

実験や調査をして得られたデータの解析手法としてすぐに思い浮かぶのは、平均値や相関係数を求めるなどといった統計的な解析であろう。しかし、統計的な分析だけではなく、人間が見える形に変換してその傾向をつかむことも重要である。例えば、アンスコムの4つ組と呼ばれる有名な例がある。この例では、4つのデータが提示される。 $k$  番目のデータは次の書式で与えられる。

$$(a_{k,1}, b_{k,1}), (a_{k,2}, b_{k,2}), \dots, (a_{k,11}, b_{k,11}) \quad (1)$$

ここで、 $a_{k,n}, b_{k,n}$  は実数である。これら4つのデータは、統計的に分析して得られる平均値や相関係数といった値は全て（ほぼ）等しい。統計的な解析だけを頼りにすると、これら4つのデータは同質のデータであると言ってしまうかもしれない。しかし、これらのデータを散布図に直してみると、驚くほどその分布は異なっていることがわかる。

図1は、アンスコムの4つ組の第一のデータを散布図にしたものである。いかにも2つの変量の間に正の相関がありそうな散布図をしていることがわかる。我々が「2つのデータに正の相関がある」と聞いたとき、真っ先に連想するであろう散布図は、図1のようなものであろう。

図2は、アンスコムの4つ組の第二のデータを散布図にしたものである。確かに正の相関はあるかもしれないが、そのデータの分布は直線よりかは放物線の方が近い形状をしている。このデータだけを見ると、直線の近似をためらう人もいるかもしれない。

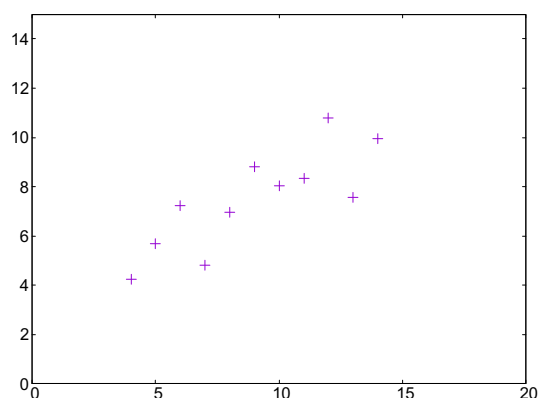


図1 アンスコムの4つ組の第一のデータを散布図にしたもの

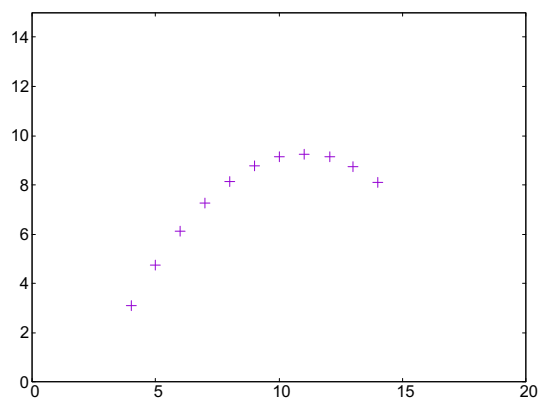


図2 アンスコムの4つ組の第二のデータを散布図にしたもの

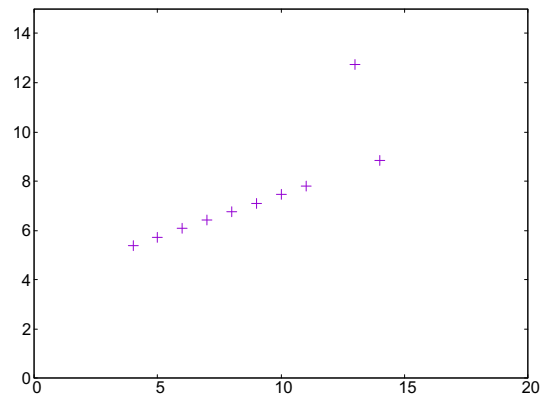


図3 アンスコムの4つ組の第三のデータを散布図にしたもの

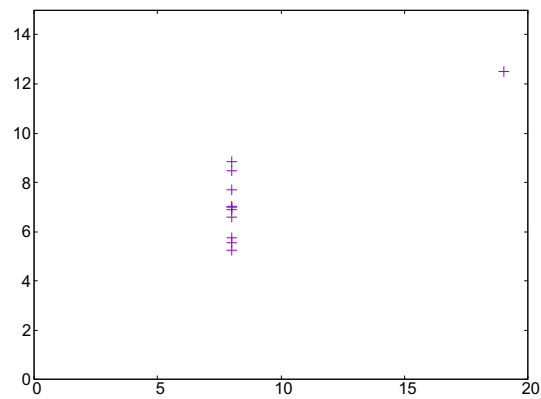


図4 アンスコムの4つ組の第四のデータを散布図にしたもの

図3は、アンスコムの4つ組の第三のデータを散布図にしたものである。確かに正の相関があると言えるデータだろうが、1つだけ外れ値があることが読み取れる。このような場合は、このままデータの解析を進めるのではなく、外れ値がなぜ得られたのかの考察をすべきだろう。

図4は、アンスコムの4つ組の第四のデータを散布図にしたものである。第三のデータと同じく1つだけ外れ値があるが、そのデータの分布は第三のデータとは大きく異なっている。このデータの場合、正の相関があるということは最早不適切であろう。