

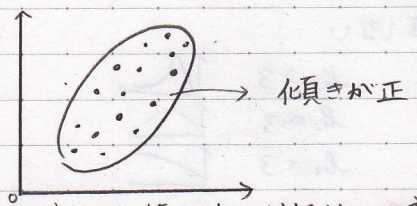
2次元データの要約

2次元データ $\rightarrow (x_1, y_1) (x_2, y_2) \dots (x_n, y_n)$

2つの種のデータが組み合わせられたもの。

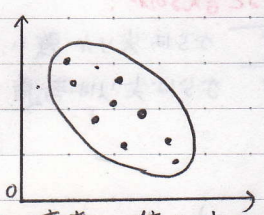
散布図

正の相関



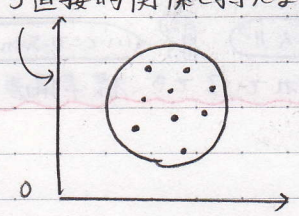
このとき、一方の変数の値が大きければ、一方の変数も大きい。

負の相関



このとき、一方の変数の値が大きければ、一方の変数は小さくなる。

2つのデータが、何ら直接的関係を持たないとき、無相関であるという。

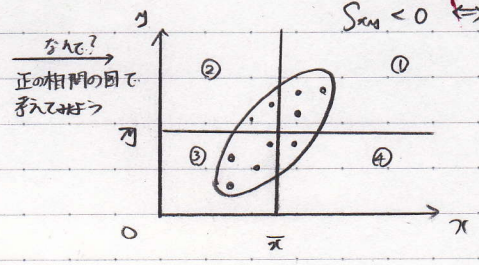


ここで、次式で表されり量を x と y の 共分散 と言う

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

ここで、 $S_{xy} > 0 \Leftrightarrow$ 正の相関

$S_{xy} < 0 \Leftrightarrow$ 負の相関が成立する



①③の領域において、 $(x_i - \bar{x})(y_i - \bar{y})$ は正となり、

②④の領域において、 $(x_i - \bar{x})(y_i - \bar{y})$ は負となるが、

この図の場合、 $(x_i - \bar{x})(y_i - \bar{y})$ が正となるデータの数が多い!

そうなるので $S_{xy} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) > 0$ となり、これは正の相関であると分かり、また、この逆も成立する。