

農場で使える統計構成 ~ 第7回 図問題だるの1 ~

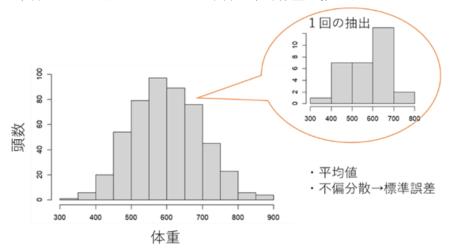
かやの

前回は平均値のバラツキをとらえる標準誤差などについて説明しました。たくさんの頭数がいる中で、その一部をサンプルして乳量や病原菌の有無などをチェックし「推定」するわけですが、そこから出した平均値なども当然のことながらバラツクわけです。そのバラツキを表したのが標準誤差でした。今回は1点という平均値(点推定)から一歩踏み込み、AからBまでというような区間推定について説明します。

区間推定(信頼区間)に必要な材料

区間推定といわれてピンと来ない人は、信頼区間を思い浮かべてください。95%信頼区間などといわれるアレです。信頼区間の推定は、調べたい対象がなんらかの分布に従うと仮定して計算するものから、ベイズ推定と呼ばれるものまで様々な手法があります。このシリーズでは、よりシンプルな方法として(一例として)t分布および正規分布を利用した区間推定をご紹介します。今回はt分布を利用して、区間推定(信頼区間)を推定しようと思います!

前回、500頭の牛群から30頭をサンプルして、群の平均体重を推定するということをやりました。



さらに、その平均値のバラツキを調べるために標準誤差というものを求めました。それが重要になってくるのですが、t分布を利用した平均値の区間推定に必要な材料を以下に書きます。

標本平均と標準誤差

この2つが分かれば、OKです。標準誤差は不偏分散と標本数(データ数)で求められるので、

標本平均と不偏分散とデータ数

となり、必要な数値はこの3つと言い換えることもできます(標準誤差などは第6回で説明しました!)。

+統計量と+分布

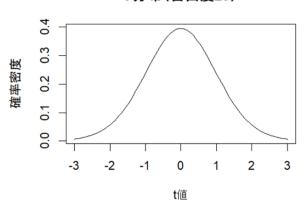
この辺りは少し難しい話になるので、興味を持てない人は飛ばしてください。いきなりですが、 t 統計量というものを考えます。

$$t = \frac{\frac{\cancel{max} + \cancel{max} + \cancel{max}}{\cancel{max}} = \frac{\cancel{max} + \cancel{max} + \cancel{max}}{\cancel{max}} = \frac{\cancel{max} + \cancel{max} + \cancel{max}}{\cancel{max}}$$

なにかよくわからないけれど、とりあえず手持ちの数値で t 統計量というものが求められそうだということだけでも感じてください。次に分布について考えてみますが、統計量 t というものは、自由度 n-1 の

t 分布という確率分布に従うということが知られています。というかそういうルールです。数学のポイントは、理由はよくわからないけど公式やルールを受け入れることです!自由度 n-1 とは、ここでは簡単に標本数-1 が自由度であると理解してください。いま、サンプル数(標本数)は30ですので、自由度は30-1=29となります。それを図にしたものが下に示すものです。

t 分布(自由度29)

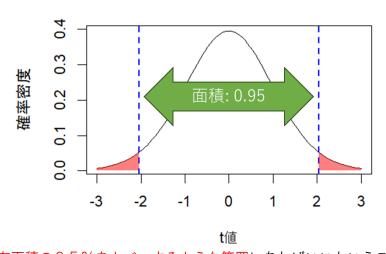


これが自由度29の t 分布です。確率分布とは曲線下の面積を全部足したら1になるということです。

95%信頼区間(+分布を利用した場合)

では次に、95%信頼区間とはどういうことかというと、下の図が直感的にわかりやすいです。

t 分布(自由度29)



t 値(t 統計量)が分布面積の9.5%をカバーするような範囲にあればいいということになります。上の図でいえば、サンプル数が3.0の t 分布の場合(自由度2.9)、だいたいですが、-2 から2 の間に t があれば 0 K ということです。厳密には、自由度2.9 の t 分布で9.5 %信頼区間(有意水準5.% ともいう)を示すような t 統計量は t 分布表というもので確認すると(Google 先生に聞いたらすぐに出てきます)、2.045 ということがわかります。ということは、t 分布は0 を挟んで左右対称なので、9.5 %信頼区間を表す t 統計量の値は、以下のようになります。

 $-2.045 \le t \le 2.045$

小休止

中途半端ですが今月はここで終わります。区間推定は、A から B までの範囲の中に求める数字があることを教えてくれます。点推定より実用性は高いかもしれません。計算に必要な推定値はこれまで触れてきたものですが、 t 分布や信頼区間の考え方は少し時間をかけて理解する必要があるかもしれません。