

# 応用システム工学

## 第三回 確率統計の基礎

平成24年05月21日  
多変数の確率分布  
一般化線形モデル  
相関分析  
回帰分析

2012/05/21

1

### 用語

- 測定値または観測値(原因)
  - 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:結果)
  - 目的変数
  - 反応変数
  - 結果変数
  - 従属変数

2012/05/21

2

# 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 $X, Y$

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

$$y \leq Y \leq y + dy \quad \text{かつ} \\ \text{となる確率が}$$

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

- $X$ と $Y$ の同時確率密度関数

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

2012/05/21

3

# 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

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

2012/05/21

4

# 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

- 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

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

2012/05/21

5

# 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( $m < n$ )

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

- $X_1, X_2, \dots, X_m$  についての周辺密度関数

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

2012/05/21

6

# 確率変数の独立性1

- 二つの確率変数 $X, Y$
- 独立性の定義

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

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

2012/05/21

7

# 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

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

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

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

- 自由度
  - 独立に選べる確率変数の数

2012/05/21

8

# 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

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

2012/05/21

9

# 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

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

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

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

$$\times f(x_1, x_2, \dots, x_m)$$

2012/05/21

10

# 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

– 事象  $Z = z$

$X, Y$ が独立の時,  $Y=z-x$ かつ $X=x$ となる確率

$$\Leftrightarrow X + Y = z$$

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

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

– 離散分布  $X \in M$

–  $Z=z$ となる確率

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

2012/05/21

11

# 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

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

線形関数

- 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

多変量解析  $\Rightarrow$  相関分析, 回帰分析

2012/05/21

12

# 相関解析 相関係数

- 二つの確率変数の中の類似の度合い(相関関係の強さ)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  $\rho_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$   $\bar{x}, \bar{y}$

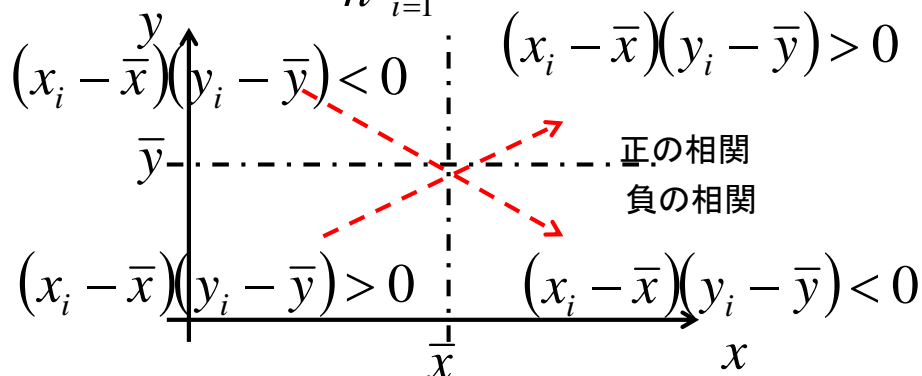
$\rho_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$  は各々の相加平均

$$|\rho_{xy}| < 1$$

# 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

$$\text{cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$



# 回帰分析とは

- データ
  - 目的変数 $y$
  - 目的変数に影響をおよぼす説明変数 $x$
- 分析
  - 予測式を求める  
(ある変数の変化をもう一方の変数の変化で説明するための関数を求める)

$$a_0 + a_1x_1 + \cdots + a_px_p \rightarrow y$$

2012/05/21

15

## 線形回帰(直線回帰)

- 説明変数の数による回帰分析の分類
  - 1個 → 単回帰分析
  - 2個以上 → 重回帰分析

$n$ 個のデータの  
目的変数 $y$ と  
説明変数 $x$ の組

例  
地区数 $n$   
地区 $i$ の  
世帯数 $x_i$ , ごみの量 $y_i$

	目的変数: $y$	説明変数: $x$
1	$y_1$	$x_1$
2	$y_2$	$x_2$
...		
$i$	$y_i$	$x_i$
...		
$n$	$y_n$	$x_n$

2012/05/21

16



# 線形回帰モデル

- モデル式

$$y_i = a_0 + a_1x_i + e_i \quad (i = 1, 2, \dots, n)$$

- 未知の定数  $a_0, a_1$

- 予測誤差  $e_i$

- 予測誤差が最小となる定数  $\hat{a}_0, \hat{a}_1$  を求める

- 予測誤差の平方和を最小にする

$$\sum_{i=1}^n e_i^2 = \sum_{i=1}^n \{y_i - (a_0 + a_1x_i)\}^2$$

- 最小二乗法 →  $\hat{a}_0, \hat{a}_1$