

相関係数と順位相関係数について

相関係数と順位相関係数の使い分けについて考えてみる。今、図1と図2のような2種類のデータを考える。(相関と順位相関.txt)

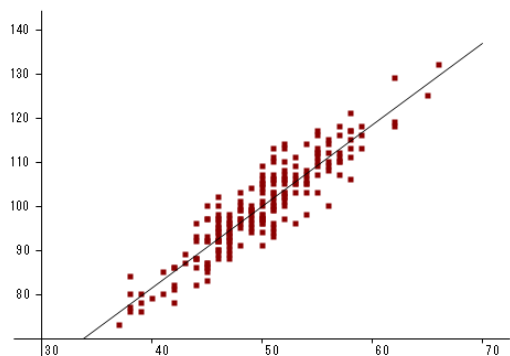


図1 データ1

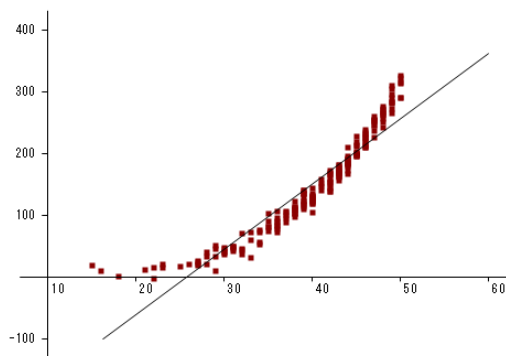


図2 データ2

この2つのデータには、 x と y に関係がありそうである。左のデータ1は2変量正規分布と呼ばれるデータで、回帰直線がよく関係を表している。それに対して右側のデータ2の関係は曲線を描いており、直線で表すには無理がある。

この2種類のデータについて2つの変数の正規性を調べてみると図3と図4のようになる。

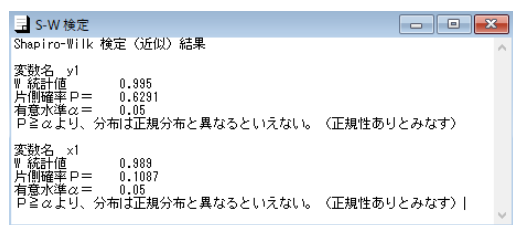


図3 データ1の正規性

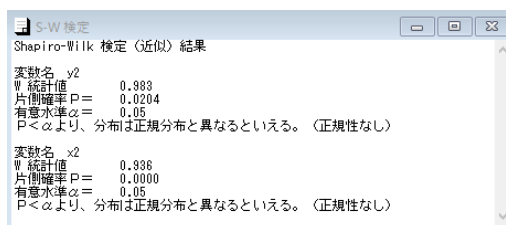


図4 データ2の正規性

データ1は2変数ともほぼ正規分布であり、データ2は2変数とも正規分布といえない(もちろん一方が正規分布する場合もある)。

次に2つのデータの相関係数と順位相関係数を求めてみよう。結果は図5と図6の通りである。

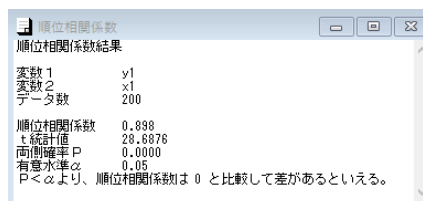
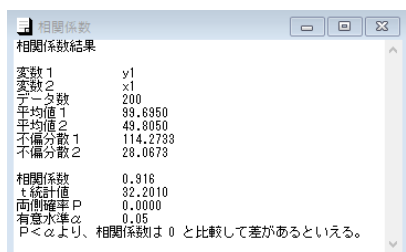


図5 データ1の相関係数と順位相関係数

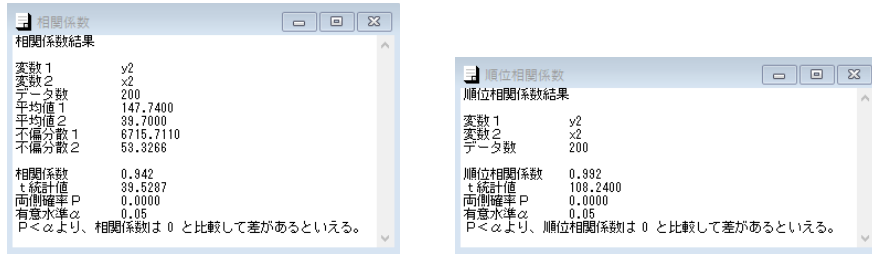


図6 データ2の相関係数と順位相関係数

図5のデータ1の場合、相関係数(0.916) > 順位相関係数(0.898)となり、どちらかと言うと相関係数が関係を表すのにふさわしいようである。これに対して図6のデータ2の場合、相関係数(0.942) < 順位相関係数(0.992)となり、順位相関係数がより強く関係を表している。散布図のデータ点の並びが曲線に沿うような場合は、2つの変数に正規性が見られないことがあり、順位相関係数を利用の方が関係を強調できる場合がある。

特に図7のデータ3のように散布点が単調増加関数上に並ぶ場合、図8のように順位相関係数は誤差の範囲で1になる。

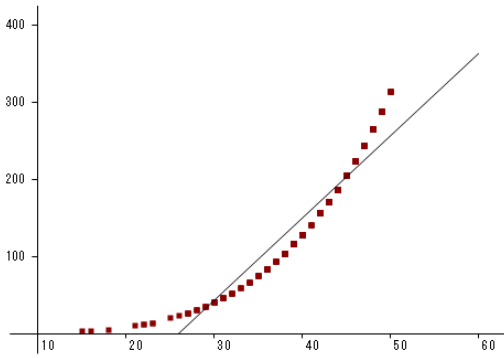


図7 データ3

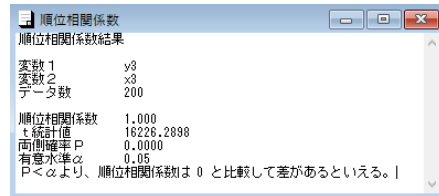


図8 データ3の順位相関係数

最後に一つ注意を与えておく。ここでは2つの変数を別々に分けて(周辺分布の)正規性を調べていることである。本来相関係数は2変量正規分布を基にしているので、直接2変量正規分布に従うかどうか調べるのが正確な方法である。そのため、ここでの議論は簡易的方法と認識しておく必要がある。