\sigma^2} \sim x^2(m_2), \dots, \frac{Q_k}{\sigma^2} \sim x^2(m_k)$$

# Cochranの定理より

- 二つの確率変数

- ただし  $m > k$   $X_1^2 \sim x^2(m), X_2^2 \sim x^2(k)$

- 差が非負定値

$$X^2 = X_1^2 - X_2^2 \geq 0$$

- なりたつ

$$X^2 \sim x^2(m - k)$$