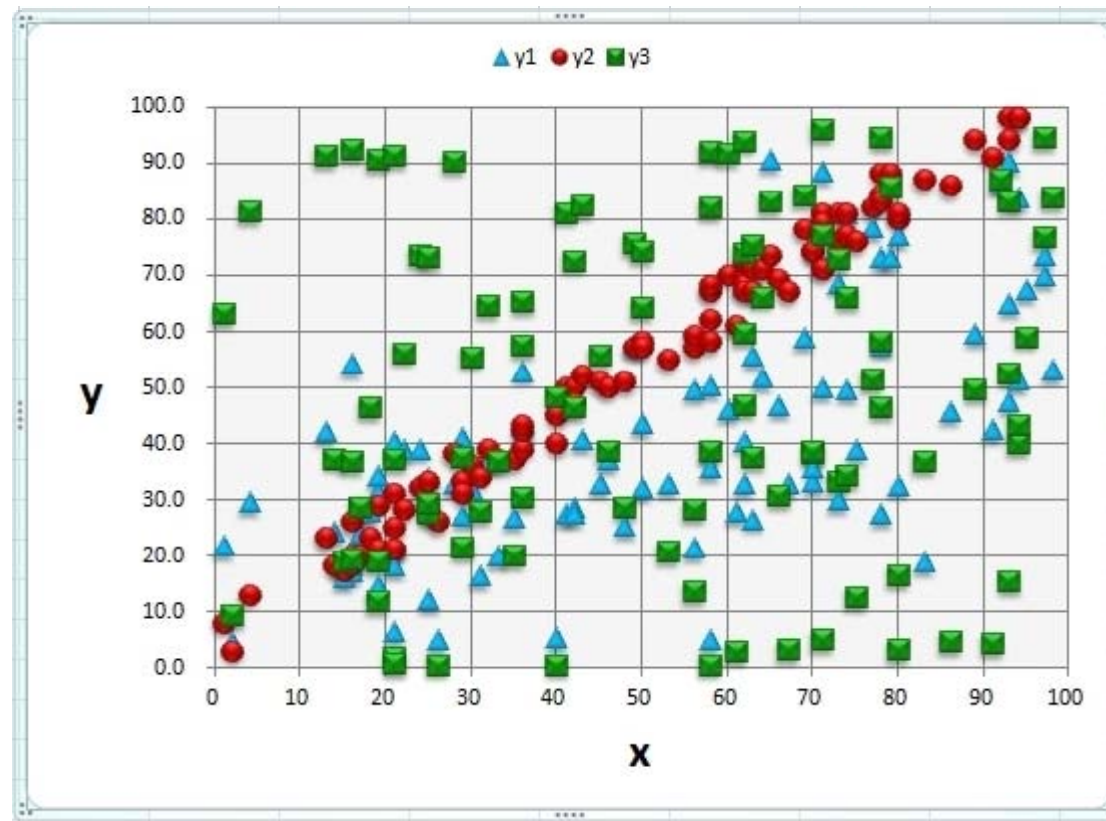


相関関係と相関係数

v0.5

舘野

相関グラフ(散布図)の例



相関関係

y

販売本数

気温と缶ビールの販売本数



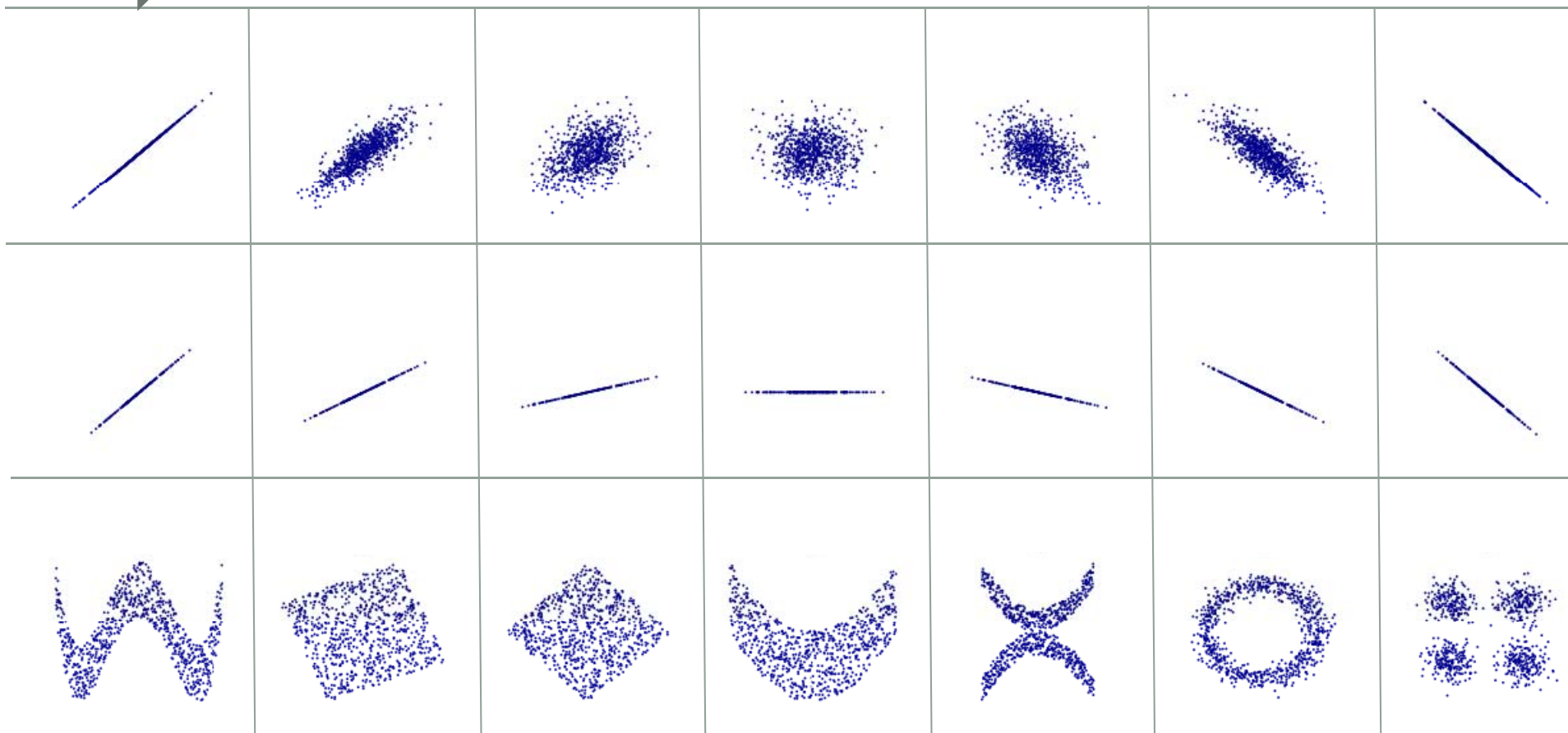
気温

x

相関関係を調べるには、
相関グラフを書いてみるのが
分かりやすい。。。
Excelの“散布図”で
簡単に作れます。

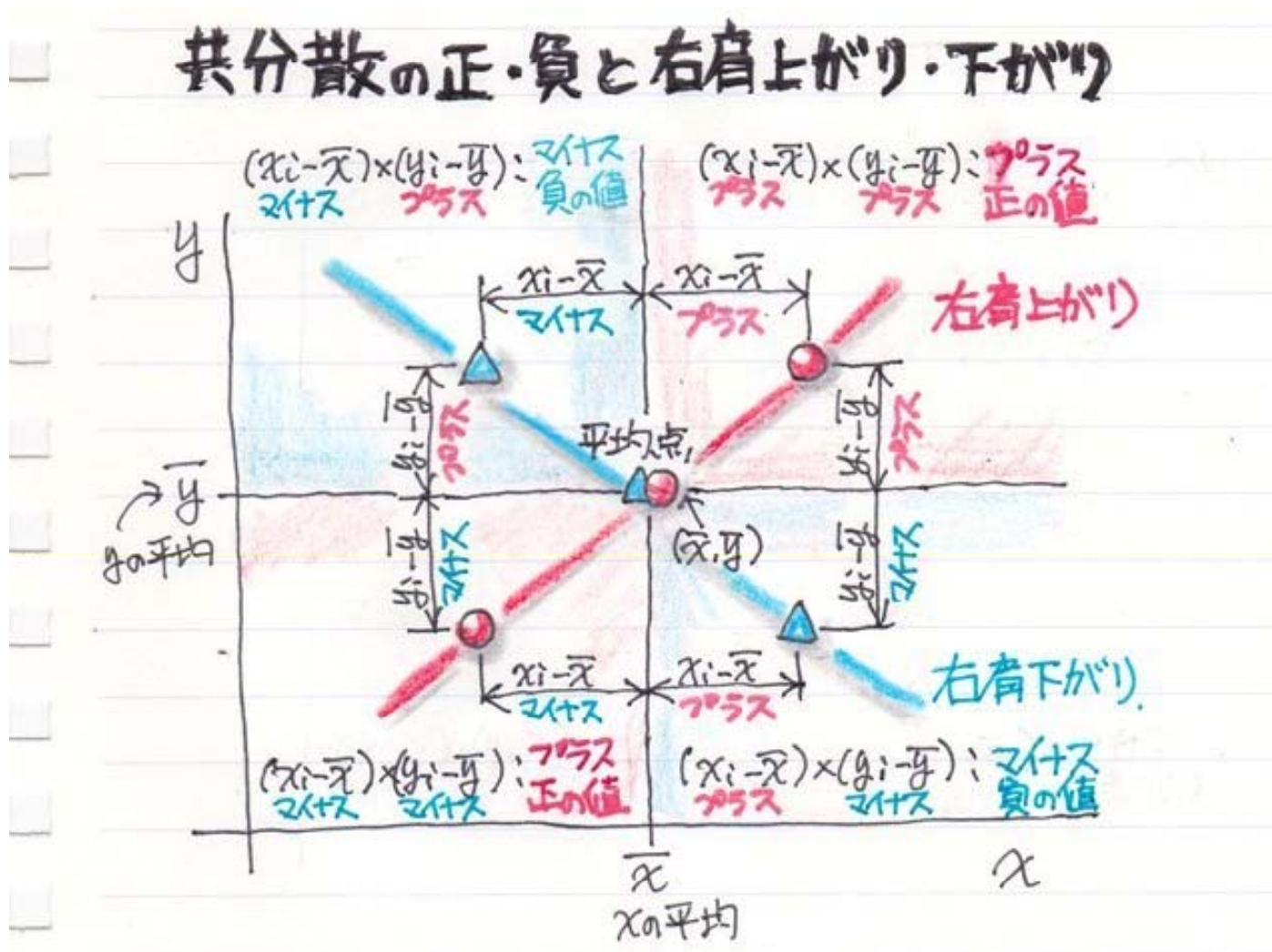


このグラフの相関係数はいくつでしょう？



共分散(きょうぶんさん、covariance)は、2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である。

<https://ja.wikipedia.org/wiki/共分散>



<http://haku1569.seesaa.net/archives/201410-1.html>

$$\text{共分散} = \left\{ \begin{aligned} & (x \text{ のデータ } 1 - x \text{ の平均}) \times (y \text{ のデータ } 1 - y \text{ の平均}) \\ & + (x \text{ のデータ } 2 - x \text{ の平均}) \times (y \text{ のデータ } 2 - y \text{ の平均}) \\ & + \dots + (x \text{ の最後のデータ} - x \text{ の平均}) \times (y \text{ の最後のデータ} - y \text{ の平均}) \end{aligned} \right\}$$

÷ データの個数



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

$\frac{1}{n}$: データの個数
 $\sum_{i=1}^n$: i が 1 ~ n までの合計
 x_i : x の i 番目のデータ
 \bar{x} : x の平均
 y_i : y の i 番目のデータ
 \bar{y} : y の平均

$$\text{相関係数} = \frac{\text{共分散}}{x\text{の標準偏差} \times y\text{の標準偏差}}$$

ここで、
共分散 \leq xの標準偏差 と yの標準偏差 の積

Σ : i = 1 から i = n までの総和
 x_m : x の平均、 $x_m = \Sigma x_i / n$
 y_m : y の平均、 $y_m = \Sigma y_i / n$
 sqrt: 平方根

共分散
 ((x_i と x_m との差)と
 ((y_i と y_m との差)の積
 の合計の平均

yの標準偏差 $\Sigma (y_i - y_m)^2$
 ((y_i と y_m との差)の2乗の合計)の平均の平方根

xの標準偏差 $\text{sqrt}[\Sigma (x_i - x_m)^2]$
 ((x_i と x_m との差)の2乗の合計)の平均の平方根

$$\text{相関係数} = \frac{\text{共分散}}{x\text{の標準偏差} \times y\text{の標準偏差}}$$

相関係数は2つの方法で求められますが、その結果は同じになります。

	偏差積和	共分散
	$S_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y})$	$s_{xy}^2 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1}$
相関係数	$r_{xy} = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}}$ <p style="text-align: center;"> $S_{xx} = \sum (x_i - \bar{x})^2$ 偏差平方和 </p>	$r_{xy} = \frac{s_{xy}^2}{\sqrt{s_x^2 s_y^2}}$ <p style="text-align: center;"> $s_x^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$ 分散 </p>

相関係数の値はなぜ $[-1 \sim +1]$ の範囲か(証明)

相関係数: r (correlation coefficient) は2つの変数間の相関すなわち類似性を示す指標で、次式で定義されます。

実数データ: (x_i, y_i) , $i = 1, 2, \dots, n$ に対して
$$r = \frac{\sum (x_i - x_m)(y_i - y_m)}{\sqrt{\sum (x_i - x_m)^2 \cdot \sum (y_i - y_m)^2}}$$

ここで、

\sum : $i = 1$ から $i = n$ までの総和

x_m : x の平均、 $x_m = \sum x_i / n$

y_m : y の平均、 $y_m = \sum y_i / n$

sqrt: 平方根

相関係数の値の範囲の証明には、**Schwartz** (シュワルツ) の不等式(下記注)や多次元ベクトルの内積を利用した方法などが知られていますが、ここでは2次関数の判別式を利用する簡単な方法を紹介します。

$|r| \leq 1$ の証明

変数 t を含む次の式:

$$Q = \sum [(x_i - x_m) + t(y_i - y_m)]^2$$

を展開すると、変数 t に関する2次式が得られます。

$$Q = \sum (x_i - x_m)^2 + 2t \sum (x_i - x_m)(y_i - y_m) + t^2 \sum (y_i - y_m)^2$$

Q の値は実数値の2乗和ゆえ、正または0です。

従って、上記2次式の判別式 D は負または0でなければなりません。

$$D = [\sum (x_i - x_m)(y_i - y_m)]^2 - \sum (x_i - x_m)^2 \sum (y_i - y_m)^2 \leq 0$$

$$[\sum (x_i - x_m)(y_i - y_m)]^2 \leq \sum (x_i - x_m)^2 \sum (y_i - y_m)^2$$

$$|\sum (x_i - x_m)(y_i - y_m)| \leq \sqrt{\sum (x_i - x_m)^2 \cdot \sum (y_i - y_m)^2}$$

$$|\sum (x_i - x_m)(y_i - y_m)| / \sqrt{\sum (x_i - x_m)^2 \cdot \sum (y_i - y_m)^2} \leq 1$$

よって、

$$-1 \leq \frac{\sum (x_i - x_m)(y_i - y_m)}{\sqrt{\sum (x_i - x_m)^2 \cdot \sum (y_i - y_m)^2}} \leq 1$$

となり、相関係数の値が $-1 \sim +1$ の範囲であることが証明されました。

主成分分析(principal component analysis:PCA)とは？

データの分散(ばらつき)が大きいところ(主成分)を見つける操作。つまり分散が大きいところ
が大事で、小さいところは気にしないようにする。

分散: $((x_i - \bar{x})^2)$ の合計の平均

