薬学情報処理演習 第2回

表計算ソフトによる統計 処理



奥菌 透 コロイド・高分子物性学



離散的なデータ

- 度数分布 F_i
- 確率 $f_i = F_i/N$

$$\bar{x} = \sum_{i=1}^{n} X_i f_i$$
 $\sigma^2 = \sum_{i=1}^{n} (X_i - \bar{x})^2 f_i$

データの分布

□ 連続的なデータ

- 度数分布→分布関数 F(x) $(N \to \infty, \Delta x = X_{i+1} - X_i \to 0)$ - 確率密度 $f(x) = F(x)/\tilde{N}$ $(\widetilde{N} = \int F(x) dx)$

 $\bar{x} = \int x f(x) dx$ $\sigma^2 = \int (x - \bar{x})^2 f(x) dx$



2

3

4

5







□ 中心極限定理

n 個の独立な確率変数 u_i (分散 s_i^2 平均値0)からなる確率変数

 $x_n = (u_1 + u_2 + \dots + u_n)/\sqrt{\sigma_n^2}$ $\sigma_n^2 = s_1^2 + s_2^2 + \dots + s_n^2$ は、 $n \to \infty$ で分散1, 平均値0の正規分布に従う。 □ ばらつき=多数の確率的事象の和 x



コンピュータ上でランダムな数(乱数)を次々に生成し、ランダムなデータを作ることができる。エクセルでは RAND()という関数が用意されている。

□ RAND()で生成される乱数は一様分布関数 f(x) = 1 $(0 \le x < 1)$

に従い、平均と分散は、

$$\overline{x} = \int_0^1 x f(x) dx = \frac{1}{2} \qquad \sigma^2 = \int_0^1 (x - \overline{x})^2 f(x) dx = \frac{1}{12}$$

となるので、平均0の一様乱数(RAND()-0.5)を12個 足し合わせたものは、近似的に、平均0分散1の正規 分布に従う乱数(正規乱数)となっている。

Excelでの正規乱数発生方法

- □ データ分析ツールの利用(関数ではない)
 - データ/分析/データ分析/乱数発生
- □ 中心極限定理の応用
 - (12個の一様乱数の和)-6
 - =rand()+rand()+...+rand()-6.0

□ 逆関数法

- =norm.inv(rand(),0,1)または =norm.s.inv(rand())を用いる
- 分布f(x)の累積分布をF(x)とする

$$y = F(x) = \int_{-\infty}^{x} f(u) du$$

$$\frac{dy}{dx} = f(x), \ dy = f(x)dx$$





- □ 上記の3つの方法で、それぞれN個の正規乱数
 を発生させる(N = 10000 程度)。
- 上で得られたデータに対する分布関数のグラフを描く。
 - 分析ツール(後述)を用いて、度数分布表を作る。
 - 規格化された分布関数のデータを計算する。
 - 得られた分布関数のデータをグラフに描き、理論曲線 と比較する。
 - (余裕があれば)平均値と分散を計算し、理論値と比 較する。



□ 分析ツールを使って正規乱数を発生

RAND()を使って12個の一様乱数を足し合わせ て正規乱数を生成

□ 逆関数法により正規乱数を発生





| □ エクセルで度数分布を作る方法に | はいろいろある |
|-------------------|----------------|
| が、ここでは「分析ツール」を使う。 | 。これを使用可 |
| 能とするには、ファイル/オプション | ン/アドイン/設 |
| 定 で「分析ツール」を選択し「OK | 、」をクリックする。 |
| データ区間(級)を作成する | データ区間 |

| □ 分析ツールを使う | -4.5 |
|-----------------------|------|
| - データ/分析/データ分析/ヒストグラム | -4.3 |
| - 入力範囲、データ区間を指定 | -4.1 |

- 出力先を選択・指定

9

-3.9

. . .

4.3



| データ区間 | 頻度 | 代表值 x | 分布 f(x) | |
|-------|------------------|--------------|-------------|---------------------------|
| -4.5 | 0 | | | $f_n = F_n / \mathcal{N}$ |
| -4.3 | 1 | =(A2+A3)/2 | =B2/\$B\$49 | — |
| | | | | |
| 4.5 | 0 | | | |
| 次の級 | 0 | | | |
| 積分値 | =SUM(B2:B47)*0.2 | | | |
| | | ∆x:データ区間 | 司の増分値 | 10 |



□ 正規分布のデータを作成する。

| 代表值 x | 正規分布 |
|-------|-----------------------------|
| -4.5 | =NORM.DIST(C3, 0, 1, FALSE) |
| | |

- 関数NORM.DIST($x, \overline{x}, \sigma^2$, FALSE) はxの値に対する平均 \overline{x} 分 散 σ^2 の正規分布 f(x)の値を返す。
- NORM.S.DIST(x, FALSE) は NORM.DIST(x, 0, 1, FALSE) と同じ。
- 最後の引数がFALSEの場合、(密度)分布関数を、TRUEの場合、 累積分布関数を返す。



□ 分布関数のグラフ

- 横軸に代表値、縦軸に確率密度分布をとる。
- 理論値と度数分布から得られたデータを比較する。



| 方法 | 平均 | 分散 |
|----------|----------|----------|
| theory | 0 | 1 |
| excel | -0.00982 | 1.007028 |
| rand | 0.01078 | 1.019344 |
| inv.func | -0.02518 | 1.015754 |

データ数: N = 10000

〜 分析ツールを使わずに分布関数の データを作成する。(補足)

□ 元データの作成(データの作成参照)

□ COUNTIF(範囲,検索条件)を使って累積度数分 布を作成。規格化して累積分布関数を得る。

□ 累積分布関数を微分して確率密度関数を得る。

| 1 正規乱数 データ区間 累積度数 累積分布 代表値 確率密度 2 -1.52061 -4.5 0 0 3 -1.37599 -4.3 0 0 -4.4 0 3 -1.37599 -4.3 0 0 -4.4 0 0 -4.4 0 3 -1.37599 -4.3 0 0 -4.4 0 0 -4.4 0 0 -4.4 0 0 -4.4 0 0 -4.4 0 0 -4.4 0 0 0 -4.4 0 0 0 -4.4 0 | |
|---|--------------------|
| 1 正規乱数 データ区間 累積度数 累積分布 代表値 確率密度 2 -1.52061 -4.5 0 0 3 -1.37599 -4.3 0 0 -4.4 0 3 -1.37599 -4.3 0 0 -4.4 0 下であるデータの数 | \$2.\$R\$ |
| 2 -1.52061 -4.5 0 0 0001,"<="&C2) | ΨΖ.ΨΟΨ |
| 3 -1.37599 -4.3 0 0 -4.4 0 下であるデータの数3 | つの値に |
| 2411254 41 0 0 42 0 [E2]-D2/D\$47 | 2の値以 を返す。 |
| | |
| 5 -0.64805 -3.9 0 0 -4 0 累積度数の最後の値 | 直(=全 112/14-1-7 |
| 6 0.113184 -3.7 0 0 -3.8 0 [G3]=(E3-E2)/0.2 | 16169 6 |
| 7 0.397722 -3.5 2 0.0002 -3.6 0.001 累積分布を微分する | 。0.2は |
| 8 0.260922 -3.3 6 0.0006 -3.4 0.002 データ区間の増分値 | 0 |