

IV-3. 統計検定

IV-3-1. Student の t 検定

t 検定とは、2組のデータがあるときに、そのデータの平均値に差があると言えるか言えないかという検定です。確率論的に言えば、平均値が同じ母集団からとられたデータである確率を検定します。

この検定が想定しているモデルを理解するために、ある値 (x) があり、その値があらかじめ予測される平均値の候補として妥当であるかを検討してみましょう。

データから計算された平均値を M 、母集団の平均値を μ 、母集団の標準偏差を σ とすると、それぞれの値の関係は以下のように表せます。

$$x = \mu \pm \alpha\sigma$$

α は、 σ を単位とした時のどのくらい離れているかを示す値で、具体的なデータがあれば、

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

という計算式で求められます。

この値を標準正規分布に当てはめて、確率水準によって決まる値より大きいか小さいかを検討し、大きければ、その値は予測される値の候補としては不適切であると結論します。たとえば、95%の信頼限界（危険率 $p=0.05$ と表すほうが一般的）で検討するのであれば、 α の値が 1.96 より大きければ、その値を候補として採用することは不適切と判断します。確かにそうですが、私たちが知っているのはデータの平均 (M) と標本集団から計算された母集団の分散の推測値 σ'^2 であって真の平均値 μ 、母集団の分散 σ^2 ではありません。 μ も σ も知りようがないから計算できません。標本サイズを大きくすれば、標本集団の平均値と分散は母集団の平均値と分散に近づくから、問題ないだろうという考え方はあります。実際、記述統計学の大成者である **Karl Pearson** はそう考えていました。実務家にとっては、実験の繰り返し数を増やすには、場所や費用の制約があります。小標本で統計学的な確率を論じたいというのは、実務家として当然の願望です。ギネスの醸造所で働いていた **William Gosset** も小標本による統計的な判定に強い興味を持ちました。ギネスでは技術者が論文を投稿することを禁止していましたから、この問題についての論文を、ピアソンの主催している **Biometrika** という雑誌に投稿しました。この時に使っていたペンネームが **Student** です。これが分散分析の創始者である **Ronald Fisher** に認められて、**Fisher** によって **Student** の t 分布として整理されます。

Gosset の考え方を整理します。

まず分散についてです。分散の値も確率的に変動するのだから、正規分布ではなくて、分散を確率変数としている χ^2 分布と正規分布を組み合わせた t 分布を確率モデルとして使え

ば良い。これが一つです。

次に μ ですが。t 検定では、平均値に差があるかないかを検定しています。つまり、平均値の差が0であることが帰無仮説ですから、 $\mu = 0$ です。

もう一つ重要なことは、平均値であるということです。データから予測された平均値は、母集団の期待値（平均値）の周りに標準誤差を偏差として分散しているはずですが、標準偏差 σ ではなくて標準誤差 $\frac{\sigma}{\sqrt{n}}$ を差の平均値Mから0の距離の単位にしなければなりません。

$$M = \mu \pm t \frac{\sigma}{\sqrt{n}}$$

t は、 $\frac{\sigma}{\sqrt{n}}$ を単位とした時のどのくらい離れているかを示す値。

これを変形して

$$t = \frac{\mu - M}{\frac{\sigma}{\sqrt{n}}}$$

$\mu = 0$ を仮定しているから

$$t = \frac{M}{\frac{\sigma}{\sqrt{n}}}$$

基本的な考え方は以上です。Gosset のオリジナルはもう少し違う形で書かれています。この形に整理したのは Fisher です。シンプルですが、これでは具体的にどうすればよいのかわかりませんね。

実際の検定では、 σ や n をどのように計算すればよいかを説明します。2つの場合が考えられます。一つは対になったデータというもので、いくつもの鉢があるときに、同じ鉢に2つの異なる植物を植えて成長を比較するというように、比較する相手が1対1に決まっているデータです。もう一つは、鉢を2つのグループに分けて、それぞれのグループの鉢には同じ種類の植物を植えて、グループ間で比較するというものです。初めの例を対になったデータのt 検定と呼び、普通のt 検定と区別します。分散の計算の考え方は、「4-1. データの扱い（情報の取り出し）」を読むと、理解できます。

対になった t 検定 (paired t test)

この検定では、対になったデータのどちらかが相手に比べ大きいか、データに差があるといつてよいかということを検定します。

対になったデータを比べるとは、たとえば、右手と左手とどちらが長いとか、手と足とどちらが長いとか、一つの鉢に違う植物をうえて、これを繰り返してどちらの植物の成長が早いかなどをくらべることをイメージすればよいでしょう。対になっているのだから、1番目の人の右手にはそれと対になった左手があります。このデータを A_1 と B_1 と表現すると、以下 A_2, B_2, A_3, B_3 のように表現できます。必ず一組になっているのだから、 C_k として $A_k - B_k$ を求めることができます。 C_k も標本から得られるデータで、確率的に変

動します。ところで、AとBに差がないということはCの平均が0ということです。つまり、帰無仮説は $\bar{C}=0$ 、差がないということです。

この場合は、データ対の数をデータ数 n をとって、差の値の平均 M 、分散 σ を計算し、次のように t の値を求めて、危険率 p を適切に決めて、自由度 $n-1$ として t 分布表を使って検定します。

$$t = \frac{M}{\frac{\sigma}{\sqrt{n}}}$$

式 42

対になっていないデータの t 検定

上記の知識を使って、対になっていないデータの t 検定の方法を考えます。 t 検定の考え方の基本は、対になった t 検定の場合と同じです。技術的に克服しなければならないポイントは差の分散をどのようにして求めるかです。

これが求まれば、式 51 を使って

$$x = \mu \pm t \frac{\sigma_{A-B}}{\sqrt{N}}$$

として

$$x = M_A \quad \mu = M_B$$

を代入して

$$t = \frac{M_A - M_B}{\frac{\sigma_{A-B}}{\sqrt{N}}}$$

を求め、この大小を検討すればよいでしょう。 M_A M_B は、それぞれ、A、B、の平均値です。ここで問題が2つあります。 σ_{A-B} と N です。 N は式 42 では n とあらわされていますが、A群のデータの数ともB群のデータの数とも違います。新たに N を決めなくてはなりません。A群のデータ数を m 、B群のデータ数を n 、 $N = \bar{n}$ とします。

$$\frac{\sigma_{A-B}}{\sqrt{\bar{n}}} = S$$

と表します。 σ_{A-B}^2 はA、B共通の分散ですが、これをA、Bの分散(σ_A^2, σ_B^2)から求めなければなりません。

ここでわれわれが証明しようとしている帰無仮説は、A,B両群が同じ母集団からとられたということです。つまり、

$$\sigma_A^2 = \sigma_B^2$$

もし、これが成り立つならば、重みをつけようが軽みをつけようがその平均である σ_{A-B} は、

$$\sigma_{A-B}^2 = \sigma_A^2 = \sigma_B^2$$

です。

すでに説明したように

$$E(M_A^2) = \frac{1}{m} \sigma_A^2$$

つまり、A群の母集団の平均値の周りに、A群のデータによって推定された平均値が $\frac{1}{m} \sigma_A^2$ の広がり分布しています。

$$E(M_A^2) = \frac{1}{m} \sigma_A^2 = \frac{1}{m} \sigma_{A-B}^2$$

$$E(M_B^2) = \frac{1}{n} \sigma_B^2 = \frac{1}{n} \sigma_{A+B}^2$$

$$SS = E(M_{A-B}^2) = E(M_A^2) + E(M_B^2) = \frac{1}{m} \sigma_A^2 + \frac{1}{n} \sigma_B^2$$

だから

$$SS = \frac{1}{m} \sigma_{A-B}^2 + \frac{1}{n} \sigma_{A-B}^2 = \left(\frac{1}{m} + \frac{1}{n} \right) \sigma_{A-B}^2$$

$$S = \sigma_{A-B} \sqrt{\frac{1}{m} + \frac{1}{n}}$$

これを $\frac{\sigma_{A-B}}{\sqrt{\bar{n}}} = S$ に入れると

$$\frac{\sigma_{A-B}}{\sqrt{\bar{n}}} = \sigma_{A-B} \sqrt{\frac{1}{m} + \frac{1}{n}}$$

$$\bar{n} = \frac{1}{\frac{1}{m} + \frac{1}{n}} = \frac{mn}{m+n}$$

$$\bar{n} = \left(\frac{mn}{m+n} \right)$$

求める t 値は

$$t = \frac{M_A - M_B}{\sigma_{A-B} \sqrt{\frac{1}{m} + \frac{1}{n}}}$$

式 43

となります。

また、「IV-2-2. 和の分散、差の分散」で説明したように、A・B の分散は、A+B の分散と同じで、A、B それぞれの分散を自由度で重みをつけた平均（下式）になります。

$$\sigma_{A-B}^2 = \sigma_{A+B}^2 = \frac{(m-1)\sigma_A^2 + (n-1)\sigma_B^2}{m+n-2}$$

以上で対になっていないデータの t 検定ができます。IV-2-2 の差の分散の説明でつかった表 13 の具体的なデータを使って実際の計算例を示します。

表 13.具体的なデータ(再掲)

	A	B
	1	1
	5	5
	6	6
		8
平均	4	5
SS	14	26

A : データ数 $m=3$ 、平均 $M_A=4$ $SS_A = 14$

B : データ数 $n=4$ 、平均 $M_B=5$ $SS_B = 26$ 分散 $\sigma^2 = 8,66667$

2つをあわせた自由度は $7 - 2 = 5$

$SS_{A+B} = SS_A + SS_B = 14 + 26 = 40$

$\sigma_{A+B}^2 = 40/5 = 8$

$\sigma_{A+B} = 2\sqrt{2} = 2.8284$

$\sqrt{\frac{1}{m} + \frac{1}{n}} = \sqrt{\frac{1}{3} + \frac{1}{4}} = 0.763763$

$M_A - M_B$ の絶対値 1

$z = 1 / (2,8284 \times 0.763763) = 0.46291$

t 表によれば $p = 0.05$ (5%の危険率) での t の臨海値は 2.571 ですから、A の母集団と B の母集団の平均値が異なっていると結論することはできません。