

生物データの取り扱い方

担当：遠藤・植山

実習場所：B11-118 (情報処理演習室)

1. 実習の目的

本実習では、実験実習や卒業研究などで取得される科学データの持つ意味を正確にわかりやすく伝える技術を身につける。科学データの整理方法を学び、わかりやすく伝えるための図表の書き方について学ぶ。

2. 使用データと表計算ソフト

使用するデータは、国内の森林において閉鎖型チャンバー法 (図1) を用いて観測された森林土壌の呼吸速度とその場所の地温である。それぞれのデータは30分毎に取得されたものであるが、今回の実習では1日平均したものが配布される。データの1列目には日にち、2列目には地温、3列目には土壌呼吸速度が保存されている。また、1行目にはデータの種類、2行目にはデータの単位がヘッダー情報として書かれており、実際のデータは3行目以降から始まる。なお、正しく測定ができなかった日については、“-9999”が異常値として記入されている。

データ整理には、表計算ソフト(Excel 2013, Microsoft®) を用いる。表計算ソフトは、平均や標準偏差等の統計値を計算することができるほか、図表としてデータを視覚的に整理することができる。



図1 森林土壌に取り付けられた土壌呼吸測定用チャンバー

使用するデータは、生態気象学研究グループのホームページ (<http://www.envi.osakafu-u.ac.jp/atmenv/>) にアップロードされている。

3. Excel を用いた作図

3-1. 散布図を用いた時系列グラフ

土壌呼吸速度と地温の季節変化を散布図で示す (図2)。散布図とは、縦軸と横軸の値を対応させた点を示した図である。今回は、横軸を日 (時間)、縦軸を土壌呼吸速度と地温とすることで、土壌呼吸速度や地温の日々の変化を図示する。

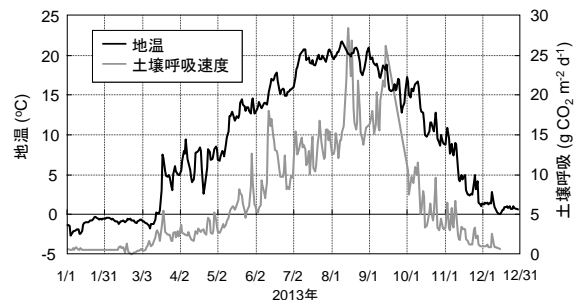


図2. 森林土壌で計測された地温と土壌呼吸速度の季節変化

■Excel 2013 による作図方法

1. 異常値を除いた作図のためのダミー列の挿入

Excel では、異常値 (-9999) を除いて作図することができない。ただし、「#N/A」値であれば、異常値を読み飛ばして作図される。そこで、IF 文を使用して異常値を「#N/A」値に置き換えた列を作成する。

例： `=IF(C3 = -9999, NA(), C3)`

2. 散布図を書く

Excel で標準出力されるグラフは、汚く見るに堪えない質であるため、見栄えよく調整する必要がある。ここでは、[図 2](#) を目指すべきグラフとして、それに対して汚いグラフの例を [図 3](#) に示す。[図 3](#) で調整すべき点は、以下の通りである。

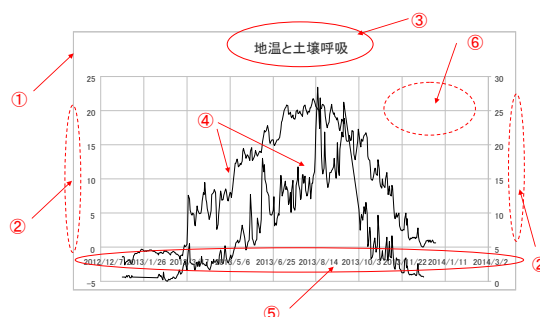


図 3. Excel 標準出力の汚いグラフの一例

- ① 図の枠線は消す。
- ② 軸見出しをつける。例えば、「地温 (°C)」や「土壤呼吸速度 (g CO₂ m² d⁻¹)」
 - CO₂の「2」は、必ず下付き文字とすること。
 - 単位にある-2乗や-1乗は、必ず上付き文字とすること。
- ③ 図のタイトルは、下に付けること。また、図表のタイトルはそれを読んだだけで何を示しているかが理解できる程度に詳しく、かつ簡潔に記載すること。
- ④ どの線が何を表しているかがわかるように、濃淡、種類（実線、破線、一点鎖線など）、太さ、色などを調整すること。
- ⑤ 目盛のタイプや区切りを標準のままにしないこと。この場合は、「月/日」がわかれば十分で年の情報は不要である。
- ⑥ 線の凡例 (Legend) を付けること。

その他、[図 3](#) の例には無いが、軸目盛の有効桁数を合わせることなど、注意が必要である。

使用する Excel 関数

IF 関数 : `=IF(条件, 条件が成り立つ場合の式あるいは値, 条件が成り立たなかった時の式あるいは値)`

例文 `=IF(C3=-9999,NA(),C3)`

NA 関数 : `=NA()`

異常値 (#N/A) を代入する関数 (N/A は、Not Available の省略語)

3-2. 回帰分析

土壌呼吸速度は微生物や植物の根の呼吸の総和である。一般に、植物や動物の呼吸速度は温度が上昇するとともに増加するとされている。このような要素間の関係性を定量化する手段として回帰分析が用いられる。回帰分析は、変数 X と変数 Y との間にどのような関係が成り立つか、あるいは関係が成り立たないかを評価する際に用いられる。変数 X と Y の間に、線形の関係（例えば、 $Y = a \times X + b$ ）がある場合は線形回帰、指数の関係（例えば、 $Y = a \times \exp(X \times b)$ ）がある場合は指数回帰を用いる。ここで、回帰分析によって得られる式は回帰式、回帰式の係数は回帰係数と呼ばれる。回帰式の当てはまりを表す指標として決定係数(R^2)や相関係数(R)といったものが用いられる。相関係数は、-1.0～1.0の値をとり、絶対値が1.0に近ければ、当てはまりが良い、あるいは変数 X と Y の間の関係性が高いことを表す（図4）。一般に、相関係数の絶対値が0.7を上回ると、両変数間に強い相関があるとされる（Wikipedia「相関係数」,2014）。

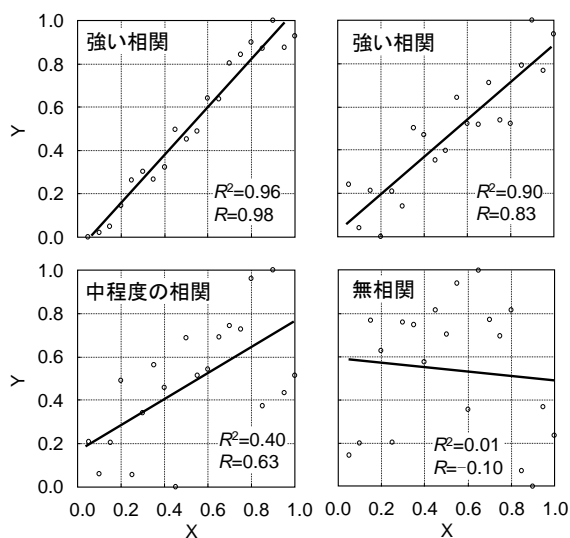


図4. データの分布と相関係数に関する散布図

■Excel 2013 による相関解析

1. 地温と土壌呼吸速度を選択して散布図を書く（図5）。

土壌呼吸速度と地温の関係を考える場合、地温の変化によって土壌呼吸速度が変化すると仮説が立てられる。逆に、土壌呼吸速度の変化で地温が変化するとは考えにくい。つまり、土壌呼吸速度は地温によって説明されるといえる。このような関係性が成立する場合、土壌呼吸速度のことを目的変数（従属変数）、地温のことを説明変数（独立変数）という。属性がはっきりしている両変数間の散布図を書く場合、横軸を説明変数、縦軸を目的変数としてデータをプロットする決まりがある。

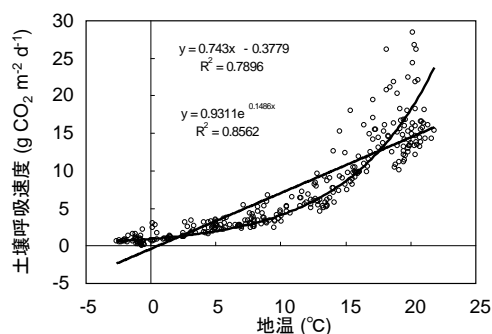


図5. 地温に対する土壌呼吸速度の関係

2. 散布図のオプションを指定して、近似線を追加する。

今回は、線形回帰式と指数回帰式の2つを図中に追加し、いずれの決定係数が大きいかで関数の当てはまりの良さを調べる。

図6は、Excelによる汚いグラフの一例である。既に図3で指摘されている点に加えて以下の点に注意して作図すること。

- ⑦ 縦軸の横目盛は、左側に詰める。横軸については、下側に詰める。
- ⑧ 回帰式の数字がグラフのプロットに重なっている。
- ⑨ 回帰式は太く見やすく書くこと。

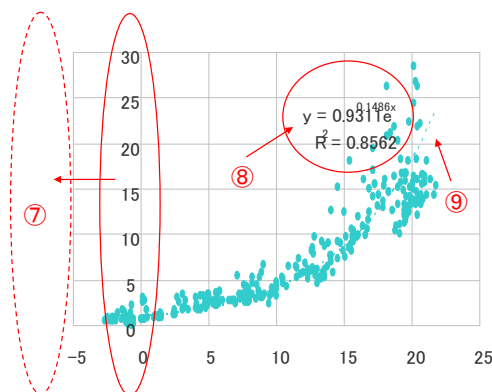


図6. 汚いグラフの例

3-3. Bin データ

回帰分析において、目的変数や説明変数が大きくばらつくために明確な関係性が評価できない場合がある。これは、データの信頼性が低くそれぞれのデータの不確か性が大きい場合、2変数以外の要因によってデータがばらついている場合（例えば、今回であれば土壌含水率などの他の要因によって土壌呼吸速度が変化している場合、地温と土壌呼吸速度との関係が不明瞭となる）などが挙げられる。このような場合、説明変数の値が近いいくつかのデータを平均して、平均された両変数間の相関解析を行うと関係性が明瞭になることが多い。このようにいくつかのデータをまとめたデータは Bin データと呼ばれる。

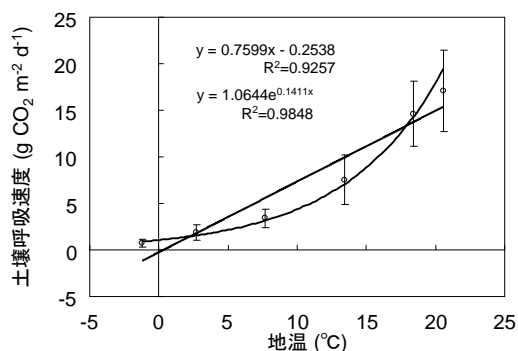


図7. 地温に対する土壌呼吸速度の関係。データは5°C毎の地温のクラスにまとめられている。

■Excel 2013 による Bin データの作成 (図7)

1. Bin データを用いた回帰分析

地温が5°C毎のクラス（例えば、0~5°C、5~10°C、10~15°C）で地温と土壌呼吸速度の平均値、標準偏差を計算する。

- 1-1. Bin データを作成するために、Excel 中の地温、土壌呼吸速度の行を別の行にコピーする。
- 1-2. コピーした行を選択して、地温についてデータの並び替えを昇順で行う。
- 1-3. 地温のクラス（例えば、0~5°C、5~10°C、10~15°C）ごとにデータセルを色分けする。
- 1-4. 同じ色のセルを選択して、地温、土壌呼吸速度のそれぞれについて平均、標準偏差を計算する。

2. Bin データの作成

上記の処理で作成された Bin データについて、散布図を書いて線形回帰、指数回帰分析を行う。ここで、各散布図のプロットに標準偏差を誤差線として書き加える。

使用する Excel 関数

Average 関数 : =Average (複数のセル; 例えば、A1:A100)

選択したセルの平均値を計算する。上例であれば、A 列の 1~100 行目までの平均を計算する。

Stdev 関数 : =Stdev (複数のセル; 例えば、A1:A100)

選択したセルの標準偏差を計算する。上例であれば、A 列の 1~100 行目までの標準偏差を計算する。
標準偏差とは、データのばらつきの程度を表す指標である。

3-4. 回帰モデルの精度評価

回帰分析によって得られた線形回帰式や指数回帰式の当てはまりの良さを調べる。ここで、回帰式から予測される変数値 (この場合は、地温から回帰式により予測される土壤呼吸速度の値) のことを推定値と呼ぶ。目的変数と推定値との比較から得られるさまざまな統計値を用いて当てはまりの良さを知ることができるが、ここでは線形回帰式から得られる傾き ($Y = a \times X + b$ の a)、切片 ($Y = a \times X + b$ の b)、相関係数 (R) について述べる。

目的変数 (観測された土壤呼吸速度; F_{obs}) と推定値 (推定された土壤呼吸速度; F_{model}) から以下の線形回帰式を得た場合 (図 8)、

$$F_{\text{model}} = \text{slope} \times F_{\text{obs}} + \text{intercept} \quad (1)$$

回帰係数 slope が 1.0 に近く、回帰係数 intercept が 0.0 に近い場合、回帰式の当てはまりが良いといえる。一方、 intercept が正の値に大きくなる、或いは slope が 1.0 よりも大きくなれば推定値が過大評価されていることを表す。逆に intercept が負の値をとる、或いは slope が 1.0 よりも小さくなれば推定値が過小評価されていることを表す。

■Excel 2013 による回帰モデルの精度評価

1. モデル値の計算

3-2 節、3-3 節で推定された線形回帰式、指数回帰式を使って土壤呼吸速度の推定値を 4 つ計算する。

2. 線形回帰分析による推定値の評価

各推定値の当てはまりの良さを、3~8 月の期間について調べる。Excel では、散布図の近似線の回帰係数を調べる関数 (SLOPE、INTERCEPT) や相関係数を計算する関数 (CORREL) があり、これらの関数を用いて当てはまりの良さを調べる。Excel で、これらの関数を使用する場合、解析対象セルに

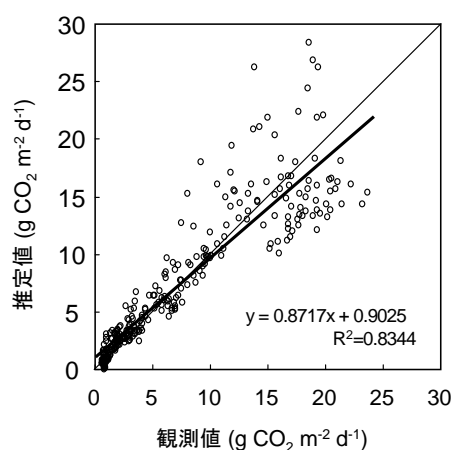


図 8. 指数回帰式による推定値と観測値による土壤呼吸速度の比較。細線は 1:1 線、太線は両者の直線回帰式を表す。

「#N/A」値が入っていると結果が「#N/A」となってしまう。そこで、「#N/A」値を空白に置き換えておく作業が必要となるが、今回のデータでは3～8月の期間に「#N/A」値が入っていないので問題としない。「#N/A」値を空白で置き換えたい場合は、3-1節で作成した異常値を除去する下記の式を

例文 =IF(C3=-9999,NA(),C3)

次のように書き換えるとよい。

例文 =IF(C3=-9999,"",C3)

今回の推定値の評価においては、1式にならって横軸(x軸)を F_{obs} 、縦軸(y軸)を F_{model} としてSlopeとintercept、そして決定係数(相関係数の2乗)を計算する。

使用する Excel 関数

EXP 関数 : =EXP(a)

eを底とする数値aのべき乗を計算する。

Slope 関数 : =Slope(既知の配列y, 既知の配列x)

$y = a \times x + b$ のaを計算する。

Intercept 関数 : =Intercept(既知の配列y, 既知の配列x)

$y = a \times x + b$ のbを計算する。

Correl 関数 : =Correl(配列a, 配列b)

配列aと配列bの間の相関係数を計算する。

Power 関数 : =Power(数値a, 指数b)

a^b を計算する。

Excel の相対参照と絶対参照

Excelでは、セルの参照形態に相対参照と絶対参照の2種類がある。例えば、以下の式は相対参照の式であり、A列1行目とB列1行目の積を計算するものである。

=A1 * B1

この式がC1に書かれている場合、式をC2にコピーすれば「=A2 * B2」となり、C3にコピーすると「=A3 * B3」となる。また、D2にコピーすれば「=B2 * C2」となり、参照するセルが相対的に移動する。一方、Aセルを絶対参照で記載した以下の式の場合、

=\$A1 * B1

表1. 観測された土壌呼吸速度に対する回帰式による推定値の関係。Slope、intercept、 R^2 は線形回帰式から算出された。1.0より大きなslope、0.0よりも小さなinterceptは、推定値の過大評価を示す。

	線形回帰	指数回帰	線形回帰.bin	指数回帰.bin
slope	0.73	1.01	0.75	0.98
intercept	2.57	0.26	2.76	0.55
R^2	0.77	0.81	0.77	0.81

表2. 推奨される書き方が表1である場合の作表の悪例。

	線形	指数	線形.bin	指数.bin
slope	0.730806	1.013356	0.747429	0.983272
intercept	2.565484	0.258013	2.756534	0.554713
R^2	0.765761	0.808898	0.765761	0.812358

表. 回帰式の推定値の評価

式をD2にコピーすれば「 $=\$A2 * C2$ 」となり、A列に対する位置は変わらない。あるいは、

$$=\$A\$1 * B1$$

と記載すると、D2にコピーしても「 $=\$A\$1 * C2$ 」となり、A1セルへの参照位置は変わらない。このようにExcelでは、セルを参照する場合に「\$」を付けると絶対参照となる。「\$」を付ける位置としては、「 $\$A\1 」、「 $A\$1$ 」、「 $\$A1$ 」の3つのパターンがあり、それぞれ「行列の固定」、「行の固定」、「列の固定」を表す。「\$」は「F4」キーを用いると簡単に挿入できる。例えば数式中のA1を選択して「F4」キーを押すごとに、 $\$A\1 (行列の固定) → $A\$1$ (行の固定) → $\$A1$ (列の固定) → A1 (相対参照) の順に切り替わる。

表の書き方

モデルの当てはまりを評価するために3~8月の期間のslope、intercept、 R^2 を表1のような体裁で作図せよ。表2にしめされる悪例とならないように気を付ける作表する。

- ⑩ 表に縦線は不要。
- ⑪ 無駄に小数点以下の桁数を示しすぎない。有効桁数を意識して作表する。
- ⑫ 表の項目は分かりやすい名前を付ける。
- ⑬ 表中のフォントは、明朝体系のフォントを使用すること。「Times New Roman」などのフォントが推奨される。

- ⑭ 抽象的な表タイトルを使用せず、タイトルから表の中身をわかるような具体的なタイトルとする。
また、表のタイトルは、表の上を書くこと。

4. レポート

1. 実習で取り扱った地温と土壌呼吸速度のデータを使って以下の図表を作成せよ。作図の際は、①～⑭で指摘された点に注意して作図すること。

地温と土壌呼吸速度の季節変化を示す散布図 (図 2)

地温に対する土壌呼吸速度の関係を示す散布図 (図 5)

(線形回帰式、指数回帰式、それぞれの決定係数を図表に含めること)

地温に対する土壌呼吸速度の関係を示す Bin データを用いた散布図 (図 6)

(線形回帰式、指数回帰式、それぞれの決定係数を図表に含めること)

回帰式の推定値の当てはまりを示す slope、intercept、決定係数を表す表 (表 1)

線形回帰、指数回帰、Bin データを用いた両回帰式による推定値の評価を行う。

評価の期間は、3-8 月についてとする。

2. 別途配布する森林における気温と蒸発散速度のデータについても、上記の地温と土壌呼吸速度のデータと同様の図表を作成せよ。ここで、気温の Bin データを作る際は、5°C 毎 (例えば、0~5°C、5~10°C、10~15°C) にクラス分けすること。また、推定値の当てはまりに関する表は、4~8 月の期間について示せ。
3. レポートは上記の図表を印刷の上、本日から 2 週間以内に植山 (B11 棟 239 号室) まで提出すること (私が居室している時間であれば最終提出時間は問わない)。図表は、課題 1、課題 2 についてそれぞれ 1 ページでまとめること。遅れたものは採点しないので、注意すること。