

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、人間が容易に特徴を把握できる形に変換してデータの特徴をつかむことも重要である。例えば、アンスコムの4つ組<sup>\*1</sup>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。各々のデータには、2つの実数の組が11個含まれている。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などの値は（ほぼ）等しいことがアンスコムによって示されている<sup>\*2</sup>。三末が指摘するように、これらの値だけに着目すると、4つのデータは同じ特徴を持っていると言ってしまうかもしれない<sup>\*3</sup>。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。

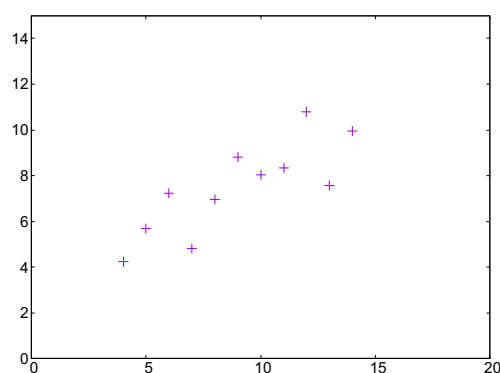


図1 アンスコムの4つ組の第一のデータを散布図にしたもの

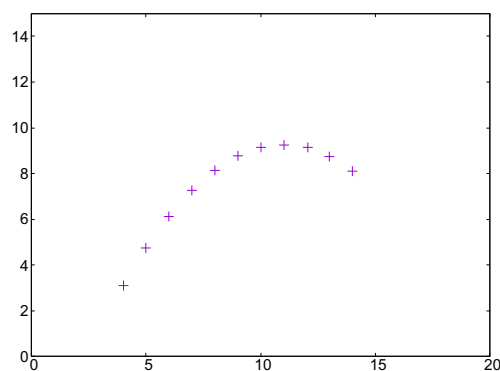


図2 アンスコムの4つ組の第二のデータを散布図にしたもの

<sup>\*1</sup> Anscombe, F. J. (1973). Graphs in Statistical Analysis. *The American Statistician*, 27(1), pp.17–21.

<sup>\*2</sup> 同上。

<sup>\*3</sup> 三末和男『情報可視化入門 人の視覚とデータの表現手法』、森北出版、2021年、p.4。

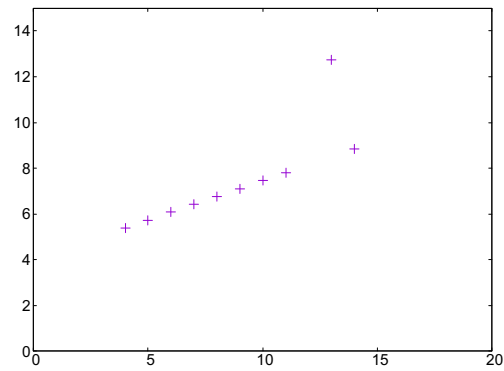


図 3 アンスコムの 4 つ組の第三のデータを散布図にしたもの

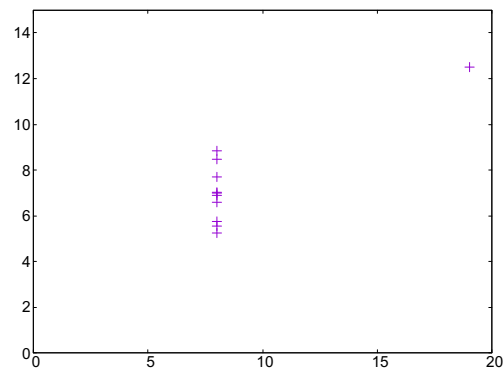


図 4 アンスコムの 4 つ組の第四のデータを散布図にしたもの

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。