

### III-2-6. Student の t 分布

t 分布は student の t 検定に使われる確率分布です。Student の t 検定は、データの差の有意性の検定です。ですから、t 分布とは正規分布すると仮定されるあるデータと、あるデータの平均値の差の分布のことです。この確率分布は正規分布に似ていて、実際自由度が十分大きければ正規分布に近似できます。左右対称というところも正規分布に似ています。正規分布と大きく違うところは、自由度によって分布が変わるところです。この点は  $\chi^2$  分布です。というのも、この確率分布が正規分布と  $\chi^2$  分布の合成によってできているからです。

これを発見したのはビール会社ギネスの技師だった W.Gosset ですが、Student の t 分布という名前がついているのは、会社から論文投稿を禁じられていたために、彼が Student という筆名でその発見を論文化したからです。おそらく、彼の発見の動機は、今、私たちが抱いている疑問と同じだったろうと思います。つまり、「正規分布することが想定されるデータについて、データから得られた平均値がどのくらい正規分布の期待値（平均値・中央値）に近いのかは、母集団の標準偏差( $\sigma$ )を尺度にして、標準化して、標準正規分布  $N(0,1)$  の分布の中でデータの平均値がどの位置にあるのかを考えれば良いというのはわかるにしても、そもそも母集団の標準偏差( $\sigma$ )を知らないのだから、標準化することができません。データから得られるのは標本分散から推定した母分散の推定値だから確率的に変動する。その変動をどう読みこむのか。」という疑問です。

我々の疑問を数式的に表し、それをどのように解決すればよいかを考えます。まず、我々は正規分布というものの存在を認めています。u が標準正規分布  $N(0,1)$  に従うのならば、その確率は

$$N_{(0,1)} = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{u^2}{2}}$$

次は u をどう求めるかですが、その作業が標準化と言われる作業で、実際のデータ、期待値、分散（偏差）が使われます。これは習った通り。

$$u = \frac{x - \mu}{\sigma}$$

なのですが、我々が疑問としているのは、「この式で使っているのは正規分布している本来の理想的なデータの期待値=平均値( $\mu$ )と分散 ( $\sigma^2$ ) なのだから、実際のデータしか知らない私たちが、そんなものは知るわけがないだろう。」ということです。仮に、母集団の平均を知っていることにして、実際に得られた値と真の平均値との差を考えられる偏差を基準値としてその平方和を求めると、

$SS = v$ として、

$$v = \sum_{i=1}^m \left( \frac{y_i - \mu}{\sigma_{se}} \right)^2$$

です。 $\sigma_{se}$ と添え字を付けたのは、まだ、 $y_i$ が何だかわからないからです。とりあえず、

something estimated だと覚えておいて下さい。そこれをよく見ると、 $\sigma_{se}$ が外から与えた期待値ならば、第一の式の総和記号の中は $\chi^2$ の定義式で、 $v$ は $\chi^2$ 分布します。このこの場合、平均値という形で合計の値が縛られている のだから、 $\chi^2$ の値を-1 個目まで決めた後の最後の値は自動的に決まります。だから全体の自由度は  $m-1$  です。つまり  $v$ は自由度  $m-1$  の $\chi^2$ 分布します。この $u$ と $v$ は互いに独立なのです（関係ない。 $u$ が変化しても $v$ は変化しない。反対に $v$ が変化して $u$ の値に変化はない。という意味です。これが W.Gosset のやった最大の発見なのです。だから、2つの変数を合成した変数を作り、その変数の確率を2つの確率の積として計算できるのです。）。私は、この2つが独立かと問われると、一瞬、言葉に詰まってわからなくなります。でも落ちつて考えると、確かに独立ですね。そこで、2つの関数を関係づけるために、この2つの変数からできる合成関数を考えます。母集団の分散  $\sigma$  は平方和を $(n-1)$ で割ったものが推定になるというのをやりましたが。この場合は、 $\mu$ と $\sigma_{se}$ を既知のものとして、外から与えているので、自由度は  $m$  です。

$$v = \sum_{i=1}^m \left( \frac{x_i - \mu}{\sigma_{se}} \right)^2 = \frac{1}{\sigma_{se}^2} \sum_{i=1}^m (x_i - \mu)^2$$

$$s_{se}^2 = \frac{1}{n} \sum_{i=1}^m (x_i - \mu)^2$$

だから

$$s^2 = \frac{1}{n} v \sigma_{se}^2$$

$$s = \sqrt{\frac{v}{n}} \sigma_{se}$$

$$\frac{s}{\sigma_{se}} = \sqrt{\frac{v}{n}}$$

$u = x - \mu \sigma$  は正規分布しますが、問題は、分母つまり単位となる物差しの長さが、データからもとめた偏差と真の偏差との間で違っているということです。その解決として「実際のデータから求めた値に、実際の値から得た偏差と真の偏差の比を掛けたものを、合成変数として、その値になる確率を考える。」というのが W.Gosset の提案内容です。

$v$ は $\chi^2$ 分布し、 $u$ は正規分布して、互いに独立ですから、2つを合成した関数変数の確率は両者の確率の積です。

具体的には、合成変数 ( $t$ : $t$  値) は次の式です。

$$t = \frac{y - \mu}{\sigma_{se}} \cdot \frac{\sigma}{s}$$

$$t = u \sqrt{\frac{n}{v}}$$

確率変数 $u$ と $v$ が同時にそれぞれ独立にある値をとって、 $t$ の値が決まるので、 $u$ になる確率

W(u)とvになる確率P(v)の積がtになる確率S(t)です。

$$S(t) = W(u)P(v)$$

確率の総和は次の重積分で計算できて、その値は1です（確率の総和は1）。

$$\int_{-\infty}^{\infty} \int_0^{\infty} W(u)P(v) dudv = 1$$

これをtとsの式に変換します。その際、 $t = u\sqrt{\frac{n}{v}}$ からヤコビアンをもとめます。ヤコビアンとは、正規分変数の変換によって生ずる。係数のようなものです。ブログに解説を載せておきました。わからなければ読んでください。

$$\begin{aligned} & \int_{-\infty}^{\infty} \int_0^{\infty} W(u)P(v) dudv \\ &= \int_{-\infty}^{\infty} \int_0^{\infty} J(u, v/t, s)W(u)P(v) dt ds \end{aligned}$$

vからsの変換ですが、sはvそのものでも良さそうだから、 $s = v$ として、

$$\frac{du}{dt} = \sqrt{\frac{v}{n}}, \quad \frac{du}{ds} = 0, \quad \frac{dv}{dt} = \frac{-2nu^2}{t^{-3}}, \quad \frac{dv}{ds} = 1$$

$$J(u, v/t, s) = \begin{vmatrix} \sqrt{\frac{v}{n}} & 0 \\ -\frac{2nu^2}{t^{-3}} & 1 \end{vmatrix} = \sqrt{\frac{v}{n}}$$

$$W(u) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{u^2}{2}}$$

$$P(v) = P\left(\chi^2_{\varphi}\right) = \frac{\chi^2_{\varphi}^{\frac{(\varphi-1)}{2}} e^{-\frac{\chi^2_{\varphi}}{2}}}{2^{\frac{\varphi}{2}} \Gamma\left(\frac{\varphi}{2}\right)}$$

なのですが、今、考えているのは、自由度は外から与えた平均値と分散を与えた時の確率分布ですから、自由度は $m-1$ ではなくて $m$ なのです。W.Gossetが、考えたt分布は、2つの群の平均値の差が、0という帰無仮説を立てて、実測された二つの群の差からt値を計算し、帰無仮説から導かれるt値=0という期待値に対して、実測されたt値がどのような確率で存在するかという確率分布です。だから、自由度は1です。

$$P(v) = \frac{\frac{-1}{v^{\frac{1}{2}}} e^{-\frac{v^2}{2}}}{2^{\frac{1}{2}} \Gamma\left(\frac{1}{2}\right)}$$

$\Gamma()$ はガンマ関数（ガンマ関数についてはカイ二乗分布を作るときに説明しました。これらを使って、以下の重積分を変形します。

$$\int_{-\infty}^{\infty} \int_0^{\infty} W(u)P(v) dudv$$

$$\int_{-\infty}^{\infty} \int_0^{\infty} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{u^2}{2}} v^{\frac{-1}{2}} e^{-\frac{v^2}{2}} \frac{1}{2^{\frac{1}{2}}\Gamma(\frac{1}{2})} dudv$$

$$\int_{-\infty}^{\infty} \int_0^{\infty} \frac{\sqrt{v}}{\sqrt{n}\sqrt{2\pi}\sigma} e^{-\frac{u^2}{2}} v^{\frac{-1}{2}} e^{-\frac{v^2}{2}} \frac{1}{2^{\frac{1}{2}}\Gamma(\frac{1}{2})} dsdt$$

$$\int_{-\infty}^{\infty} \int_0^{\infty} \frac{1}{\sqrt{2n\pi}\sigma} e^{-\frac{u^2}{2}} \frac{e^{-\frac{v^2}{2}}}{2^{\frac{1}{2}}\Gamma(\frac{1}{2})} dsdt$$

かなり長い計算になるので、省略します。興味があれば、ブログを読んでください。少し誤りがあるので、近々、修正しておきます。

最終的に、

$$\int_{-\infty}^{\infty} \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi}\Gamma(\frac{n}{2})\left(\frac{t^2}{n^2}+1\right)^{\frac{n}{2}+1}} dt = 1$$

となりますが

個々の t の値をとる確率は、 $-\infty$  から t の値までの定積分を微分すれば良いので、積分記号の中の間数になって、

$$S(t) = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi}\Gamma(\frac{n}{2})\left(\frac{t^2}{n^2}+1\right)^{\frac{n}{2}+1}}$$

です。この複雑な式は、 $\beta$ 関数を使えば、簡略化して表現できます。 $\beta$ 関数はガンマ関数の積で以下の式で表せます。

$$\beta(p, q) = \frac{\Gamma(p)\Gamma(q)}{\Gamma(p+q)} \quad \text{式 22}$$

$\Gamma(\frac{1}{2}) = \sqrt{\pi}$  ですから

$$\beta\left(\frac{n}{2}, \frac{1}{2}\right) = \frac{\Gamma(\frac{n}{2})\Gamma(\frac{1}{2})}{\Gamma(\frac{n}{2} + \frac{1}{2})} = \frac{\pi\Gamma(\frac{n}{2})}{\Gamma(\frac{n}{2} + \frac{1}{2})}$$

で

$$S(t) = \frac{1}{\sqrt{n\pi}\beta\left(\frac{t^2}{n^2}+1\right)^{\frac{n}{2}+1}} \quad \text{式 23}$$

となります。一般に student の t 値の確率密度関数として紹介されているのは、この式です。

話を元に戻します。

$$v = \sum_{i=1}^m \left( \frac{y_i - \mu}{\sigma_{se}} \right)^2$$

のところで、 $y_i$ は何だかわからないとしましたが、今はわかります。 $y_i$ は t 値です。 $\mu$ はその期待値で0です。つまり、グループ A とグループ B の平均値の差です。 $\mu$ は0なんだから、 $\left( \frac{y_i - \mu}{\sigma_{se}} \right) = \left( \frac{y_A - y_B}{\sigma_{se}} \right)$ です。

A の平均値 $y_A$ が期待値0だとすれば、 $y_B$ は期待値0の周辺で確率的に変動するでしょう。

また、B の平均値 $y_B$ が期待値0だとすれば、 $y_A$ は期待値0の周辺で確率的に変動するでしょう。これは、真の平均値がデータから得られた平均値の周りで変動するときの積率でもあります。データから推定された平均値 $m$ が真の平均値の周りで変動するときの積率を平均誤差と言うということは、すでに説明しました。平均誤差は $\frac{\sigma}{\sqrt{n}}$ です。つまり、2つの群の

平均値から推定した、t 値の分散は $\frac{\sigma^2}{n}$ で偏差は $\frac{\sigma}{\sqrt{n}}$ です。 $\sigma_{se} = \frac{\sigma}{\sqrt{n}}$ つまり、 $\sigma_{se}$ (Something estimated)は S tandard error だという話です。t 値の定義式を書き換えておきます。

$$t = \frac{y - \mu}{\sigma_{se}} \cdot \frac{\sigma}{s} = \frac{y - \mu}{\frac{\sigma}{\sqrt{n}}} \cdot \frac{\sigma}{s} = \frac{y - \mu}{\frac{s}{\sqrt{n}}} \quad \text{式 24}$$

sはサンプル集団の偏差

ポイントを要約すると、

t 分布は、正規分布する分散の等しい二つのデータ群の平均値の差から、t 値を作って、その値がカイ二乗分布するとして、求めた確率分布で、データ群の差を標準化するために、データ群から求めた標準誤差で差を割ったもので。