

統計的仮説検定 その1

1 t 検定

2つの平均値に有意な差が存在するのか、という疑問については古くから統計的な処理の必要な多くの問題が指摘されていた。例えば、従来の飼料と新しく開発した飼料で家畜の成長が変わるのか、従来の薬剤と新しい薬剤とで治療成果が異なるのか、従来の製造法と新しい製造法で材料の剛性が変わるのか、など、幅広い分野で検討がなされてきた。19世紀のころ、統計的に扱うためには標本の数が多いほど正確である、という常識が作業において支配的であった。生物学の分野では標本が多く集められるという観点からショウジョウバエやマウスが今でも多く使われるのは周知のことであろう。

しかし、工業分野、特に、製品の出荷前検査などでは多くの製品を標本として確保することは生産性の低下にもつながりかねない重要な問題である。当時アイルランドのギネスビール社で醸造と大麦の生産に統計的な手法を適用することを業務としていたオックスフォード大学出身の化学者・数学者であったウィリアム・ゴセットは小標本の問題に着目し、正規分布としての母数を推定できないような少ない標本でも使用可能な扱いを考案し、論文として発表した。しかしながら、ギネスビール社では従業員の論文発表を認めていなかったため、Student というペンネームで論文を投稿したこと、この小標本を想定した分布はそのペンネームから Student の t 分布と呼ばれ、現在でも t 分布として知られている。ここでは、 t 分布を用いた仮説検定を紹介する。

1.1 t 分布

標本数を無数に多く用意することができると、正規分布を仮定した様々な操作が可能となるが、母集団の性質がよくわかっていないものに対して少数の標本からその分布を正規分布とみなして母数を求めることには危険が伴う。そこで、少数の標本から母数を推定するときに t 分布というものを用いる。詳細は後述するが、中心極限定理の時に説明した以下の式が関係する。

$$u = \frac{\bar{x} - m}{\sigma/\sqrt{n}} \quad (1)$$

上式は、測定値の分布が正規分布に従わない時でもその平均値の分布は正規分布になるというものであった。しかし、母標準偏差は通常知ることが難しいことが多い。そこで、母標準偏差の代わりに標準偏差（不偏分散の平方根） v を用いて、以下のように書き換える。

$$t = \frac{\bar{x} - m}{v/\sqrt{n}} = \frac{\bar{x} - m}{s/\sqrt{n-1}} \quad (2)$$

ここで、 \bar{x} は変数 x の平均値、 m は母平均、 v は不偏分散の平方根（標準偏差）、そして、 n は標本の個数である。この t という統計量は、不偏分散を用いていることからわかるように、標本を採るたびに変動する。そこで、あらたに分布を形成するが、この分布を自由度 $\phi = n - 1$ の t 分布という。図1に自由度 $\phi = 3$ のときの t 分布と標準正規分布を比較して示す。 t 分布は自由度が大きくなるにつれて標準正規分布に近づいていく。 t の値は自由度 ϕ によって変化するので、必ず自由度とセットで扱うことに注意が必要である。

1.2 t 分布表

自由度 ϕ の t 分布について、ある $t(\phi, \alpha) > 0$ の値に対し、 $\pm t(\phi, \alpha)$ の外側に t が現れる確率 α について、 α と $t(\phi, \alpha)$ の関係を与える t 分布表が作られている。今回配布したこの t 分布表も今後試験も含めて使用していくので、書き込みなどをしないように注意すること。

式からわかるように変数 t は標本の個数によって変化する量であるので、この t の値をまとめた t 分布表では標本の個数から 1 を引いた自由度に対して、両端の確率 α に該当するときの境界の値が示されている。

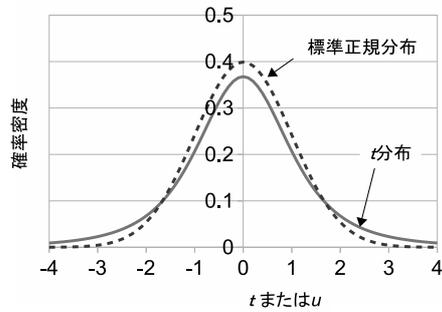


図 1: 標準正規分布と t 分布 (自由度 $\phi = 3$) の比較

表 1: t 検定サンプル

番号	x	y	x^2	y^2
1	2.5	2.9	6.25	8.41
2	2.6	3.1	6.76	9.61
3	2.1	2.8	4.41	7.84
4	3.1	3.6	9.61	12.96
5	2.7	3.3	7.29	10.89
6	2.4	3.4	5.76	11.56
和	15.4	19.1	40.08	61.27
平均値	2.567	3.183		
S_x	0.5533			
S_y	0.4683			

1.3 実際の t 検定

さて、実際に 2 つの平均値を比較する作業であるが、以下に示す表 1 のデータをもとに作業を紹介する。

表 1 において、変数 x と y の平均値はそれぞれ 2.567 と 3.183 となっており異なっているが、それが意味のある差であるのか、たまたま偶然生じたものであるのかを調べる。

まず、式 (2) の t を 2 つの変数に対応するもの書き換える。その際に問題となるのは分散の決め方となるが、両方が同じ分散に従うと仮定できる場合には以下の式 (3) で共分散 $\hat{\sigma}^2$ を求めることができる。

$$\hat{\sigma}^2 = \frac{S_x + S_y}{n_x + n_y - 2} \cdot \left(\frac{1}{n_x} + \frac{1}{n_y} \right) \quad (3)$$

ここで、 n_x および n_y はそれぞれ変数 x と y のデータの個数である。これを使うと、平均値の差に対応する t の値 t_0 を以下の式で求めることができる。

$$t_0 = \frac{|\bar{x} - \bar{y}|}{\sqrt{\frac{S_x + S_y}{n_x + n_y - 2} \cdot \left(\frac{1}{n_x} + \frac{1}{n_y} \right)}} \quad (4)$$

表 1 のデータに対して t_0 を求めると、3.342 ということになる。この場合の自由度 ϕ は $n_x + n_y - 2 = 10$ であるので、 t 分布表の自由度 10 の行を見ると、

ϕ	0.10	0.05	0.02	0.01	0.001
10	1.8125	2.2281	2.7638	3.1693	4.5869

となっている。今求めた t_0 は危険率 0.01 と 0.001 の間に入る値となっている。そこで、以下のように検定を進める。

- 帰無仮説の決定: 2 つの平均値 \bar{x} と \bar{y} に差は無い。(同じ母集団に属する標本と言える。)

- 対立仮説：2つの平均値 \bar{x} と \bar{y} には差がある。
- t_0 を求める。(今回は $t_0 = 3.342$)
- t 分布表の危険率を確認する。(自由度 $\phi = 10$ の行において, $t_0 > t(10, 0.01)$ が確認された.)
- 帰無仮説を棄却する際に, 棄却が間違っている確率(危険率)は1%未満であることが確認された。
- よって, 帰無仮説を棄却することが可能である。(対立仮説の採用)
- 結論：2つの平均値 \bar{x} と \bar{y} には有意な差がある。($p < .01$)

実際にここまでの作業で確認したことは, 図2に示すように2つの平均値 \bar{x} と \bar{y} の差の t 値を求めたところ, 右端の確率1%未満の領域に含まれるので, これらを異なる母集団に属する標本とみなしても問題は無いということである。

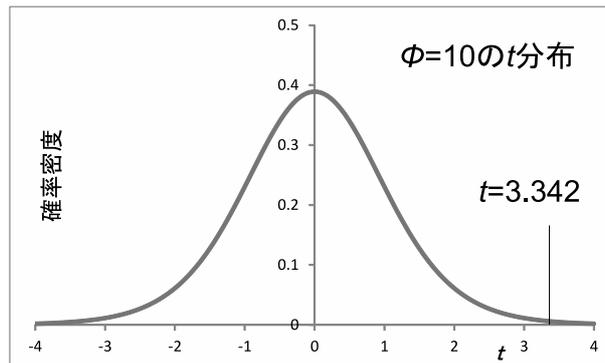


図 2: $\phi = 10$ の t 分布

なお, 2つの平均値の分散が同じとみなせない時には Welch の方法など他の t 検定の手法を用いる必要があるが, ここでは省略する。

【演習】資料にそって t 検定の演習を行う。