

統計学・データ分析勉強会／第4回宿題解答例

第4回宿題の解答例を以下に示します。今回の問題は、時系列データの分析をする際に陥りがちな「見せかけの回帰」について知っておくことを目的としています。

問1

まず `corp_stats` という名前でデータを読み込みます。

```
setwd('/Users/stanaka/Rstats')  
corp_stats <- read.csv('演習データ (時系列) .csv')
```

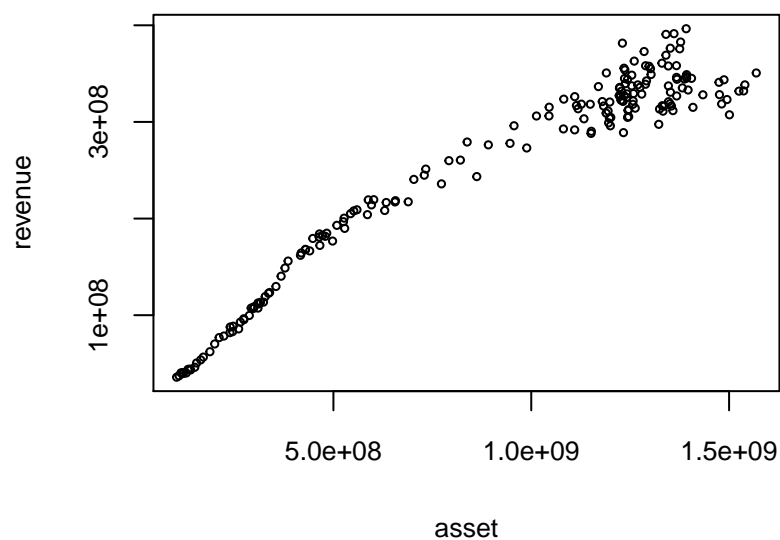
中身を確認しましょう。`corp_stats` は2列目に期間（四半期）、3列目に総資産、4列目に売上高が入ったデータフレームです。

```
head(corp_stats, 5)
```

```
##           期間      総資産   売上高  
## 1 1970年1 - 3 月 103761211 35826142  
## 2 1970年4 - 6 月 109870599 37248437  
## 3 1970年7 - 9 月 114662139 40312069  
## 4 1970年10-12月 118759543 40617531  
## 5 1971年1 - 3 月 120910968 39666742
```

最初に、総資産（`asset`）を説明変数、売上高（`revenue`）を被説明変数として回帰分析を試みましょう。`corp_stats` から両データを抽出し、散布図と相関係数を表示し、回帰式を求めます。

```
asset <- corp_stats[,2]  
revenue <- corp_stats[,3]  
plot(asset, revenue, cex=0.5, cex.axis=0.8, cex.lab=0.8)
```



```
cor(asset, revenue)
```

```
## [1] 0.9668794
```

```
reg <- lm(revenue~asset)
summary(reg)
```

```
##
## Call:
## lm(formula = revenue ~ asset)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -73558700 -19102528  2233536  21738975  59773623
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 5.244e+07  4.318e+06  12.14  <2e-16 ***
## asset       2.189e-01  4.225e-03  51.80  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 26630000 on 187 degrees of freedom
## Multiple R-squared:  0.9349, Adjusted R-squared:  0.9345
## F-statistic: 2684 on 1 and 187 DF, p-value: < 2.2e-16
```

散布図を見ると売上高と総資産はきれいな直線に乗っており、相関係数も約 0.97 と非常に高い値を示しています。当然、回帰分析は検定をパスしており ($p < 2e-16$)、決定係数も約 0.93 と申し分ありません。この結果から、ある四半期における企業の売上高は、その四半期における総資産の値によって高い精度で説明できると結論してよいでしょうか？

結論から言うと、**この分析は誤りです**。問2でその理由を考えます。

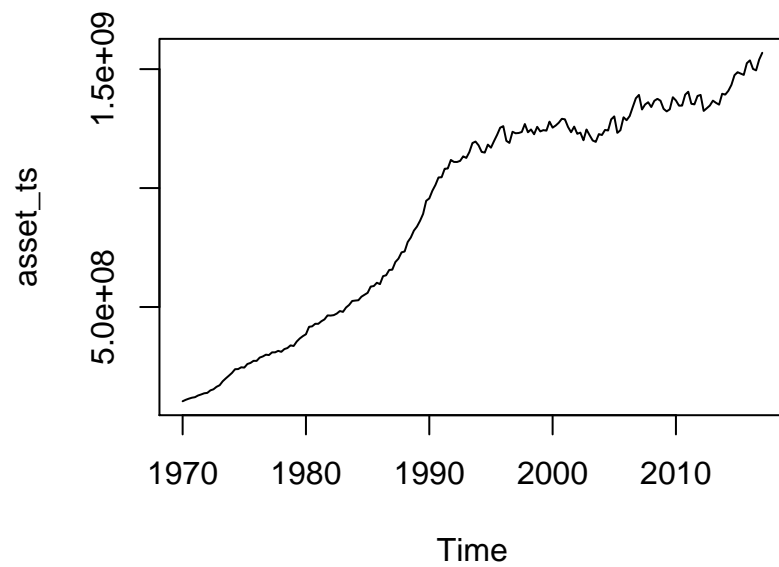
問2

今回扱っているデータセット `corp_stats` は、1970～2017 年まで四半期ごとに企業の売上高と総資産を集計した時系列のデータでした。データを明示的に時系列データとして扱うときは、`ts()` 関数を使うのが便利です。

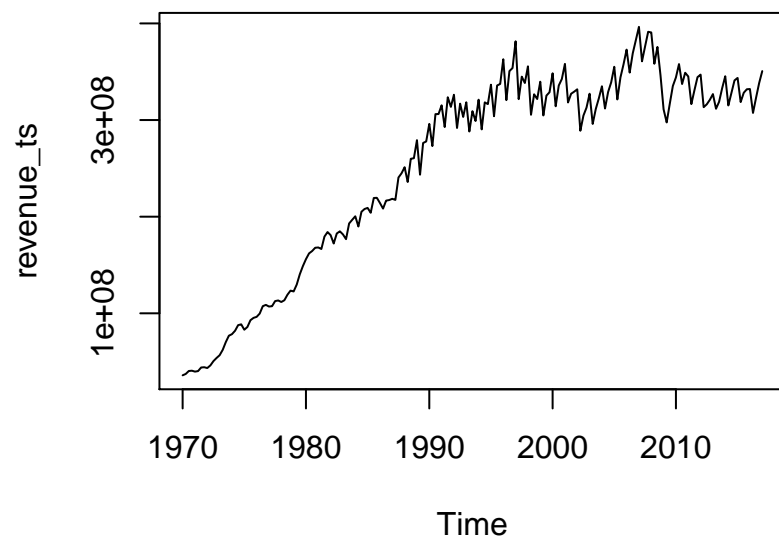
```
asset_ts <- ts(asset, frequency = 4, start=c(1970, 1))
revenue_ts <- ts(revenue, frequency = 4, start=c(1970, 1))
```

この段階で、総資産と売上高をそれぞれグラフにしてみましょう。

```
plot(asset_ts, cex=0.8, cex.axis=0.8, cex.lab=0.8)
```

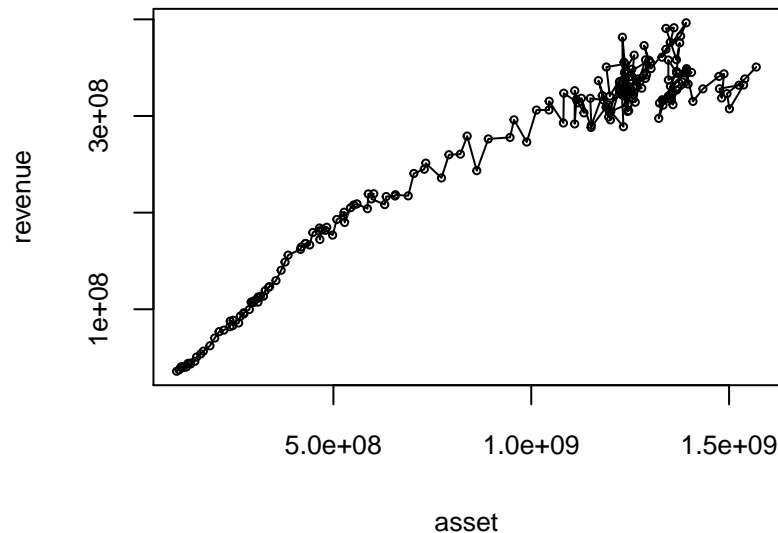


```
plot(revenue_ts, cex=0.8, cex.axis=0.8, cex.lab=0.8)
```



総資産、売上高とも年（四半期）を追うごとに増加しています。さらに、先ほど描いた散布図において、各データ点を登場する順番に線で結んでみましょう。plot() 関数のオプションとして `type='o'` と指定すれば、データを時系列順に結んだ散布図を作成できます。

```
plot(revenue~asset, type='o', cex=0.5, cex.axis=0.8, cex.lab=0.8)
```



これを見ると、散布図上で互いに近接して現れるデータ点は、時期（四半期）が近い（多くの場合は隣り合っている）データであることが分かります。散布図の左下ほど昔のデータ点、右上ほど最近のデータ点を表しています。従って、ある時点の売上高を説明する要因としては、総資産の値よりもまず時間の要素を疑うべきということになります。

勉強会の第1～3回では、「時間に依存せず、同一の分布から独立に抽出したデータ」を前提として、仮説検定や回帰分析などの手法を学んできました。一方、今回のような時系列データを分析するうえでも、重要な前提が存在します。それは「**データが定常性をもつ（ざっくり言うと、データの平均や分散が時間によらず常に一定）**」という条件を満たすことです。定常性を持たない2つの時系列データ（厳密には単位根過程と呼ばれます）同士を単純に回帰分析すると、それらが互いに全く無関係の系列でも検定が有意になったり、決定係数が高くなったりといった「**見せかけの回帰**」という現象が生じてしまうことが知られています。

上のグラフや散布図から分かるように、今回扱っている総資産や売上高のデータは四半期を追うごとに値が増加しており、明らかに定常性の条件を満たしません。このため、2つの変量を単純に回帰分析することはできません。

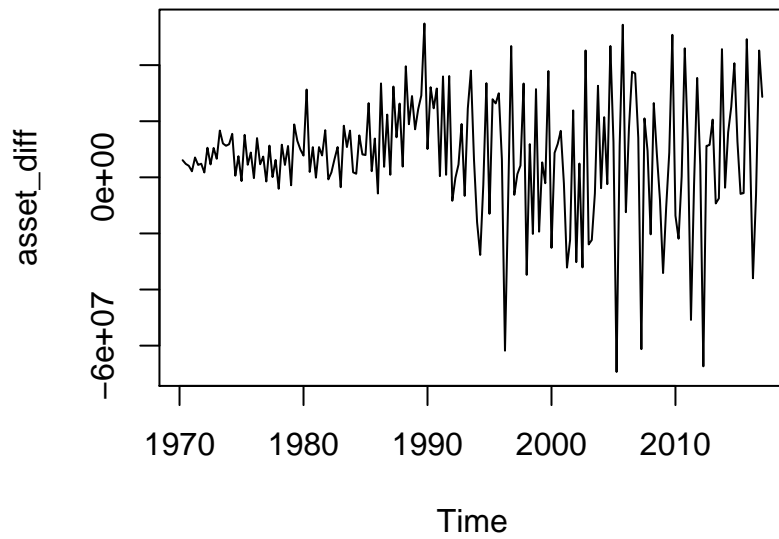
問3

見せかけの回帰を回避する方法の一つは、差分系列（時系列データにおいて、1時点前との差分を取った系列）を分析することです。diff() 関数を使って、総資産と売上高それぞれのデータについて、1四半期前との差分を計算します。

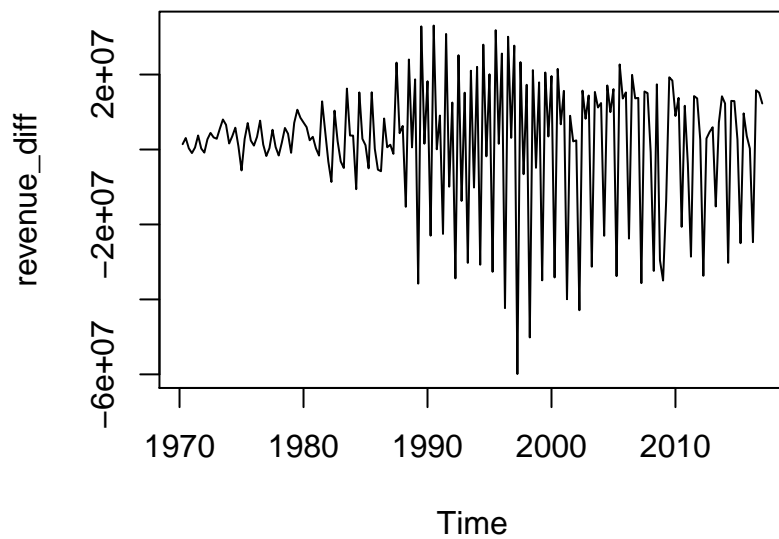
```
asset_diff = diff(asset_ts)
revenue_diff = diff(revenue_ts)
```

先ほどと同じように、それぞれの系列のグラフと散布図を作成し、相関係数を計算しましょう。

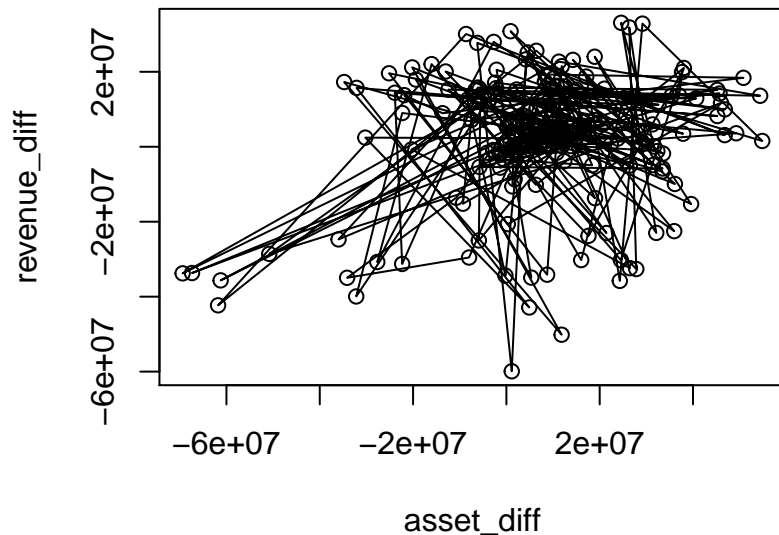
```
plot(asset_diff, cex=0.8, cex.axis=0.8, cex.lab=0.8)
```



```
plot(revenue_diff, cex=0.8, cex.axis=0.8, cex.lab=0.8)
```



```
plot(asset_diff, revenue_diff, type='o', cex=0.8, cex.axis=0.8, cex.lab=0.8)
```



```
cor(asset_diff, revenue_diff)
```

```
## [1] 0.2335809
```

まだ完全に平均や分散が一定とは言えませんが、先ほどと比べてかなり定常なデータと言えます。散布図上のデータの現れ方も、時間に強く依存せずランダムに近づきました。一方で、相関係数をみると売上高と総資産の間の直線関係はだいぶ弱くなってしまいました。この状態でも「売上高を総資産で説明できる」と言えるでしょうか。改めて回帰分析を実行してみましょう。

問4

差分系列で回帰分析を再度実行します。

```
reg_new <- lm(revenue_diff~asset_diff)
summary(reg_new)
```

```
##
## Call:
## lm(formula = revenue_diff ~ asset_diff)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -60331775 -5894921  1131411  11325090  31506354
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.065e+05  1.323e+06   0.156  0.87610
## asset_diff   1.884e-01  5.751e-02   3.276  0.00125 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 17060000 on 186 degrees of freedom
## Multiple R-squared:  0.05456,    Adjusted R-squared:  0.04948
```

F-statistic: 10.73 on 1 and 186 DF, p-value: 0.001255

傾きの検定は有意水準 5% でパスしていますが、決定係数は 0.05 程度と非常に低い値になっています。従って、売上高の変動を総資産で説明するのは無理があるという結論になりました。残念ながら問 1 で現れた 2 変量の強い関係性は見せかけであり、他のデータセットや分析手法を試す必要があります。