

データ解析

Rによる多変量解析入門

(4) 回帰分析 (1)

lec20021031 下平英寿 shimo@is.titech.ac.jp

1

```
> ### 散布図
> ax <- "E09504"
> X2000$item[ax]
[1] "最終学歴が大学・大学院卒の割合"
> x <- X2000$item[ax]
> ay <- "A05203"
> X2000$item[ay]
[1] "合計特殊出生率"
> y <- X2000$item[ay]
> rbind(x,y)
      Hokkaido Aomori Iwate Miyagi Akita Yamagata Fukushima Ibaraki Tochigi Gumma
x      7.70  5.50  6.10  9.60  5.60  6.30  6.50  9.30  8.20  8.16
y      1.23  1.47  1.56  1.39  1.45  1.62  1.65  1.47  1.48  1.5
Saitama Chiba Tokyo Kanagawa Niigata Toyama Ishikawa Fukui Yamagashi Nagano
x     14.2  15.9  21.20  19.50  6.20  8.80  9.30  8.3  9.10  8.16
y      1.3  1.3  1.07  1.28  1.51  1.45  1.45  1.6  1.51  1.55
Gifu Shizuoka Aichi Mie Shiga Kyoto Osaka Hyogo Nara Wakayama Tottori
```

単回帰モデル (simple regression model)

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i; \quad i = 1, 2, \dots, n$$

$$y_1 = \beta_0 + \beta_1 x_1 + \epsilon_1$$

$$y_2 = \beta_0 + \beta_1 x_2 + \epsilon_2$$

⋮

$$y_n = \beta_0 + \beta_1 x_n + \epsilon_n$$

y_i 目的変数, 従属変数, 応答変数

x_i 説明変数, 独立変数, 予測変数

ϵ_i 誤差

β_k 回帰係数, 偏回帰係数

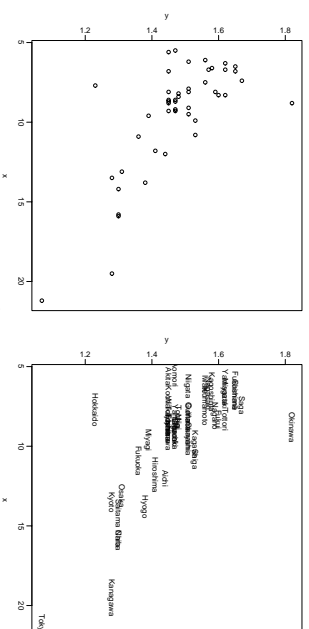
5

直線のあてはめ

```
x 8.70  9.20 12.00 8.40 10.80 13.50 13.10 13.80 15.8  8.10  8.30
y 1.47  1.44 1.44 1.48 1.53 1.28 1.31 1.38 1.3  1.45  1.62
Shimane Okayama Hiroshima Yamaguchi Tokushima Kagawa Ehime Kochi Fukuoka S
x  6.80  9.50  11.80  8.60  9.90  8.70  6.80  10.90 7
y  1.65  1.51  1.41  1.47  1.45  1.53  1.45  1.45  1.36 1
Nagasaki Kumamoto Oita Miyazaki Kagoshima Okinawa
x  6.70  7.50  7.90  6.70  6.60  8.80
y  1.57  1.56  1.51  1.62  1.58  1.82
> ## 散布図を点で描く
> plot(x,y)
> ## 散布図を県名で描く
> myplot(x,y)
```

2

散布図 (scatter plot)



$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

データ

3

最小二乗法 (least squares method)

$$\text{誤差 } \epsilon_i = y_i - (\beta_0 + \beta_1 x_i)$$

$$\text{誤差の二乗和 } S = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n \{y_i - (\beta_0 + \beta_1 x_i)\}^2$$

$$= A_{00}\beta_0^2 + 2A_{01}\beta_0\beta_1 + A_{11}\beta_1^2$$

Sが最小になるように β_0 と β_1 を調節する

$$\bar{x} = \frac{1}{n} \sum x_i, \quad \bar{y} = \frac{1}{n} \sum y_i;$$

$$S_{xx} = \sum (x_i - \bar{x})^2, \quad S_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y})$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}, \quad \hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$$

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

6

平均, 分散, 標準偏差, 共分散, 相関

```
> mymean1
function(x) sum(x)/length(x)
> mymean1(x)
[1] 9.540426
> mymean1(y)
[1] 1.472979
> myvar1
function(x,y=x) sum((x-mymean1(x))*(y-mymean1(y)))/(length(x)-1)
> sqrt(myvar1(x))
[1] 3.438950
> sqrt(myvar1(y))
[1] 0.1331380
> myvar1(x,y)/sqrt(myvar1(x)*myvar1(y))
[1] -0.729628
> mycorr1
function(x,y) myvar1(x,y)/sqrt(myvar1(x)*myvar1(y))
> mycorr1(x,y)
[1] -0.729628
```

4

数値例 1

```
> b1 <- myvar1(x,y)/myvar1(x,x)
> b0 <- mymean1(y) - b1 * mymean1(x)
> coef <- c(b0,b1)
> coef
[1] 1.74248324 -0.02824869
> plot(x,y)
> ## 散布図に回帰直線を描く
> abline(b0,b1)
> ## 予測
> pred <- b0 + b1 * x # この計算では、スカラ(長さ1のベクトル)と長さ4.7のベクトルを足し算している。長さの違うベクトルを足し算すると、短いほうが繰り返し用いられる。
> rbind(pred,y)
      Hokkaido Aomori Iwate Miyagi Akita Yamagata Fukushima Ibaraki Tochigi Gumma Saitama Chiba Tokyo Kanagawa Niigata Toyama
pred 1.524988 1.587115 1.570166 1.471286 1.584291 1.564516 1.558867 1.479770
y      1.230000 1.470000 1.560000 1.390000 1.450000 1.620000 1.650000 1.470000
      Tochigi Gumma Saitama Chiba Tokyo Kanagawa Niigata Toyama
pred 1.510844 1.513669 1.341352 1.293329 1.143611 1.191634 1.567341 1.493895
```

7

重回帰モデル (multiple regression model)

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i \quad i = 1, \dots, n$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad x_k = \begin{bmatrix} x_{1k} \\ \vdots \\ x_{nk} \end{bmatrix}, \quad \epsilon = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}$$

$$y = \beta_0 1_n + \beta_1 x_1 + \dots + \beta_p x_p + \epsilon$$

$$X = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix}$$

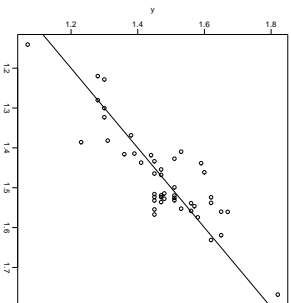
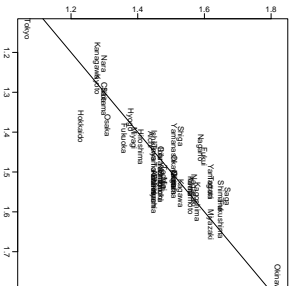
$$y = X\beta + \epsilon$$

14

数値例2

```
> ## 多変量で線形
> ax <- c("E09504", "A0410302", "C01301", "B02101"); x <- X2000$x[,ax]
> ay <- "A05203"; y <- X2000$y[,ay]
> X2000$item[c(ax,ay)]
      E09504      A0410302
"最終学歴が大学・大学院卒の者の割合" "未婚者割合 [20~24歳・女]"
      C01301      B02101
"県民1人当たり県民所得" "年平均気温"
      A05203
"合計特殊出生率"
```

```
> abline(0,1)
> plot(pred,y)
> abline(0,1)
```



射影

$$\|e\|^2 = \|y - X\beta\|^2$$

を最小にするには

$$\hat{\beta} = (X'X)^{-1}X'y$$

$$\hat{y} = X\hat{\beta}$$

ハット行列 (射影行列)

$$H = X(X'X)^{-1}X'$$

$1_n, x_1, \dots, x_p$ の張る空間への y の射影

$$\hat{y} = Hy$$

もし $X'X$ が退化している場合には

$$\hat{y} = XX^+y$$

14

```
A0410302 -2.727500e-02 7.438516e-03 -3.6667262 6.850351e-04
C01301 1.228960e-05 4.635656e-05 0.2651098 7.922217e-01
B02101 2.012511e-02 5.098161e-03 3.9475245 2.951848e-04
> f$summary
Mean Sum Sq R Squared F-value Df 1 Df 2 Pr(>F)
Y 0.07061925 0.7431179 30.37478 4 42 6.678249e-12
> x1 <- cbind(rep(1, length(y)), x)
> coef <- solve(t(x1) %*% x1) %*% (t(x1) %*% y)
> coef
      [1]
      3.636118e+00
E09504 -1.655659e-02
A0410302 -2.727500e-02
C01301 1.228960e-05
B02101 2.012511e-02
> pred <- x1 %*% coef
> t(cbind(pred,y))
```

```
> ## 必ず単回帰, それからダミー変数 (1)
> ## 単回帰
> ax <- "C01301"; ay <- "C04602"
> X2000$item[c(ax,ay)]
      C01301
"県民1人当たり県民所得"
```

数値例3

```
> ## 個人預貯金残高(人口1人当たり)
> x <- X2000$x[,ax, drop=F]
> y <- X2000$y[,ay, drop=F]
> f <- mylsfit(x,y)
> f$summary
      Estimate Std.Err t-value Pr(>|t|)
Intercept 26.9369511 102.88756993 0.2618096 7.946630e-01
```

```
> f$summary
Mean Sum Sq R Squared F-value Df 1 Df 2 Pr(>F)
C04602 90.60172 0.3606658 25.38572 1 45 8.093436e-06
> myplot(x,y)
> abline(f)
> ## 中心化してダミー変数を使う
> y <- x - mean(x)
> x <- y - mean(y)
> x <- cbind(x, nhomregion)
      C01301 tohoku kanto shinsetsu tokai kinki chugoku shikoku kyushu
```

18

導出

$$\|e\|^2 = (y - X\beta)'(y - X\beta) = y'y - 2\beta'X'y + \beta'(X'X)\beta$$

これを β で微分して

$$\frac{\partial \|e\|^2}{\partial \beta} = -2X'y + 2X'X\beta = 0$$

すなわち正規方程式 (normal equation)

$$X'X\beta = X'y$$

これを解いて

$$\hat{\beta} = (X'X)^{-1}X'y$$

16

```
Hokkaido Aomori Iwate Miyagi Akita Yamagata Fukushima Ibaraki
1.385854 1.518977 1.538630 1.414345 1.518523 1.524002 1.619096 1.467897
Y 1.230000 1.470000 1.560000 1.390000 1.450000 1.620000 1.650000 1.470000
Tochigi Gunma Saitama Chiba Tokyo Kanagawa Niigata Toyama
1.513416 1.531548 1.323061 1.300365 1.140643 1.219962 1.525388 1.464334
Y 1.480000 1.510000 1.300000 1.070000 1.280000 1.510000 1.450000
Ishikawa Fuku Yamanaishi Nagano Gifu Shizuoka Aichi Mie
1.433561 1.461327 1.427325 1.438374 1.453923 1.535835 1.418394 1.527675
Y 1.450000 1.600000 1.510000 1.590000 1.470000 1.440000 1.480000
Shiga Kyoto Osaka Hyogo Nara Wakayama Tottori Shimane
1.403378 1.280022 1.381912 1.368422 1.228403 1.554402 1.538035 1.559977
Y 1.530000 1.280000 1.310000 1.380000 1.450000 1.620000 1.650000
Okayama Hiroshima Yamaguchi Tokushima Kagawa Ehime Kochi Fukuoka
1.499213 1.436992 1.52368 1.523004 1.52604 1.531791 1.567354 1.415729
Y 1.510000 1.410000 1.47000 1.450000 1.530000 1.450000 1.450000 1.360000
Saga Nagasaki Kumamoto Oita Miyazaki Kagoshima Okinawa
1.560404 1.546608 1.558161 1.519602 1.631219 1.574310 1.768123
Y 1.670000 1.570000 1.560000 1.510000 1.620000 1.580000 1.820000
> myplot(pred,y)
```

```
C01301 0.1804065 0.03580613 5.0384246 8.093436e-06
```

```
> f$summary
Mean Sum Sq R Squared F-value Df 1 Df 2 Pr(>F)
C04602 90.60172 0.3606658 25.38572 1 45 8.093436e-06
> myplot(x,y)
> abline(f)
> ## 中心化してダミー変数を使う
> y <- x - mean(x)
> x <- y - mean(y)
> x <- cbind(x, nhomregion)
      C01301 tohoku kanto shinsetsu tokai kinki chugoku shikoku kyushu
Hokkaido -118.659574 1 0 0 0 0 0 0
Aomori -360.659574 1 0 0 0 0 0 0
Iwate -230.659574 1 0 0 0 0 0 0
Miyagi -73.659574 1 0 0 0 0 0 0
Akita -275.659574 1 0 0 0 0 0 0
Yamagata -220.659574 1 0 0 0 0 0 0
```



```

Y 28.30 29.30 30.50 29.50 31.80 32.60 32.10 30.60 31.90 32.60
Saitama Chiba Tokyo Kanagawa Niigata Toyama Ishikawa Fukui Yamaguchi Nagano
Y 33.01 33.1 33.05 33.05 31.84 32.27 32.27 32.31 32.94 31.83
> mycor1(f1$pred, y)
[1] 0.686967
> mycor1(f2$pred, y)^2
[1] 0.4719236
> f2$summary
Mean Sum Sq R Squared F-value Df 1 Df 2 Pr(>F)
Y 1.070696 0.4719236 19.66064 2 44 7.931007e-07
> ## 数値例4 (3次式)
> mycor1(f3$pred, y)
[1] 0.5437978
> f3$summary
Mean Sum Sq R Squared F-value Df 1 Df 2 Pr(>F)
Y 1.006674 0.5437978 17.08548 3 43 1.877689e-07
> f1$summary
Mean Sum Sq R Squared F-value Df 1 Df 2 Pr(>F)

```

```

Y 1.304014 0.1988971 11.17256 1 45 0.001678815
> ## 数値例4 (2次式)
> mycor1(f2$pred, y)
[1] 0.686967
> mycor1(f2$pred, y)^2
[1] 0.4719236
> f2$summary
Mean Sum Sq R Squared F-value Df 1 Df 2 Pr(>F)
Y 1.070696 0.4719236 19.66064 2 44 7.931007e-07
> ## 数値例4 (3次式)
> mycor1(f3$pred, y)
[1] 0.5437978
> f3$summary
Mean Sum Sq R Squared F-value Df 1 Df 2 Pr(>F)
Y 1.006674 0.5437978 17.08548 3 43 1.877689e-07

```

残差 (residual)

$$e_i = y_i - \hat{y}_i \quad i = 1, \dots, n$$

$$e = \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix}, \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \hat{y} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix}$$

$$e = y - \hat{y}$$

ハット行列を用いると

$$e = (I_n - H)y$$

$HX = X$ なので $X'e = 0$. とくに

$$1_n'e = 0, \quad \hat{y}'e = 0$$

ピタゴラスの定理

$$\|y\|^2 = \|\hat{y}\|^2 + 2\hat{y}'e + \|e\|^2 = \|\hat{y}\|^2 + \|e\|^2$$

```

> ### 残差
> e <- y-f1$pred
> t(round(e,2))
Hokkaido Aomori Iwate Miyagi Akita Yamagata Fukushima Ibaraki Tochigi Gun
Y -2.25 -1.72 -0.69 -2.36 0.65 1.18 0.17 -2.39 -1.04 -0.
Saitama Chiba Tokyo Kanagawa Niigata Toyama Ishikawa Fukui Yamaguchi Nag
Y 0.19 -1.5 -0.65 -1.75 1.06 0.83 0.23 1.29 0.36 0
Gifu Shizuoka Aichi Mie Shiga Kyoto Osaka Hyogo Nara Wakayama Tottori
Y 2.28 -2.08 1.95 -0.31 0.93 2.29 2.18 0.72 1.58 -0.04 1.2
Shimane Okayama Hiroshima Yamaguchi Tokushima Kagawa Ehime Kochi Fukuoka
Y 0.81 1.23 0.6 0.02 0.24 0.99 0.01 -0.83 -0.12
Saga Nagasaki Kumamoto Oita Miyazaki Kagoshima Okinawa

```

```

Y 0.23 -0.89 0.51 -0.57 -1.76 -0.68 -2.38
> sum(e)
[1] 9.947598e-14
> sum(f1$pred * e)
[1] 3.312767e-12
> mymean1(e)
[1] 2.116510e-15
> mycor1(f1$pred, e)
[1] 1.675907e-15
> ## ピタゴラスの定理
> sum(y^2)
[1] 49980.06
> sum(f1$pred^2)
[1] 49903.54
> sum(e^2)
[1] 76.52093
> sum(f1$pred^2) + sum(e^2) - sum(y^2)
[1] -7.275958e-12

```

重相関係数と残差の関係

まず中心化 $y \leftarrow y - \bar{y}1_n, \hat{y} \leftarrow \hat{y} - \bar{y}1_n$ と中心化しておくとき,

$$R = \frac{\hat{y}'y}{\|\hat{y}\|\|y\|}$$

$e'y = 0$ なので $y'y = (\hat{y} + e)'y = \|\hat{y}\|^2$, すなわち

$$R^2 = \frac{\|\hat{y}\|^2}{\|y\|^2} = \frac{S_{\hat{y}\hat{y}}}{S_{yy}}$$

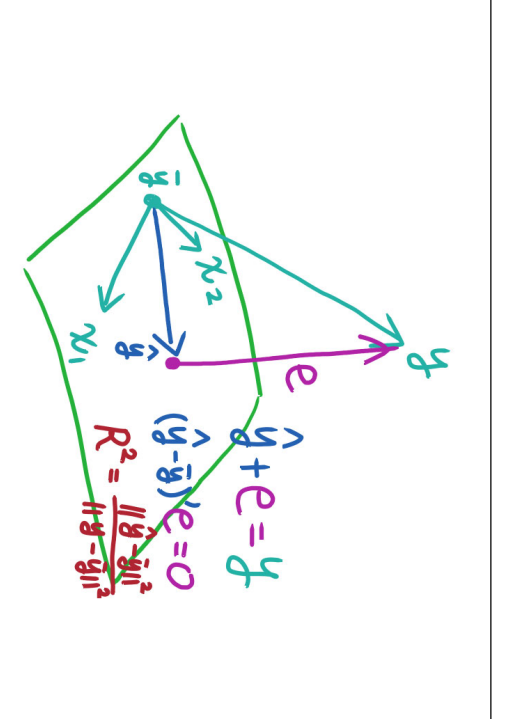
ピタゴラスの定理 $\|y\|^2 = \|\hat{y}\|^2 + \|e\|^2$ もしくは $S_{yy} = S_{\hat{y}\hat{y}} + S_{ee}$ より

$$R^2 = 1 - \frac{\|e\|^2}{\|y\|^2} = 1 - \frac{S_{ee}}{S_{yy}}$$

```

> ### 重相関係数と残差の関係
> e <- y-f1$pred # 残差
> cor(f1$pred, y)^2 # R^2
[1]
Y 0.1988971
> sum(f1$pred^2)/sum(y^2) # 中心化しないピタゴ
[1] 0.998469
> sum((f1$pred-mymean1(f1$pred))^2)/sum((y-mymean1(y))^2) # R^2
[1] 0.1988971
> myvar1(f1$pred)/myvar1(y) # これでもOK
[1] 0.1988971
> 1-sum(e^2)/sum(y^2) # 中心化しないピタゴ
[1] 0.998469
> 1-sum(e^2)/sum((y-mymean1(y))^2) # R^2
[1] 0.1988971
> 1-myvar1(e)/myvar1(y) # これでもOK
[1] 0.1988971

```



部分回帰 (Subset Regression)

モデル (k)

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + e_i \quad i = 1, \dots, n$$

すなわち

$$\beta_{k+1} = \dots = \beta_p = 0$$

$$X^{(k)} = [1_n, x_{1\cdot}, \dots, x_{k\cdot}], \quad \beta^{(k)} = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_k \end{bmatrix}$$

$$y = X^{(k)}\beta^{(k)} + e$$

$$\hat{\beta}^{(k)} = (X^{(k)})^{-1}X^{(k)'}y$$

一般に $\hat{\beta}^{(k)}$ と β の対応する要素は一致しない: $\hat{\beta}_i^{(k)} \neq \beta_i$

直交化とQR分解

```
> ## 部分回帰
> t(f1$coef)
Intercept      X
Y 31.62375 -0.01808129
> t(f2$coef)
Intercept      X      x^2      x^3
Y 31.62375 0.1227429 -0.002052218 7.450801e-06

X = [x0, ..., xn], Q = [q0, ..., qn], Q'Q = I_{p+1}
X = QR, r_{ij} = 0, i > j, γ = Rβ
y = Xβ + ε = Qγ + ε
γ̂ = Q'y, β̂ = R^{-1}γ̂ = R^{-1}Q'y
c.f. (X'X)^{-1}X' = (RR')^{-1}R'Q' = R^{-1}Q'
γ_k = r_{kk}β_k + ... + r_{kp}β_p に注意すると
γ_{k+1} = ... = γ_p = 0 ⇔ β_{k+1} = ... = β_p = 0
モデル(k)
γ̂^{(k)} = Q^{(k)}y, γ̂_k^{(k)} = γ̂_k
```

```
> ## 直交化とQR分解
> xx <- cbind(1,x,x^2,x^3)
> dimnames(xx)[2,] <- c("1","x","x^2","x^3")
> q <- qr.Q(qr(xx))
> R <- qr.R(qr(xx))
> cbind(xx,q)
      1      x      x^2      x^3
Hokkaido 1 145 21025 3048625 -0.145865 -0.46469866 0.58635783 0.54555490
Aomori 1 119 14161 1685159 -0.145865 -0.35684268 0.16939012 -0.16646877
Iwate 1 110 12100 1331000 -0.145865 -0.31950791 0.05955370 -0.25620527
Miyagi 1 7 5329 389017 -0.145865 -0.16602055 -0.20035745 -0.12673425
Akita 1 112 12544 1404928 -0.145865 -0.32780453 0.08229878 -0.24191600
Yamagata 1 97 9409 912673 -0.145865 -0.26557992 -0.06654976 -0.28255814
Fukushima 1 69 4761 328509 -0.145865 -0.14942732 -0.21017241 -0.08803184
Ibaraki 1 10 1000 10000 -0.145865 0.09532280 0.05944675 -0.01420506
Tochigi 1 13 169 2197 -0.145865 0.08287788 0.02700782 0.03316958
Gunma 1 8 64 512 -0.145865 0.10361942 0.08218756 -0.05056613
Saitama 1 9 81 729 -0.145865 0.09947111 0.07070567 -0.03189642
```

Chiba	1	4	16	64	-0.145865	0.12021265	0.13034481	-0.13536963
Tokyo	1	7	49	343	-0.145865	0.10776772	0.09389924	-0.07023124
Kanagawa	1	7	49	343	-0.145865	0.10776772	0.09389924	-0.07023124
Miagata	1	74	5476	405224	-0.145865	-0.17016885	-0.19734628	-0.13616848
Toiyama	1	50	2500	125000	-0.145865	-0.07060948	-0.20807451	0.08707592
Ishikawa	1	50	2500	125000	-0.145865	-0.07060948	-0.20807451	0.08707592
Fukui	1	48	2304	110592	-0.145865	-0.06231287	-0.20317131	0.10168360
Yamanashi	1	13	169	2197	-0.145865	0.08287788	0.02700782	0.03316958
Nagano	1	71	5041	357911	-0.145865	-0.15772393	-0.20571087	-0.10754201
Gifu	1	25	625	15625	-0.145865	0.03309820	-0.08269060	0.14536604
Shizuoka	1	5	25	125	-0.145865	0.11606434	0.11797104	-0.11261593
Aichi	1	18	324	5832	-0.145865	0.06213635	-0.02259767	0.09415285
Mie	1	20	400	8000	-0.145865	0.05383973	-0.04087907	0.11265321
Shiga	1	39	1521	59319	-0.145865	-0.02497810	-0.17006988	0.14998770
Kyoto	1	37	1369	50653	-0.145865	-0.01668149	-0.16026134	0.15599117
Osaka	1	14	196	2744	-0.145865	0.07872957	0.01664078	0.04711818
Hyogo	1	22	484	10648	-0.145865	-0.04564312	-0.05826859	0.12799090
Nara	1	25	625	15625	-0.145865	0.03309820	-0.08269060	0.14536604
Nakayama	1	13	169	2197	-0.145865	0.08287788	0.02700782	0.03316958

Tottori	1	43	1849	79507	-0.145865	-0.04157133	-0.18701103	0.13242784
Shimane	1	38	1444	54872	-0.145865	-0.02082980	-0.16527710	0.15321596
Okayama	1	17	289	4913	-0.145865	0.06628465	-0.01312251	0.08367402
Hiroshima	1	21	441	9261	-0.145865	0.04969142	-0.04968832	0.12070886
Yamaguchi	1	22	484	10648	-0.145865	0.04554312	-0.05826859	0.12799090
Tokushima	1	12	144	1728	-0.145865	0.08702619	0.03759783	0.01831088
Kagawa	1	9	81	729	-0.145865	0.09947111	0.07070567	-0.03189642
Ehime	1	10	100	1000	-0.145865	0.09532280	0.05944675	-0.01420506
Kochi	1	8	64	512	-0.145865	0.10361942	0.08218756	-0.05056613
Fukuoka	1	14	196	2744	-0.145865	0.07872957	0.01664078	0.04711818
Saga	1	17	289	4913	-0.145865	0.06628465	-0.01312251	0.08367402
Nagasaki	1	10	100	1000	-0.145865	0.09532280	0.05944675	-0.01420506
Kumamoto	1	10	100	1000	-0.145865	0.09532280	0.05944675	-0.01420506
Oita	1	6	36	216	-0.145865	0.11191603	0.10582024	-0.09090882
Miyazaki	1	1	1	1	-0.145865	0.13265757	0.16880393	-0.21008090
Kagoshima	1	5	25	125	-0.145865	0.11606434	0.11797104	-0.11261593
Okinawa	1	0	0	0	-0.145865	0.13690587	0.18206958	-0.23719159

```
[1,] -223.3485 4.368715 -5.106775 2.620177
> t(tr %*% g3) # betaに等しいV\X
      1      x      x^2      x^3
[1,] 31.62375 0.1227429 -0.002052218 7.450801e-06

> ## 平方和の分解
> sum((f1$pred-mean(f1$pred))^2) # 予測値の平方和 (1次)
[1] 18.99839
> sum(g3[2,]^2) # g[2]^2
[1] 18.99839
> sum((f2$pred-mean(f2$pred))^2) # 予測値の平方和 (2次)
[1] 45.07754
> sum(g3[2:3,]^2) # g[2]^2+g[3]^2
[1] 45.07754
> sum((f3$pred-mean(f3$pred))^2) # 予測値の平方和 (3次)
[1] 51.94287
> sum(g3[2:4,]^2) # g[2]^2+g[3]^2+g[4]^2
[1] 51.94287
> cor(f1$pred,y)^2 # R^2 (1次)
[1,]
Y 0.1988971
> sum(g3[2,]^2)/sum((y-mean(y))^2)
[1] 0.1988971
```

```
> t(f3$coef) # beta
Intercept      X      x^2      x^3
Y 31.62375 0.1227429 -0.002052218 7.450801e-06
> t(tr[1:2,1:2]) %*% g3[1:2]) # gamma (1次)
      1      x
[1,] 33.17502 -0.01808129
> t(f1$coef)
Intercept      X
Y 33.17502 -0.01808129
[1,] 32.24523 0.05023272 -0.0005693289
> t(f2$coef)
Intercept      X      x^2      x^3
Y 32.24523 0.05023272 -0.0005693289
```

```
平方和の分解
γ̂ = Qγ̂ = γ̂_0q_0 + ... + γ̂_pq_p
||γ̂||^2 = ||γ̂_0q_0||^2 + ... + ||γ̂_pq_p||^2 = γ̂_0^2 + ... + γ̂_p^2
q_0 = 1/√n 1_n, γ̂_0 = √n γ̄
モデル(k)
||Q^{(k)}γ̂^{(k)} - γ̄_{1:n}||^2 = γ̂_1^2 + ... + γ̂_k^2
||e||^2 = ||y||^2 - γ̂_0^2 - ... - γ̂_k^2 = ||y - γ̄_{1:n}||^2 - γ̂_1^2 - ... - γ̂_k^2
R^2 = (γ̂_1^2 + ... + γ̂_k^2) / ||y - γ̄_{1:n}||^2
```

```
> ## 平方和の分解
> sum((f1$pred-mean(f1$pred))^2) # 予測値の平方和 (1次)
[1] 18.99839
> sum(g3[2,]^2) # g[2]^2
[1] 18.99839
> sum((f2$pred-mean(f2$pred))^2) # 予測値の平方和 (2次)
[1] 45.07754
> sum(g3[2:3,]^2) # g[2]^2+g[3]^2
[1] 45.07754
> sum((f3$pred-mean(f3$pred))^2) # 予測値の平方和 (3次)
[1] 51.94287
> sum(g3[2:4,]^2) # g[2]^2+g[3]^2+g[4]^2
[1] 51.94287
> cor(f1$pred,y)^2 # R^2 (1次)
[1,]
Y 0.1988971
> sum(g3[2,]^2)/sum((y-mean(y))^2)
[1] 0.1988971
```

第 4 回 課題

```
> cor(f2$pred,y)^2 # R^2 (2次)
[1]
y 0.4719236
> sum(g3[2:3]^2)/sum((y-mean(y))^2)
[1] 0.4719236
> cor(f3$pred,y)^2 # R^2 (3次)
[1]
y 0.5437978
> sum(g3[2:4]^2)/sum((y-mean(y))^2)
[1] 0.5437978
```

4. `kaiki2`, `jyusokansq` を使い, `X2000` から適当な項目を選んで重回帰分析する. 係数 β と重相関係数 R を計算する. `myfunc20020919.R` にある `my1stfit` をつかって同じ分析をして, 結果が同じになるかどうか確認する.
5. 上で得られた結果について, `pred` を X 軸, Y を Y 軸とするプロットをする. X 軸 Y 軸となる直線を描く (`abline(0,1)` をつかう). さらに県名を使ったプロットをする. `myfunc20020919.R` の `myplot` 関数を参考にせよ. プロットは `myfunc20020919.R` にある `psinit` 関数などを使い `eps` ファイルとして出力し, それをプリンタで印刷する.


```
psinit("ファイル名") # これ以後のプロットの結果をファイルに eps 形式で書き出す
                        ここでプロットをおこなう...
dev.off() # ファイルをクローズする
```

1. 直線当てはめ (単回帰) の関数 `kaiki1` を作れ.

```
kaiki1 <- function(x,y) {
  # x,y は同じ長さの実数ベクトル
  # y = coef[1] + coef[2]*x + resid の形の単回帰分析を行う
  # 以下の coef, pred, resid を計算する
  # coef (係数) は 2 次元ベクトル
  # resid (残差) は y と同じ長さのベクトル
  # pred (予測値) = coef[1] + coef[2]*x は y と同じ長さのベクトル
  # 次の行は結果をリストとして返す.
  list(coef,pred,resid)
}
```

2. 重回帰分析の関数 `kaiki2` を作れ.

```
kaiki2 <- function(x,y) {
  # x は n * p 次元の行列
  # y は長さ n のベクトル
```

```
# y = coef[1] + coef[2]*x[,1] + ... + coef[p+1]*x[,p] + resid
# の形の重回帰分析を行う
# 以下の coef, pred, resid を計算する
# coef (係数) は p+1 次元ベクトル
# resid (残差) は y と同じ長さのベクトル
# pred (予測値) は y と同じ長さのベクトル
# 次の行は結果をリストとして返す.
list(coef,pred,resid)
}
```

3. `kaiki2` の返す値から重相関係数の二乗 R^2 を計算する関数 `jyusokansq` を作れ.

```
jyusokansq <- function(kout) {
  # kout$pred は予測値, kout$resid は残差
  # これらから重相関係数の二乗を計算し rsq に代入
  rsq
}
```