

医学統計勉強会

東北大学病院循環器内科・東北大学臨床研究推進センター 共催

東北大学大学院医学系研究科EBM開発学寄附講座

宮田 敏

回帰分析

回帰分析 (regression analysis) は、一つの連続数（実数）の値を複数の変数によって説明、予測する**多変量解析 (multivariate analysis)** の一つ。

$$Y = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k + \varepsilon$$

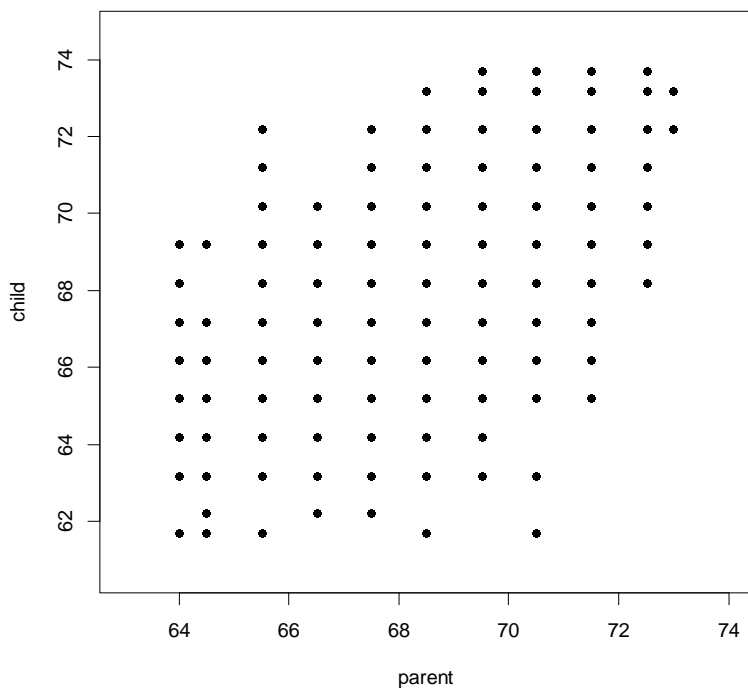
Y : response variable, 従属変数, 被説明変数

x_1, \dots, x_k : independent variable, 独立変数, 説明変数

$\beta_0, \beta_1, \dots, \beta_k$: regression coefficient, 回帰係数

ε : error term, 攪乱項、誤差項

Example 1: 親子の身長



205組の夫婦から生まれた、928人の成人した子供の身長（インチ）。

child: 子の身長

parent: 両親の身長の平均

データは0.1インチ刻みに丸められている。そのためデータ点が重なっている。

Galton, F. (1886). Regression Towards Mediocrity in Hereditary Stature
Journal of the Anthropological Institute, 15, 246-263

2014/10/9

東北大学 医学統計勉強会

3

データの要約

データが手に入ったら、まずデータを要約して、その傾向、特徴を把握する。

個々の変数の数量的要約：Galtonの親子の身長

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.		SD	IQR
parent	64	67.50	68.5	68.31	69.5	73		1.79	2
child	61.7	66.20	68.2	68.09	70.2	73.7		2.52	4

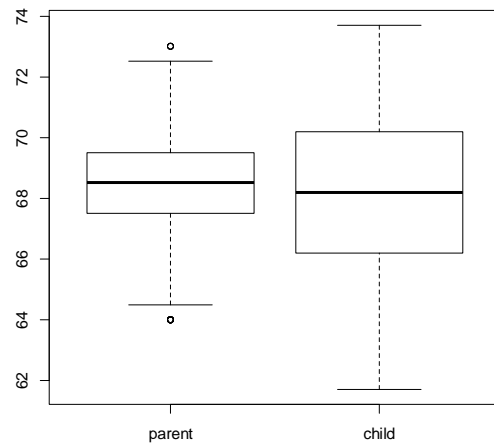
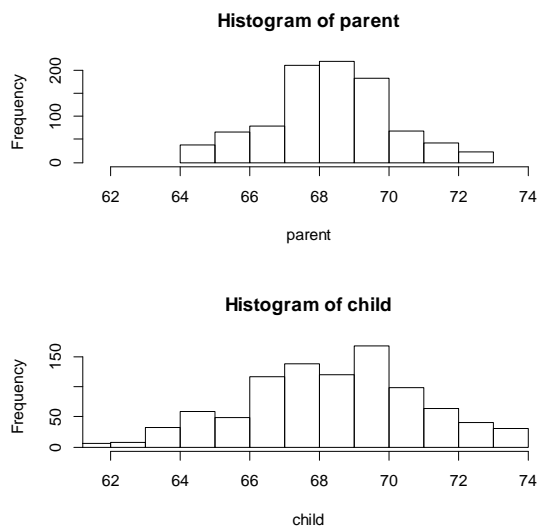
- 平均，中央値はほぼ等しく，分布の位置は同じ。
- 標準偏差（SD），四分位点間距離（IQR）は childの方が大きく，子の分布の散らばりの方が大きい。

2014/10/9

東北大学 医学統計勉強会

4

個々の変数の視覚的要約：Galtonの親子の身長



グラフを比較するときは、軸を揃えることがコツ。

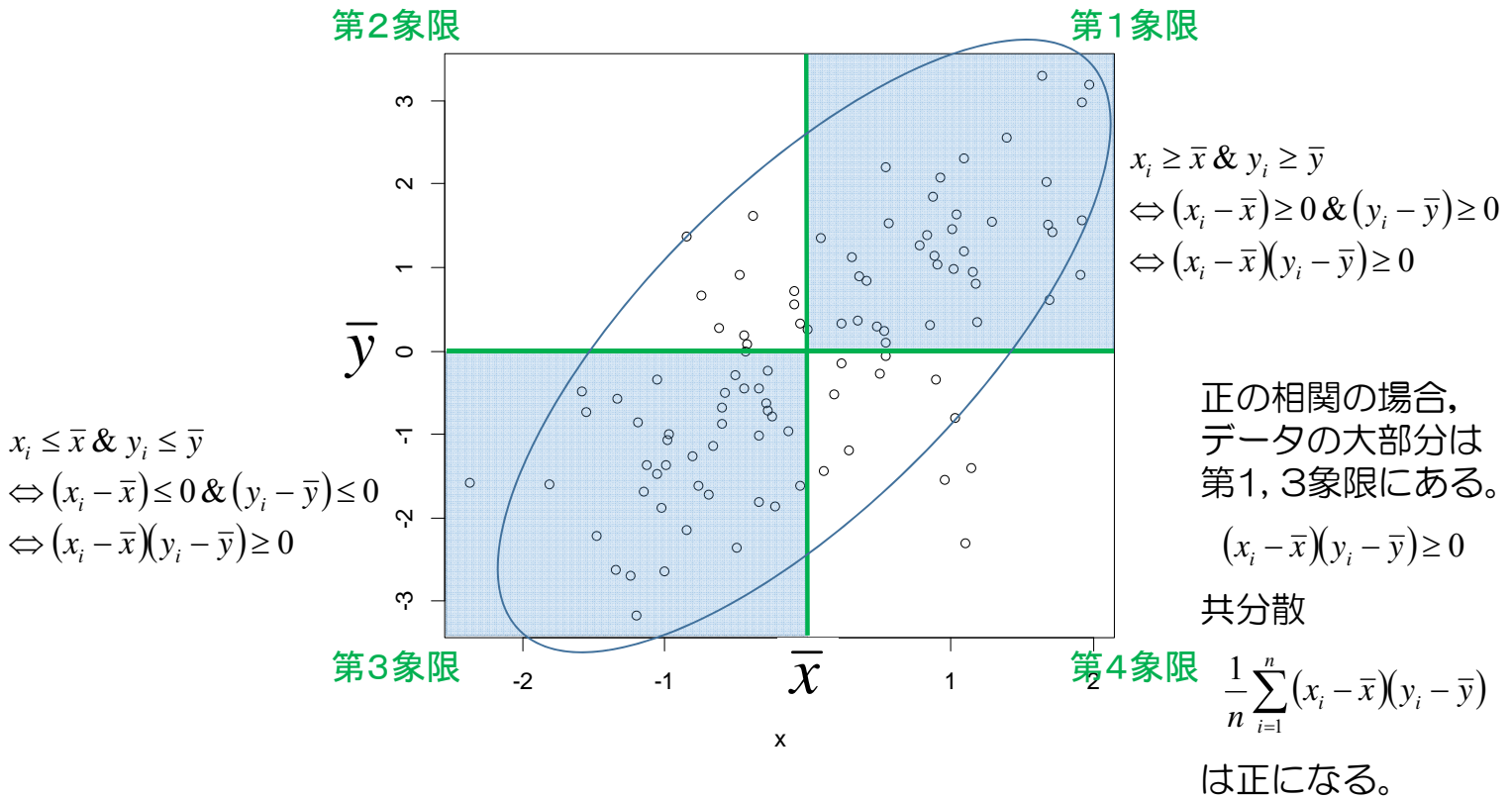
標本共分散

二つの変数 x と y が互いに影響し合っているとき、 x と y が如何に強く関係しあっているか知りたい。このとき x と y の**共分散(covariance)**を以下で定義する。

$$\text{Cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$$

x と y の間に**正の(負の)**相関があるとき、 $\text{Cov}(X, Y)$ はそれぞれ**正(負)**になる。

正の相関の場合



2014/10/9

東北大学 医学統計勉強会

7

標本相関係数

二つの変数 x と y の間の **相関係数 (correlation coefficient)** を、以下で定義する。相関係数は x と y の **線形関係の強さ** を測る量である。

$$\text{Corr}(X, Y) = r = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

- 相関係数 r は x と y の線形関係の強さを測る。
- $-1 \leq r \leq 1$
- $r = +(-) 1$: 正(負)の完全な相関, 線形関係

2014/10/9

東北大学 医学統計勉強会

8

回帰分析

二つの変数 x と y の関係が線形（直線）で近似できるとする。このとき x と y の関係を以下の回帰式 (regression equation) でモデル化する。

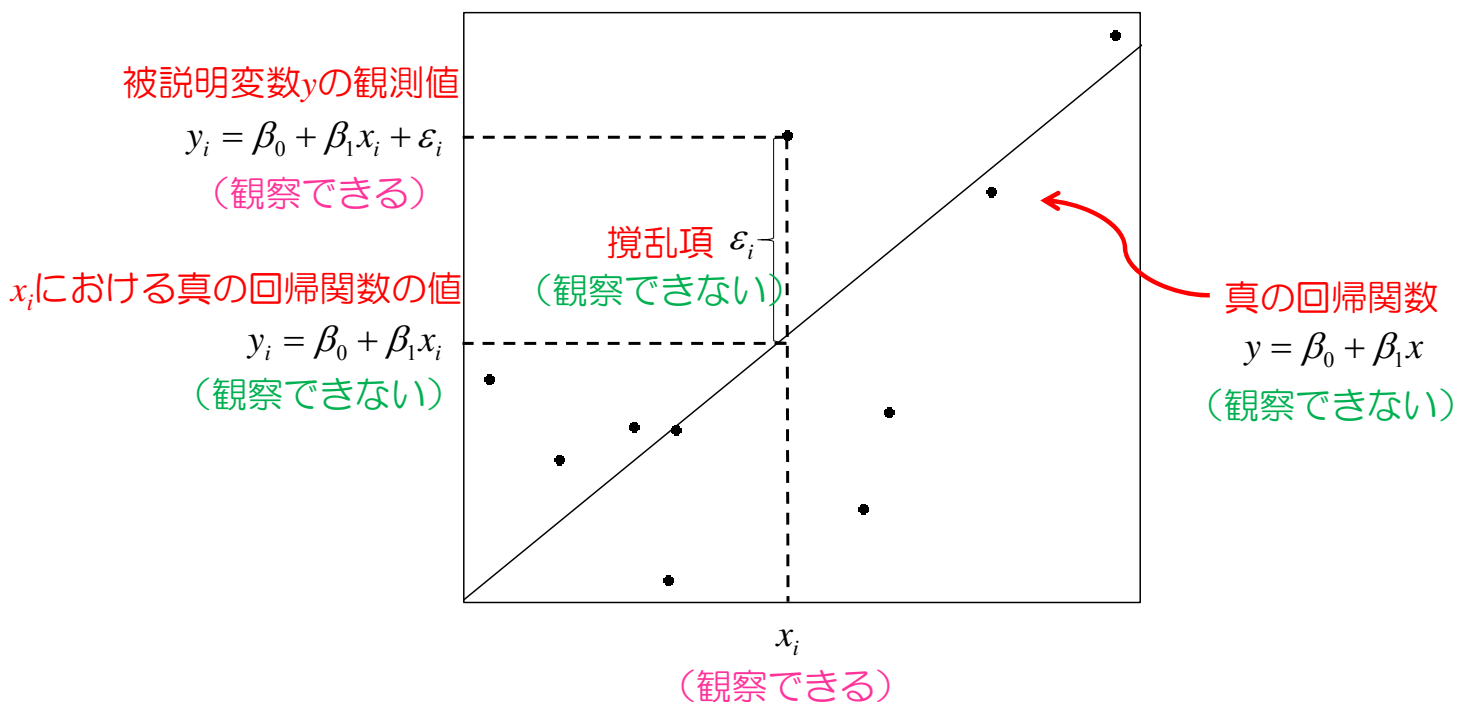
$$Y = \beta_0 + \beta_1 x + \varepsilon$$

Y : response variable, 従属変数, 被説明変数

x : independent variable, 独立変数, 説明変数

β_0, β_1 : regression coefficient, 回帰係数

ε : error term, 攪乱項、誤差項



回帰分析のモデルの仮定

- **線形性 (Linearity)** : $Y = \beta_0 + \beta_1 x + \varepsilon$
被説明変数 y と説明変数 x の関係は直線で近似できる。
- **独立性 (Independence)** $\{(x_i, Y_i)\}_{i=1}^n$ は互いに独立である。
あるサンプルの値が他のサンプルの値に影響しない。
- **正規性 (Normality)** : $\varepsilon_i \sim N(0, \sigma^2)$, iid
攪乱項 ε_i は正規分布に従う。
- **等分散性 (homoskedasticity)** : σ^2 . 分散一定

正規分布 : Normal distribution

$$f(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}, -\infty < x < \infty$$

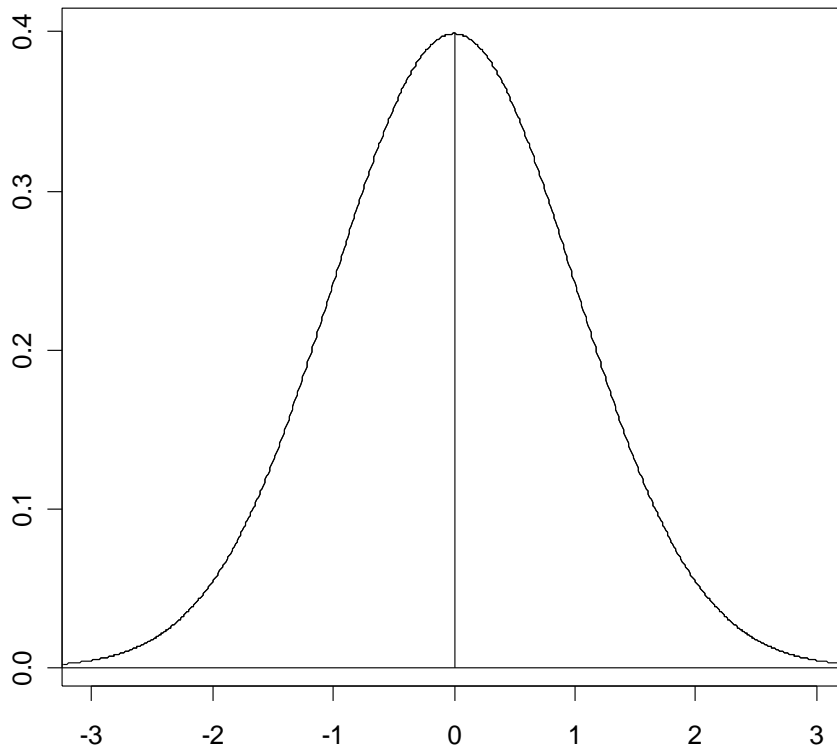
- $E(X) = \mu, V(X) = \sigma^2$.
- 分布の形状は μ と σ^2 によって特徴づけ (parameterized) される。
- **標準正規分布** (standard normal distribution)

$$f(z; 0, 1) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}, -\infty < z < \infty,$$

$$\Phi(z) = P(Z < z) = \int_{-\infty}^z f(y; 0, 1) dy.$$

- **釣り鐘型 (bell-shaped)** で、**左右対称**な分布。

Standard Normal Distribution



2014/10/9

東北大学 医学統計勉強会

13

回帰係数の推定

最小二乗法 (ordinary least squares estimation, OLSE)

$$\begin{aligned} & \min_{\beta_0, \beta_1} \sum_{i=1}^n \{y_i - (\beta_0 + \beta_1 x_i)\}^2 \\ \Rightarrow & \begin{cases} \frac{\partial}{\partial \beta_0} \sum_{i=1}^n \{y_i - (\beta_0 + \beta_1 x_i)\}^2 = 0 \\ \frac{\partial}{\partial \beta_1} \sum_{i=1}^n \{y_i - (\beta_0 + \beta_1 x_i)\}^2 = 0 \end{cases} \Rightarrow \begin{cases} n\beta_0 + \beta_1 \left(\sum_{i=1}^n x_i\right) = \sum_{i=1}^n y_i \\ \beta_0 \left(\sum_{i=1}^n x_i\right) + \beta_1 \left(\sum_{i=1}^n x_i^2\right) = \sum_{i=1}^n x_i y_i \end{cases} \\ \Rightarrow & \begin{cases} \hat{\beta}_1 = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) / \sum_{i=1}^n (x_i - \bar{x})^2 \\ \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \end{cases} \end{aligned}$$

推定された回帰直線： $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

2014/10/9

東北大学 医学統計勉強会

14

最小二乗推定量の性質

- $E(\hat{\beta}_0) = \beta_0, E(\hat{\beta}_1) = \beta_1$, **不偏性**
- $\sigma_{\hat{\beta}_0}^2 = \text{var}(\hat{\beta}_0) = \frac{\sigma^2 \sum x_i^2}{n \sum (x_i - \bar{x})^2}$,
- $\sigma_{\hat{\beta}_1}^2 = \text{var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$,
- $\hat{\beta}_0 \sim N(\beta_0, \sigma_{\hat{\beta}_0}^2), \hat{\beta}_1 \sim N(\beta_1, \sigma_{\hat{\beta}_1}^2)$, **正規性**
- $s^2 = \hat{\sigma}^2 = \frac{\sum (y_i - \hat{y}_i)^2}{n-2} = \frac{\text{SSE}}{n-2}, E(\hat{\sigma}^2) = \sigma^2$,

- 信頼区間 (Confidence Interval):

$$T = \frac{\hat{\beta}_1 - \beta_1}{s_{\hat{\beta}_1}} \sim t_{n-2}, s_{\hat{\beta}_1} = \frac{s}{\sqrt{\sum (x_i - \bar{x})^2}}$$

$$\Rightarrow \text{C.I. of } \beta_1 : \left(\hat{\beta}_1 - s_{\hat{\beta}_1} t_{\alpha/2, n-2}, \hat{\beta}_1 + s_{\hat{\beta}_1} t_{\alpha/2, n-2} \right)$$

$$\text{同様に, C.I. of } \beta_0 : \left(\hat{\beta}_0 - s_{\hat{\beta}_0} t_{\alpha/2, n-2}, \hat{\beta}_0 + s_{\hat{\beta}_0} t_{\alpha/2, n-2} \right)$$

- 仮説検定 (Hypothesis Testing):

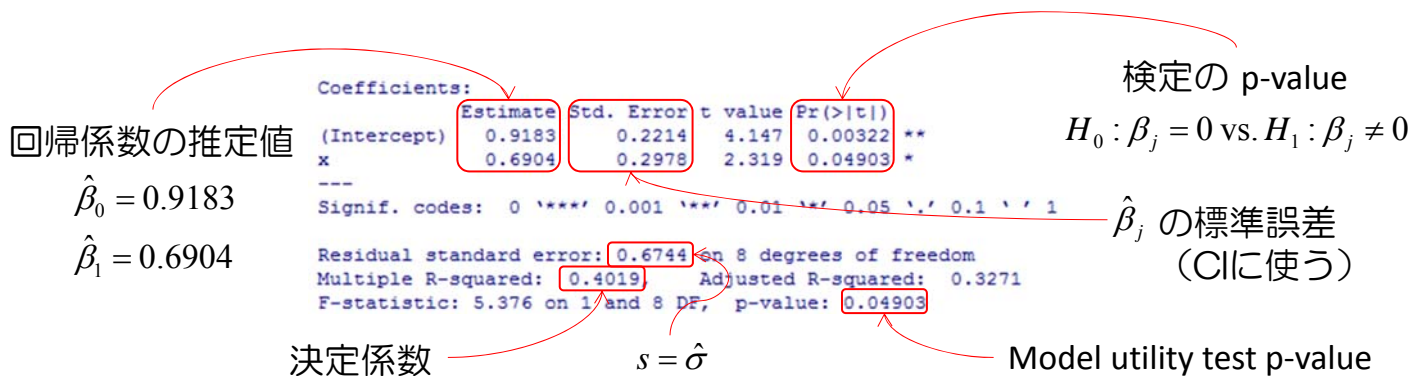
$$H_0 : \beta_1 = \beta_{10} \text{ vs. } H_a : \beta_1 \neq (\text{or } < \text{ or } >) \beta_{10}$$

$$\text{検定統計量: } T = \frac{\hat{\beta}_1 - \beta_{10}}{s_{\hat{\beta}_1}} \sim t_{n-2} \text{ under } H_0$$

$H_a : \beta_1 > \beta_{10}$	$t > t_{\alpha, n-2}$
$H_a : \beta_1 < \beta_{10}$	$t < -t_{\alpha, n-2}$
$H_a : \beta_1 \neq \beta_{10}$	$ t > t_{\alpha/2, n-2}$

回帰分析の結果

- 回帰係数の推定値 $\hat{\beta}_0 = 0.9183, \hat{\beta}_1 = 0.6904$
- 回帰係数の有意性検定のp値
- 決定係数（被説明変数の変動のうち回帰によって説明された変動の割合） y の変動の40.19%が説明された
- Model utility test（回帰モデル全体の有意性検定。後でもう一度触れます）のp値 $p=0.049$
- 攪乱項の標準誤差 $s=0.6744$



多変量回帰分析

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \varepsilon, \varepsilon \sim N(0, \sigma^2)$$

Y : response variable, 従属変数, 被説明変数

x_1, \dots, x_k : independent variable, 独立変数, 説明変数

$\beta_0, \beta_1, \dots, \beta_k$: regression coefficient, 回帰係数

ε : error term, 攪乱項、誤差項

モデルの仮定

1. パラメーターに関する線形性 (Linearity)
2. 攪乱項の独立性 (Independence)
3. 攪乱項の正規性 (Normality)
4. 攪乱項の等分散性 (homoscedasticity)

最小二乗法によるパラメータの推定

- $\beta_0, \beta_1, \dots, \beta_k$ は最小二乗法で推定される:

$$\min_{\beta_0, \dots, \beta_k} \sum_{j=1}^n (y_j - (\beta_0 + \beta_1 x_{1j} + \dots + \beta_k x_{kj}))^2$$

$$\hat{\beta}_{k \times 1} = (X'X)^{-1} X'Y,$$

$$\text{where } X = [x_{ij}]_{n \times (k+1)}, Y = (Y_1, \dots, Y_n)'$$

- $E(\hat{\beta}_j) = \beta_j, j = 1, \dots, k$: unbiased.

- $V(\hat{\beta}) = \sigma^2 (X'X)^{-1}$.

- $\hat{\beta} \sim N(\beta, \sigma^2 (X'X)^{-1})$

パラメータの推測、信頼区間、検定

- $s^2 = \hat{\sigma}^2 = \frac{\sum (y_i - \hat{y}_i)^2}{n - (k + 1)} = \frac{\text{SSE}}{n - (k + 1)}$. $E(s^2) = \sigma^2$.

- $T = \frac{\hat{\beta}_j - \beta_j}{s_{\hat{\beta}_j}} \sim t_{n-(k+1)}, s_{\hat{\beta}_j} = \left([s^2 (X'X)^{-1}]_{jj} \right)^{1/2}, j = 1, \dots, k$.

- C.I. of β_j : $(\hat{\beta}_j - s_{\hat{\beta}_j} t_{\alpha/2, n-(k+1)}, \hat{\beta}_j + s_{\hat{\beta}_j} t_{\alpha/2, n-(k+1)})$

- Hypothesis testing : $H_0 : \beta_j = \beta_{j0}$ vs. $H_a : \beta_j \neq \beta_{j0}$,

$$\text{Test statistic : } T = \frac{\hat{\beta}_j - \beta_{j0}}{s_{\hat{\beta}_j}} \sim t_{n-(k+1)} \text{ under } H_0.$$

決定係数 (coefficient of determination)

回帰係数の有意性検定は個々の係数の検定。回帰モデル全体のパフォーマンスを測る方法が必要。

SST (Total Sum of Squares): $\sum_{i=1}^n (y_i - \bar{y})^2$
yの全変動

SSE (Error Sum of Squares): $\sum_{i=1}^n (y_i - \hat{y}_i)^2$
回帰で説明されなかった変動

SSR (Regression sum of Squares): $\sum_{i=1}^n (\hat{y}_i - \bar{y})^2$
回帰で説明された変動

$R^2 = \text{決定係数} = \frac{SSR}{SST} = \text{回帰によって説明された変動の割合}$

Model utility test

回帰モデル全体の有意性を検定するために、以下の **Model Utility Test** が知られている。

$$H_0 : \beta_1 = \dots = \beta_k = 0 \text{ vs. } H_1 : \text{not } H_0$$

検定統計量：

$$F = \frac{R^2/k}{(1-R^2)/[n-(k+1)]} = \frac{SSR/k}{SSE/[n-(k+1)]} \sim F_{k, n-(k+1)} \text{ under } H_0$$

棄却域：Reject H_0 if $F > F_{\alpha, k, n-(k+1)}$

H_0 は、 $Y = \beta_0 + \varepsilon$ と同義。回帰が全く無効。

回帰分析の結果のチェックポイント

- 回帰係数の**推定値**
- 回帰係数の有意性**検定のp値**
- 回帰係数の**信頼区間**
- **決定係数**（被説明変数の変動のうち回帰によって説明された変動の割合）
 - ✓ 決定係数は、いくつ以上なら良いか？
 - もちろん、 R^2 は大きいに越したことはない。 R^2 が小さいということは、 y に影響する、今のモデルに含まれない**変動要因**が存在する、ということ。もし、 R^2 が小さくてもなお有意な説明変数があれば、上記の limitationの下で、回帰は有効である。
- **Model utility test**（回帰モデル全体の有意性検定）の**p値**
- 攪乱項の**標準誤差**

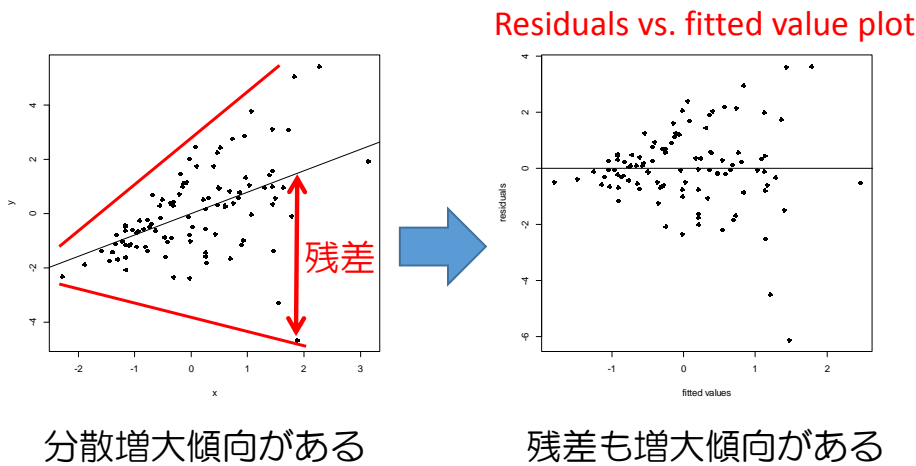
回帰診断

回帰モデルの仮定の確認：

- **線形性**: Y と x の間の線形関係。 Y と x の間には非線形な関係がない。 x 同士の間には線形関係がない。
 - Multiple scatter plots.
- **独立性**: $\{(x_i, Y_i)\}_{i=1}^n$ は互いに独立
 - residual vs. fitted value plot
- **正規性**: $\varepsilon \sim N(0, \sigma^2)$
 - 残差のQQ-norm plot （後述する）
- **等分散性**: $\text{Var}(\varepsilon) = \sigma^2$
 - residual vs. fitted value plot

回帰診断 (続き)

独立性, 正規性, 等分散性の仮定は, いずれも攪乱項についての仮定. 攪乱項そのものは観察できないため, 残差 (residuals): $e_i = y_i - \hat{y}_i$ をレプリカとして使う.



- 残差が均一ならば、等分散性の仮定は満たされる。
- 独立性の仮定が満たされる場合、残差プロットには特異なパターンがない。

分布の正規性の確認

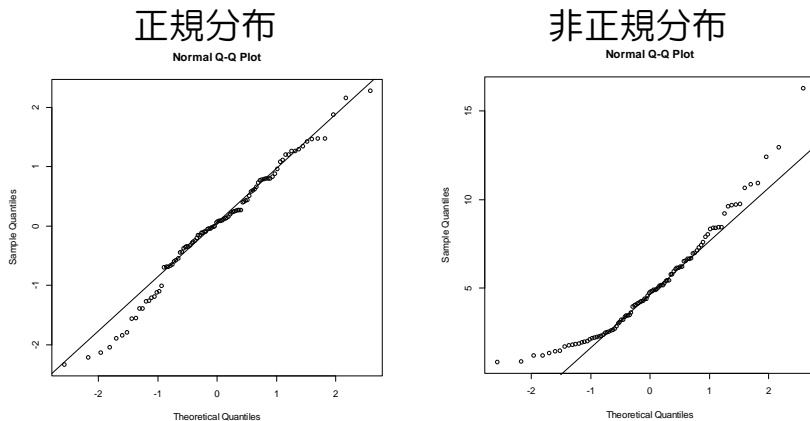
標本分布の正規性の確認は, 適切なモデルを選択する上で重要.

Definition: n 個の標本を大きさ順に並べたとき, i 番目に小さな標本は $[100(i-.5)/n]$ 標本パーセント点 (sample percentile) であるという.

例えば標本が, 正規分布など特定の確率分布から抽出されたとする. このとき, その特定の分布の理論上の $[100(i-.5)/n]$ パーセント点は, データの $[100(i-.5)/n]$ 標本パーセント点の近くにあるはずである.

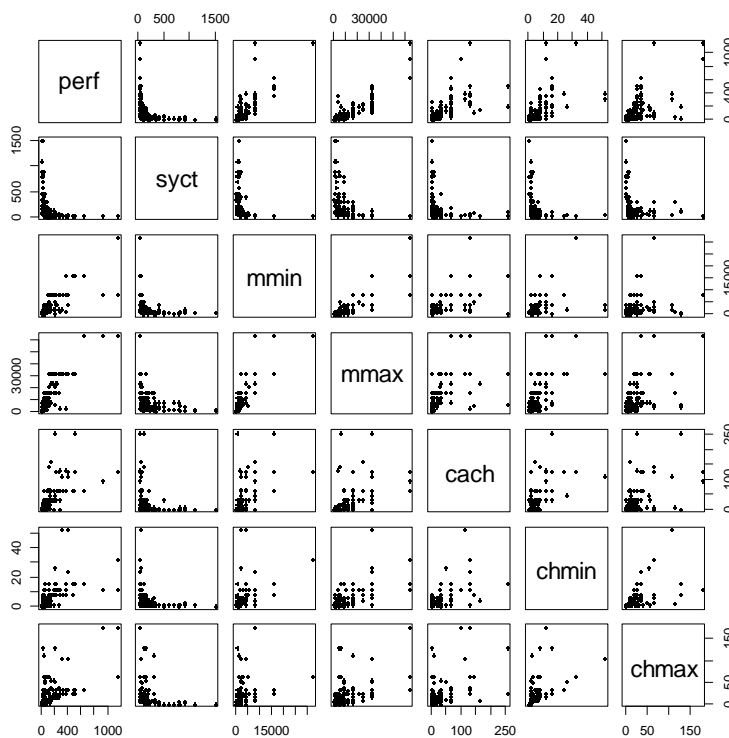
QQ-norm plot (Normal probability plot)

Definition: n 個の標本が得られたとき、標準正規分布の $[100(i-.5)/n]$ パーセント点と、 i 番目に小さな観測値 = $[100(i-.5)/n]$ 標本パーセント点のプロットを、**QQ-norm plot** という。



- 元のデータが**正規分布**から得られた場合、QQ-norm plotは**直線**上にプロットされる。
- 元データが正規分布に従わない場合、直線から外れる。(右図)

CPUデータ

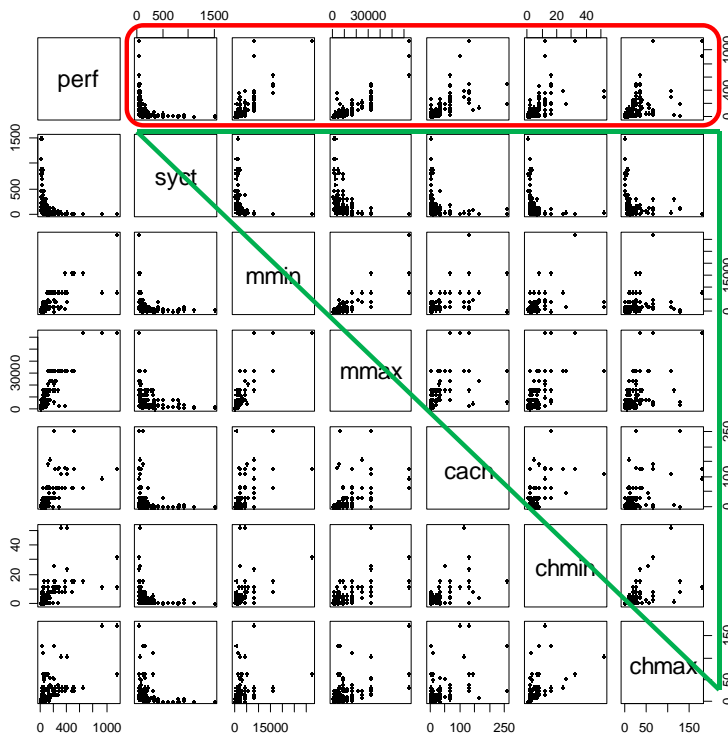


209のコンピュータのCPUの持つ、性能と各種特性値。

- 'name' Manufacturer and model
- 'syct' cycle time in nanoseconds
- 'mmin' minimum main memory in kilobytes
- 'mmax' maximum main memory in kilobytes
- 'cach' cache size in kilobytes
- 'chmin' minimum number of channels
- 'chmax' maximum number of channels
- 'perf' published performance on a benchmark mix relative to an IBM 370/158-3

P. Ein-Dor and J. Feldmesser (1987) Attributes of the performance of central processing units: a relative performance prediction model. *Comm. ACM.* **30**, 308–317.

CPUデータ（続き）



- y と x の間に、非線形な関係が存在しないか？
 - perf と syct の間に、明らかな非線形関係がある。
- 誤差項の分散は一定か？
 - mmin 等に、明らかな分差増大傾向がある。
- x 同士の間、線形な関係が存在しないか？
 - 例えば、mmin と mmax の間に明らかに線形関係がある。

予備的な視覚的要約の段階でも、線形回帰モデルを当てはめるのは不適切であることがわかる。

でも、とにかく回帰モデルを当てはめてみる。

P. Ein-Dor and J. Feldmesser (1987) Attributes of the performance of central processing units: a relative performance prediction model. *Comm. ACM*, 30, 308–317.

CPUデータ（元データの回帰分析）

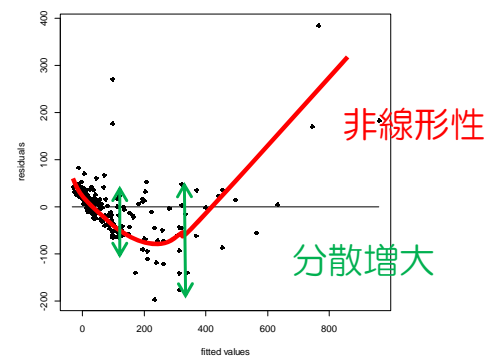
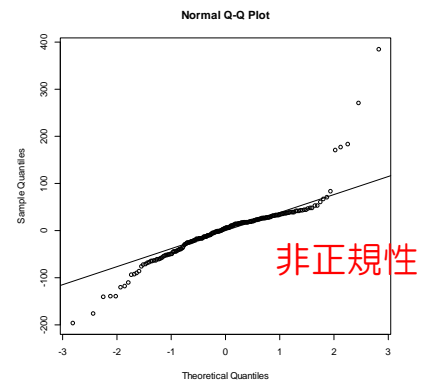
```

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) -5.590e+01  8.045e+00  -6.948 4.99e-11 ***
syct         4.886e-02  1.752e-02   2.789 0.00579 **
mmin        1.529e-02  1.827e-03   8.371 9.42e-15 ***
mmax        5.571e-03  6.418e-04   8.680 1.33e-15 ***
cach        6.412e-01  1.396e-01   4.594 7.64e-06 ***
chmin       -2.701e-01  8.557e-01  -0.316 0.75263
chmax       1.483e+00  2.201e-01   6.738 1.64e-10 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 59.99 on 202 degrees of freedom
Multiple R-squared:  0.8649,    Adjusted R-squared:  0.8609
F-statistic: 215.5 on 6 and 202 DF,  p-value: < 2.2e-16
    
```

- 説明変数の有意性検定は、chminを除き、ほとんどが強く有意。
- 決定係数： $R^2=0.8649$ 被説明変数の変動の86.49%が説明できた。
- Model utility test: $p\text{-value} < 2.2 \times 10^{-16}$

回帰分析は、成功しているとはいえない。
回帰診断の結果、モデルの仮定は破綻している。



変数変換

線形回帰モデルの仮定（線形性, 正規性, 等分散性）が満たされないとき，変数に何らかの変換を施すことで，モデルを改善できる場合がある。

例えば，攪乱項の分散が説明変数の値とともに大きくなる場合，**logarithmic/power** 変換が有効であることが多い。

被説明変数の予測値を得るには，まず変換された変数に対して線形回帰モデルを当てはめ，次にもとのモデルに逆変換する。最もよい変換を選ぶため，いくつかの変換を試してみる必要がある。

Box-Cox変換

対数変換，冪変換を組み合わせたBox-Cox変換により，**分散の安定化**と**正規性の改善**を同時に達成できる場合がある。

$$\text{Box-Cox変換} : y^{(\lambda)} = \begin{cases} (y^\lambda - 1)/\lambda : \lambda \neq 0 \\ \log(y) : \lambda = 0 \end{cases}$$

Box-Cox変換は，パラメター λ によって特徴付けられる。パラメター λ は，モデルの適合度を最適化するように，ソフトウェアにより自動的に選択される。

（統計解析ソフトRなどが，Box-Cox変換を実装している）

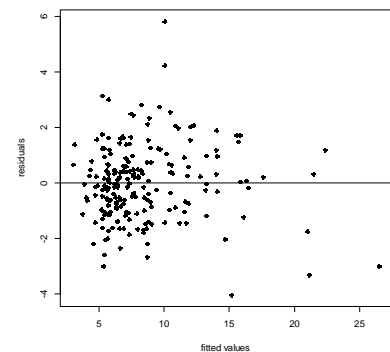
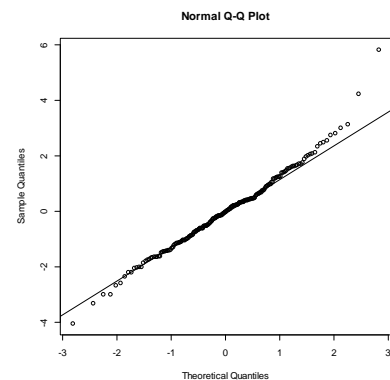
CPUデータ (Box-Cox変換後)

```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  5.214e+00  1.843e-01  28.284 < 2e-16 ***
sycy         -1.681e-03  4.014e-04  -4.187 4.21e-05 ***
mmin         1.868e-04  4.186e-05   4.463 1.34e-05 ***
mmax         1.607e-04  1.471e-05  10.924 < 2e-16 ***
cach         2.792e-02  3.198e-03   8.731 9.56e-16 ***
chmin        2.774e-02  1.961e-02   1.415  0.159
chmax        8.330e-03  5.042e-03   1.652  0.100
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.375 on 202 degrees of freedom
Multiple R-squared:  0.8821,    Adjusted R-squared:  0.8786
F-statistic: 251.8 on 6 and 202 DF,  p-value: < 2.2e-16
```

- chmin, chmax は有意ではない。
- 決定係数: $R^2=0.8821$ 被説明変数の変動の 88.21 (> 86.49) %が説明できた。
- Model utility test: p-value < 2.2×10^{-16}

回帰分析は、成功している。



変数選択

多変量回帰モデルを考えると、多数の説明変数の候補が考えられる場合がある。その時、被説明変数の変動を説明する、説明変数の最適な組み合わせを探索する必要がある。

- モデル選択基準 (赤池の情報量基準 (AIC), ベイズ情報量基準 (BIC), etc.) を最小にする、説明変数の組み合わせを求める。
- 当初のモデルに新たな変数を投入もしくは除去したとき、モデルの適合度の差が有意なものか検定する。(ワルド検定, 尤度比検定, ほか)

ステップワイズ法

説明変数の候補が k 個あるとき、可能な変数の組み合わせは 2^k 通り。総当たりの探索は無理。

- **変数増加法 (forward selection)**：定数項のみのモデルから出発して、有用な変数を加えていく方法。
- **変数減少法 (backward elimination)**：すべての変数を用いたモデルから出発して、不要な変数を除去していく方法。
- **変数増減法 (stepwise procedure)**：増加法と減少法の組み合わせ。

変数選択の実際（私の場合）

医学統計の性質上、説明変数は取りこぼしなく広めに選択したい。しかし、候補の数がサンプル数より大きければ、最初に全ての変数を用いたモデルは推定できない。また、欠測が多いとサンプル数が減る。

1. まず**単変量解析**で、変数ごとの p 値を出す。
2. ある程度 p 値が小さい変数に、**候補を絞る**。
($p < 0.2$ 程度で、あまりきつく絞らない)
3. 絞った候補から、**変数減少法**で選択する。（割合多めの変数が選択される）

変数選択の問題

- 変数選択の結果選ばれたモデルに、有意でない変数が含まれる。

変数選択は、被説明変数の変動をもっともよく説明する変数の組み合わせを探索する。個々の変数が有意でなくても、組み合わせ全体として最適であると解釈する。

- 変数選択の結果、興味のある変数がモデルから除かれてしまった。

上記に矛盾するようだが、回帰の目的は被説明変数を説明することだけではない。変数間の関係を推測するため、興味ある変数を強制投入してもよい。

Take Home Message

1. 回帰分析
2. 共分散と相関係数
3. 線形回帰モデル
 - 回帰係数の推定. 最小二乗推定量の性質
4. 回帰診断：回帰モデルの仮定の確認
 - 散布図：線形性の確認
 - QQ-normプロット：残差の正規性の確認
 - 残差プロット：等分散性, 独立性の確認
5. Box-Cox変換：分散の安定化と正規性の向上
6. 変数選択