

盛岡三高数学科通信

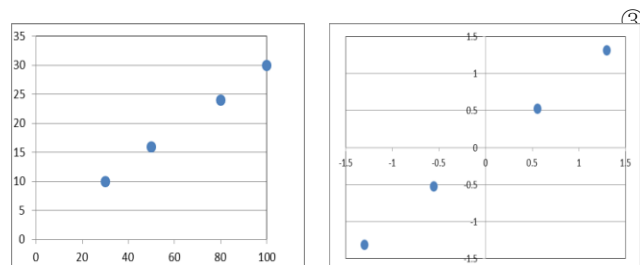
How do you solve? How do you teach?

第19号

発行責任者
盛岡第三高等学校
下町壽男

データの相関 ②

	テスト1	テスト2	テスト5	テスト6
学者うさぎ	100	30	1.30	1.31
マドンナうさぎ	80	24	0.56	0.53
元気ネコ	50	16	-0.56	-0.53
うかれぎつね	30	10	-1.30	-1.31
平均	65	20		



テスト1 vs テスト2

テスト5 vs テスト6

前回は、テスト1とテスト2の相関を考える代わりに、テスト5とテスト6の相関を考えても良さそうだという話までいきました。

テスト5とテスト6は、テスト1、テスト2をそれぞれ「平均0、標準偏差1」に平行移動し、拡大縮小したものであることに注意してください。このようなデータの変換を「標準化」または「正規化」する、と呼ぶことがあります。

ここで、少し脱線します。偏差値について考えてみたいと思います。

例えば、模試で数学の得点が0点だったのに、偏差値が30もついてくることがあります。この、偏差値の計算のメカニズムはどうなっているのか、おさらいしておきましょう。

ある集団でテストを行ったとき、得点を X とし、それに対応する偏差値を Z とします。全体の平均点を \bar{X} 、標準偏差を σ とおけば

$$Z = \frac{X - \bar{X}}{\sigma} \times 10 + 50 \quad \text{と表されます。}$$

これが偏差値です。

この式の意味を考えてみましょう。

① $X - \bar{X}$ ……平均点を0点にする

② $\frac{X - \bar{X}}{\sigma}$ ……標準偏差を1にする

$\frac{X - \bar{X}}{\sigma} \times 10$ ……標準偏差を10にする

$\frac{X - \bar{X}}{\sigma} \times 10 + 50$ ……平均を50にする

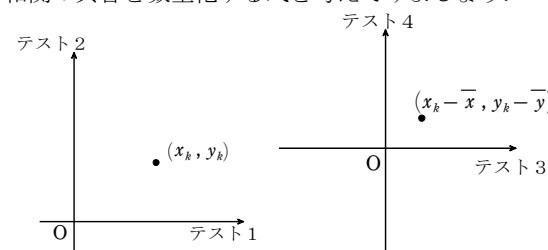
つまり、偏差値とは、得点の分布を、平均50、標準偏差10になるように平行移動したり、拡大縮小した値にすぎないということです。

このような処理を日本中で行えば、他と比較して自分の位置がどのくらいのところにあるかわかる一つの指標になるというわけですね。

話をもどします。相関について、次のようにまとめておきましょう。

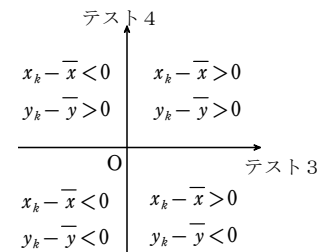
変数 x, y の相関を考えることは、各変数を平均0、標準偏差1になるように標準化したデータの相関を考えることと同じである。

では、これまで述べてきたテスト1とテスト2の相関の具合を数量化する式を考えてみましょう。



テスト1とテスト2に関するデータは、左下の図の点 (x_k, y_k) で表されます。一方、それぞれの平均点を引いた得点に変換した、テスト3とテスト4のデータは、平均点分それぞれ平行移動しているの、左下図の、 $(x_k - \bar{x}, y_k - \bar{y})$ という点に対応しています。

このとき、次のことに注意します。



上図で、第1象限は2つの変数がともに平均を超えている領域、第2象限はテスト3の平均は越えず、テスト4の平均は超えている領域、第3象限は、2つの変数がともに平均を下回っている領域、第4象限は、テスト3の平均は越え、テスト4の平均が下回っている領域ということになります。

このように、散布図は、2つの変数の平均のところを線を書いて、4分割して考えると分析しやすくなります。ここで、あるデータ (x_k, y_k) に対し、平均との差の積 $(x_k - \bar{x})(y_k - \bar{y})$ を考えると、

データが第1・第3象限にあるとき、積は正になり、第2・第4象限にあるとき、積は負になります。

第1・第3象限にデータがあるということは、正の相関が強いときであり、逆に、データが第2・第4象限にあるときは負の相関が強いときですね。

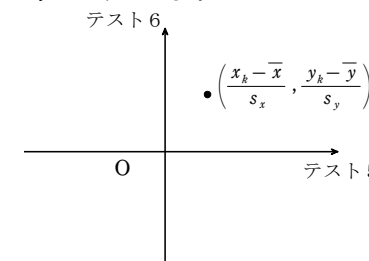
そこで、すべてのデータに対して、平均との差の積を加えて、その値が大きいほど、データが第1・第3象限にたくさんある、つまり正の相関が強いということがいえます。この考えを利用すると、相関の強さを数量で表すことができるわけです。平均との差の積の平均 $\frac{1}{n} \{ (x_1 - \bar{x})(y_1 - \bar{y}) + (x_2 - \bar{x})(y_2 - \bar{y}) + \dots + (x_n - \bar{x})(y_n - \bar{y}) \}$ を共分散といって、 s_{xy} という文字で表します。

まとめましょう。

$$s_{xy} > 0 \Rightarrow \text{①③にデータがある} \Rightarrow \text{正の相関が強い}$$

$$s_{xy} < 0 \Rightarrow \text{②④にデータがある} \Rightarrow \text{負の相関が強い}$$

さて、ではこの共分散を、テスト5とテスト6で考えてみましょう。



$$\frac{1}{n} \cdot \frac{1}{s_x s_y} \cdot \{ (x_1 - \bar{x})(y_1 - \bar{y}) + (x_2 - \bar{x})(y_2 - \bar{y}) + \dots + (x_n - \bar{x})(y_n - \bar{y}) \}$$

これは、 $\frac{s_{xy}}{s_x s_y}$ というシンプルな式になりますね。

この値を「相関係数」といい、 r で表します。

このように表すことによって、相関の強さを-1から1までの数で標準化できるので、イメージがしやすくなります。

相関係数をベクトルから眺めてみましょう。

$$\vec{a} = (x_1 - \bar{x}, x_2 - \bar{x}, x_3 - \bar{x}, \dots, x_n - \bar{x})$$

$$\vec{b} = (y_1 - \bar{y}, y_2 - \bar{y}, y_3 - \bar{y}, \dots, y_n - \bar{y})$$

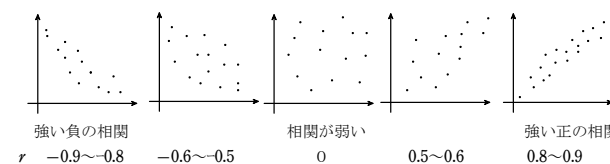
$$\text{とすると、} r = \frac{\vec{a} \cdot \vec{b}}{|\vec{a}| |\vec{b}|} \quad \text{ということです。}$$

つまり、2つのベクトルのなす角を θ としたとき、相関係数は2つのベクトルの $\cos \theta$ を表すものと考えることができます。

コサインは、2つベクトルの「近さ加減」を表す量なので、これを相関係数と考えても納得がいきますね。

基本的に、相関係数の算出はコンピュータ等にやらせればよくて、問題は、相関係数の見方です。

散布図と相関係数の関係を図に表しておきましょう。



これを基に、データの相関の強弱を判断できればよいと思います。大切なのは、相関や相関係数とは、データ分析の目を持つために重要な事柄であるということで、決してセンター試験のためにあるものではありません。