

回帰分析の基礎

- 直線のあてはめ
- 回帰分析
- あてはまりのよさ
- 部分回帰

線形代数の知識，とくに「射影」を前提とする

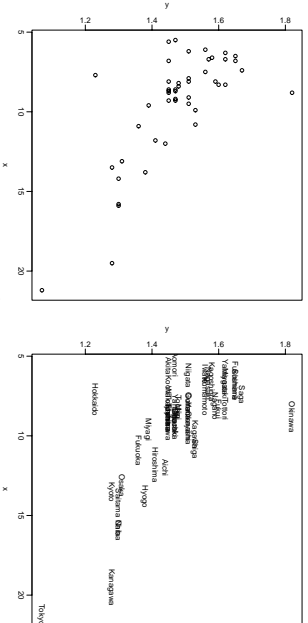
直線のあてはめ

データ解析

Rによる多変量解析入門

回帰分析 (1)

散布図 (scatter plot)



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

セッションファイル

```

edu: ~shimo/class/gakuhubu200209/note20021020.Rr
> ## 散布図
> ax <- "E09504"
> X2000$item[ax]
E09504
"最終学歴が大学・大学院卒の者の割合 "
> x <- X2000$xl[ax]
> ay <- "A05203"
> X2000$item[ay]
A05203
"合計特殊出生率 "
> y <- X2000$xl[ay]
> rbind(x,y)
Hokkaido Aomori Iwate Miyagi Akita Yamagata Fukushima Ibaraki Tochigi Gunma
x 7.70 5.50 6.10 9.60 5.60 6.30 6.50 9.30 8.20 8.10
y 1.23 1.47 1.56 1.39 1.45 1.62 1.65 1.47 1.48 1.57
Saitama Chiba Tokyo Kanagawa Niigata Toyama Ishikawa Fukui Yamagashi Nagano
    
```

単回帰モデル (simple regression model)

$$y = \beta_0 + \beta_1 x + \epsilon$$

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

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

⋮

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

y 目的変数, 従属変数, 応答変数

x 説明変数, 独立変数, 予測変数

ϵ 誤差

β_0, β_1 回帰係数, 偏回帰係数

平均，分散，標準偏差，共分散，相関

$$\bar{x} = \frac{1}{n} \sum x_i$$

$$\bar{y} = \frac{1}{n} \sum y_i$$

$$s_{xy} = \frac{1}{n-1} \sum (x_i - \bar{x})(y_i - \bar{y})$$

$$s_x^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2 = s_{xx}$$

$$s_x = \sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2}$$

$$r_{xy} = \frac{s_{xy}}{s_x s_y} = \frac{s_{xy}}{\sqrt{s_{xx} s_{yy}}}$$

x_i の平均

y_i の平均

x_i と y_i の共分散

x_i の分散

x_i の標準偏差

相関係数

[1] 1.472979

```

> myvari
function(x,y=x) sum((x-mymeanl(x))*(y-mymeanl(y)))/(length(x)-1)
> sqrt(myvari(x))
[1] 3.438950
> sqrt(myvari(y))
[1] 0.1331380
> myvari(x,y)/sqrt(myvari(x)*myvari(y))
[1] -0.7296628
> mycor1
function(x,y) myvari(x,y)/sqrt(myvari(x)*myvari(y))
> mycor1(x,y)
[1] -0.7296628
    
```

```

> mymeanl
function(x) sum(x)/length(x)
> mymeanl(x)
[1] 9.540426
> mymeanl(y)
    
```

最小二乗法 (least squares method)

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

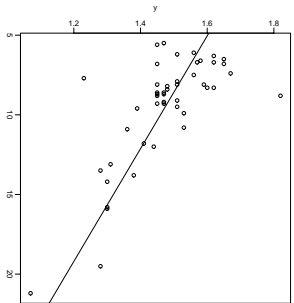
誤差の二乗和 $S = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n \{y_i - (\beta_0 + \beta_1 x_i)\}^2$
 $= \sum \{\beta_0^2 + 2\beta_0\beta_1 x_i + \beta_1^2 x_i^2 - 2y_i\beta_0 - 2y_i x_i \beta_1 + y_i^2\}$
 $= n \{ \beta_0^2 + 2\bar{x}\beta_0\beta_1 + \bar{x}^2\beta_1^2 - 2\bar{y}\beta_0 - 2\bar{y}\bar{x}\beta_1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

Sが最小になるように β_0 と β_1 を調節すると...

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

これで予測

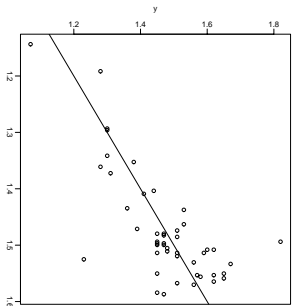
$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$



(x_i, y_i) の散布図

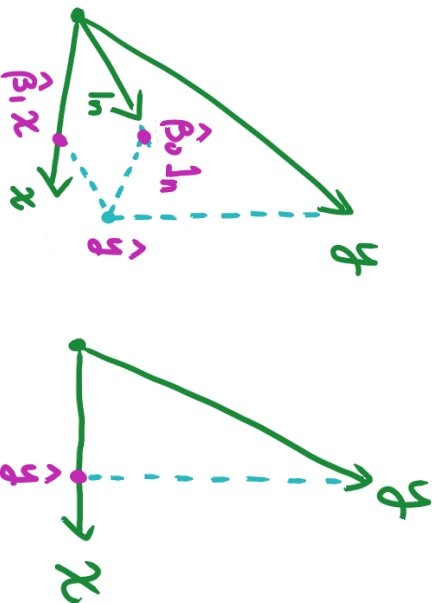
$y = \beta_0 + \beta_1 x$

$\beta_0 = 1.74, \beta_1 = -0.028$



(\hat{y}_i, y_i) の散布図

$\hat{y} = y$



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
> ## 散布図に回帰直線を描く
> plot(x,y)
> abline(b0,b1)
> ## 予測
> pred <- b0 + b1 * x # この計算では、xカラム(長さ1のベクトル)と長さ47のx
カラムを足し算している。長さの違うベクトルを足し算すると、短い方が繰り返し使
れる。
> rbind(pred,y)
Hokkaido Amori Iwate Miyagi Akita Yamagata Fukushima Ibaraki
pred 1.524968 1.587115 1.570166 1.471296 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 Gunma Saitama Gaiba Tokyo Kanagawa Niigata Toyama
pred 1.510844 1.513669 1.341352 1.293329 1.143611 1.191634 1.567341 1.493895
```

最小二乗推定量の導出

$$\frac{\partial S}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\frac{\partial S}{\partial \beta_1} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) x_i = 0$$

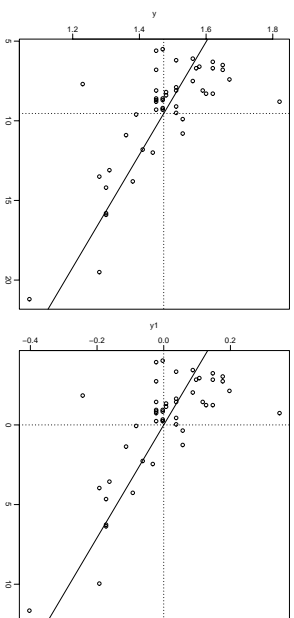
$$\sum_{i=1}^n y_i = n\beta_0 + \beta_1 \sum_{i=1}^n x_i$$

$$\sum_{i=1}^n x_i y_i = \beta_0 \sum_{i=1}^n x_i + \beta_1 \sum_{i=1}^n x_i^2$$

$$\begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix} = \begin{bmatrix} \sum 1 & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$$

$$\begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} = \begin{bmatrix} \sum 1 & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix}^{-1} \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}$$

中心化 (centering)



$\hat{x}_i \leftarrow x_i - \bar{x}, \quad \hat{y}_i \leftarrow y_i - \bar{y}$

```
y 1.480000 1.510000 1.300000 1.070000 1.280000 1.510000 1.450000
Ishikawa Fukui Yamanaishi Nagano Gifu Shizuoka Aichi Mik
pred 1.479770 1.508019 1.485420 1.513669 1.496720 1.482595 1.403499 1.50519
y 1.450000 1.600000 1.510000 1.590000 1.470000 1.470000 1.440000 1.480000
Shiga Kyoto Osaka Hyogo Nara Makayama Tottori Shimane
pred 1.437397 1.361126 1.372425 1.352651 1.296154 1.513669 1.508019 1.550392
Okayama Hiroshima Yamaguchi Tokushima Kagawa Ehime Kochi Fuku
pred 1.474121 1.409149 1.499545 1.499545 1.462821 1.496720 1.550392 1.434
y 1.510000 1.410000 1.470000 1.450000 1.530000 1.450000 1.450000 1.3600
Saga Nagasaki Kumamoto Oita Miyazaki Kagoshima Okinawa
pred 1.533443 1.532117 1.530618 1.519319 1.553217 1.556042 1.493895
y 1.670000 1.570000 1.560000 1.510000 1.620000 1.580000 1.820000
```

ベクトル表現

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad e = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$y = \beta_0 \mathbf{1}_n + \beta_1 x + e$$

$$= \begin{bmatrix} 1_n, x \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + e$$

$$= X\beta + e$$

$$X = \begin{bmatrix} 1_n, x \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}, \quad S = \|e\|^2 = \|y - X\beta\|^2 \text{の最小化}$$

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

$y_i = \beta_1 x_i + \epsilon_i$
 $y = \beta_1 x + e$

$\|y - \beta_1 x\|^2$ の最小化はyのxへの射影
 $\hat{y} = \hat{\beta}_1 x = \left(\frac{x}{\|x\|} \right) \left(\frac{x}{\|x\|} \right)' y = \left(\frac{x'y}{\|x\|^2} \right) x$

$S_{xx} = \sum x_i^2 = x'x = \|x\|^2, \quad S_{xy} = \sum x_i y_i = x'y$
 $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{x'y}{\|x\|^2}$

```
> n <- length(y)
> n
[1] 47
> i <- rep(1,n)
> i
[1] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
[39] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
> i %**% x / n
      [,1]
[1,] 9.540426
      [,1]
[1,] 1.472979
> x <- x - i * (i %**% x / n)
> y <- y - i * (i %**% y / n)
> mymean1(x)
[1] -1.530690e-15
> mymean1(y)
```

13

重回帰モデル (multiple regression model)

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$$

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

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

$$X = \begin{bmatrix} 1_n & x_1 & \dots & x_p \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix}$$

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

15

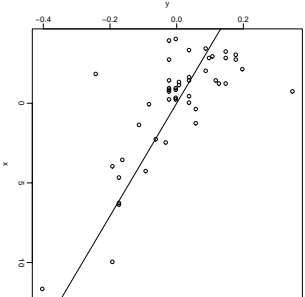
回帰分析のプログラミング

- 入力: X と y
- 回帰係数の推定: $\hat{\beta} = (X'X)^{-1}X'y$ または $\hat{\beta} = X^+y$
- 数値計算の方法: 逆行列, QR分解, 特異値分解などいろいろありえる
- 予測値: $\hat{y} = X\hat{\beta}$
- 残差: $e = y - \hat{y}$
- 出力: $\hat{\beta}, \hat{y}, e$

演習問題 (myfun20020919.Rのmy1sfitを参考にしてもよい。)

18

```
[1] -1.842498e-16
> (x %**% y)/(x %**% x) # = b1
      [,1]
[1,] -0.02824869
> plot(x,y)
> abline(0,(x %**% y)/(x %**% x))
```



射影

を最小にするには

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

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

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

ハット行列 (射影行列)

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

$1_n, x_1, \dots, x_p$ の張る空間への射影

$$\hat{y} = Hy$$

もし $X'X$ が退化している場合は

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

16

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$item[cbind(ax,ay)]
      E09504
"最終学歴が大学・大学院卒の者の割合" "未婚者割合 [20~24歳・女]"
      C01301
"県民1人当たり県民所得"
      B02101
"年平均気温"
      A05203
"合計特殊出生率"
      ,      = y
      ,      ,
> f <- my1sfit(x,y)
> f$summary
```

```
Estimate Std.Err t-value Pr(>|t|)
Intercept 3.636118e+00 6.108140e-01 5.9529052 4.642519e-07
E09504 -1.656659e-02 6.570354e-03 -2.5198932 1.562857e-02
```

19

回帰分析

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

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

もし $X'X$ が退化している場合は

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

17

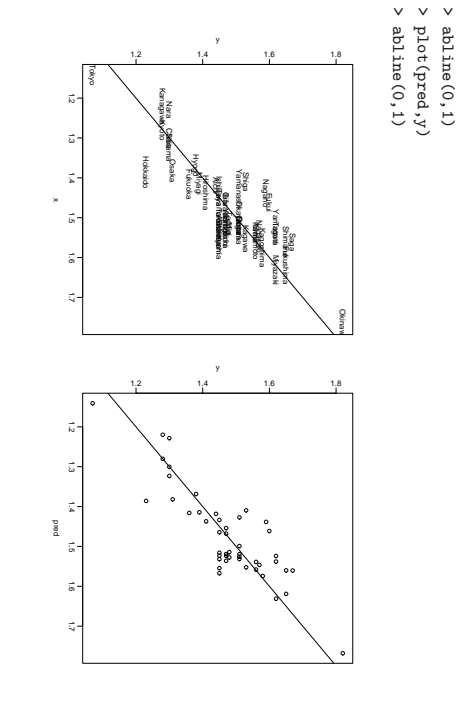
```
> f$summary
      A0410302  -2.727500e-02  7.438516e-03  -3.6667262  6.850351e-04
      C01301    1.228960e-05  4.635665e-05  0.2651098  7.922217e-01
      B02101    2.012511e-02  5.098161e-03  3.9475245  2.951848e-04
      ,      ,      ,      ,      ,      ,
> f$summary
      Mean Sum Sq R Squared 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]
[1,] 3.636118e+00
      E09504 -1.656659e-02
      A0410302 -2.727500e-02
      C01301 1.228960e-05
      B02101 2.012511e-02
      ,      ,      ,      ,
> pred <- x1 %**% coef
      ,      ,
> t(cbind(pred,y))
```

14

```

Hokkaido Amori Iwate Miyagi Akita Yamagata Fukushima Ibaraki
1.385854 1.518977 1.538630 1.414345 1.515823 1.524002 1.619096 1.467897
1.230000 1.470000 1.560000 1.390000 1.620000 1.650000 1.470000
Tochigi Gumma Saitama Chiba Tokyo Kanagawa Niigata Toyama
1.514316 1.531544 1.323061 1.300365 1.140643 1.219962 1.525388 1.464334
1.480000 1.510000 1.300000 1.300000 1.070000 1.280000 1.510000 1.450000
Ishikawa Fukui Yamaguchi Nagano Gifu Shizuoka Aichi Mie
1.433561 1.461327 1.427325 1.438374 1.453923 1.535835 1.418394 1.527675
1.450000 1.600000 1.510000 1.590000 1.470000 1.470000 1.440000 1.480000
Shiga Kyoto Osaka Hyogo Nara Wakayama Tottori Shimane
1.409378 1.280022 1.381912 1.368422 1.228403 1.554402 1.538035 1.559977
1.530000 1.280000 1.310000 1.380000 1.300000 1.450000 1.620000 1.650000
Okayama Hiroshima Yamaguchi Tokushima Kagawa Ehime Kochi Fukuoka
1.499213 1.436992 1.52368 1.523004 1.552604 1.531791 1.567354 1.415729
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 1.519622 1.631219 1.574310 1.768123
1.670000 1.570000 1.560000 1.510000 1.620000 1.580000 1.820000
> myplot(pred,y)

```



```

> ## ます単回帰, それからダミー変数 ( 2 )
> # 単回帰
> ax <- "C01301", ay <- "C04602"
> X2000$item[c(ax,ay)]
C01301
"個人 預貯金 残高 (人口 1 人当たり)"
C04602
" (千円:thousand yen) " (万円:10 thousand yen) "
C01301
> X2000$item[c(ax,ay)]
C04602
" (千円:thousand yen) " (万円:10 thousand yen) "
C01301
> f$summary
Mean Sum Sq R Squared F-value Df 1 Df 2 Pr(>|t|)
Intercept -69.1667194 92.49169483 -0.6396983 5.26857e-01
C01301 0.2087250 0.03192249 6.4445142 7.486358e-08
shikoku 163.9675326 42.22411062 3.8832679 3.422263e-04
C04602 79.0721 0.5238522 24.20414 2 44 8.1343e-08

```

```

Estimate Std.Err t-value Pr(>|t|)
Intercept 26.9369511 102.88756993 0.2618096 7.946630e-01
C01301 0.1804065 0.03580613 5.0384246 8.093436e-06
> f$summary
Mean Sum Sq R Squared F-value Df 1 Df 2 Pr(>|t|)
C04602 90.60172 0.3606658 25.38572 1 45 8.093436e-06
> myplot(x,y)
> abline(f)
> ## ダミー変数
> nhrefregion["shikoku"]
Hokkaido Amori Iwate Miyagi Akita Yamagata Fukushima Ibaraki
0 0 0 0 0 0 0 0
Tochigi Gumma Saitama Chiba Tokyo Kanagawa Niigata Toyama
0 0 0 0 0 0 0 0
Ishikawa Fukui Yamaguchi Nagano Gifu Shizuoka Aichi Mie
0 0 0 0 0 0 0 0
Shiga Kyoto Osaka Hyogo Nara Wakayama Tottori Shimane

```

```

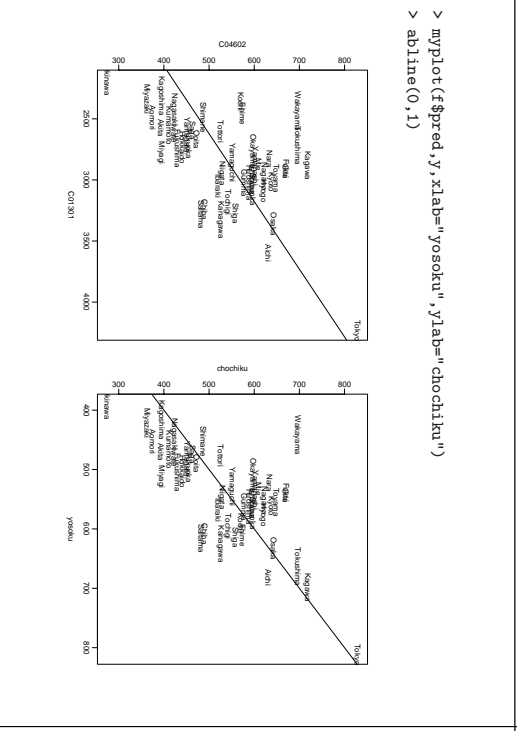
0 0 0 0 0 0 0 0
Okayama Hiroshima Yamaguchi Tokushima Kagawa Ehime Kochi Fuku
0 0 0 0 1 1 1 1
Saga Nagasaki Kumamoto Oita Miyazaki Kagoshima Okinawa
0 0 0 0 0 0 0 0
> x <- cbind(X2000$x[,ax,drop=F],nhhrefregion[, "shikoku",drop=F])
> t(x)
Hokkaido Amori Iwate Miyagi Akita Yamagata Fukushima Ibaraki Tochigi
C01301 2731 2489 2619 2776 2574 2629 2737 3047 3183
shikoku 0 0 0 0 0 0 0 0 0
Gumma Saitama Chiba Tokyo Kanagawa Niigata Toyama Ishikawa Fukui
C01301 3022 3280 3243 4230 3326 2941 2982 2996 2904
shikoku 0 0 0 0 0 0 0 0 0
Yamanashi Nagano Gifu Shizuoka Aichi Mie Shiga Kyoto Osaka Hyogo Har
C01301 2885 2969 2931 3073 3598 2874 3271 3015 3359 3088 287
shikoku 0 0 0 0 0 0 0 0 0 0 0
Wakayama Tottori Shimane Okayama Hiroshima Yamaguchi Tokushima Kagawa
C01301 2436 2604 2485 2764 3019 2855 2716 2885
shikoku 0 0 0 0 0 0 0 0 1

```

```

Estimate Std.Err t-value Pr(>|t|)
Intercept -69.1667194 92.49169483 -0.6396983 5.26857e-01
C01301 0.2087250 0.03192249 6.4445142 7.486358e-08
shikoku 163.9675326 42.22411062 3.8832679 3.422263e-04
C04602 79.0721 0.5238522 24.20414 2 44 8.1343e-08
> f$summary
Mean Sum Sq R Squared F-value Df 1 Df 2 Pr(>|t|)
C04602 90.60172 0.3606658 25.38572 1 45 8.093436e-06
> myplot(x,y)
> abline(f)
> ## 中心化したダミー変数を使う
> x <- x - mean(x)
> y <- y - mean(y)
> x <- cbind(x,nhrefregion)
C01301 tohoku kanto shintoshu tokai kinki chugoku shikoku kyushu
Hokkaido -118.659574 1 0 0 0 0 0 0
Amori -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 1 0 0 0 0 0
Fukuoshima -112.659574 1 0 0 0 0 0 0

```



```

> ## ます単回帰, それからダミー変数 ( 1 )
> # 単回帰
> ax <- "C01301"; ay <- "C04602"
> X2000$item[c(ax,ay)]
C01301
"個人 預貯金 残高 (人口 1 人当たり)"
C04602
" 県民 1 人当たりの県民所得 "
C01301
> x <- X2000$x[,ax,drop=F]
> y <- X2000$x[,ay,drop=F]
> f <- mylsfit(x,y)
> f$summary
Estimate Std.Err t-value Pr(>|t|)
Intercept 26.9369511 102.88756993 0.2618096 7.946630e-01
C01301 0.1804065 0.03580613 5.0384246 8.093436e-06

```

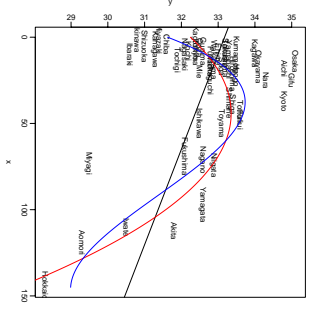
```

Estimate Std.Err t-value Pr(>|t|)
Intercept 26.9369511 102.88756993 0.2618096 7.946630e-01
C01301 0.1804065 0.03580613 5.0384246 8.093436e-06

```



```
> f3$fsummary
Mean Sum Sq R Squared F-value Df 1 Df 2 Pr(>F)
Y 1.006674 0.5437978 17.08548 3 43 1.877689e-07
> xs <- seq(min(x), max(x), length=300)
> ys <- cbind(1, xs, xs^2, xs^3) %>% f3$coefficients
> lines(xs, ys, col=4)
> dev.off()
null device
1
```



```
> ## 重相関係数
> t(round(cbind(f1$pred, y), 2))
Hokkaido Aomori Iwate Miyagi Akita Yamagata Fukushima Ibaraki Tochigi Gumma
Y 30.55 31.02 31.19 31.86 31.15 31.42 31.93 32.99 32.94 33.04
Y 28.30 29.30 30.50 29.50 31.80 32.60 32.10 30.60 31.90 32.66
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 31.88
Y 33.20 31.6 32.40 31.30 32.90 33.10 32.50 33.60 33.30 32.66
Gifu Shizuoka Aichi Mie Shiga Kyoto Osaka Hyogo Nara Wakayama Tottori
Y 32.72 33.08 32.85 32.81 32.47 32.51 32.92 32.78 32.72 32.94 32.4
Y 35.00 31.00 34.80 32.50 33.40 34.80 35.10 33.50 34.30 32.90 33.6
Shimane Okayama Hiroshima Yamaguchi Tokushima Kagawa Ehime Kochi Fukuoka
Y 32.49 32.87 32.8 32.78 32.96 33.01 32.99 33.03 32.92
Y 33.30 34.10 33.4 32.80 33.20 34.00 33.00 32.20 32.80
Saga Nagasaki Kumamoto Oita Miyazaki Kagoshima Okinawa
Y 32.87 32.99 32.99 33.07 33.16 33.08 33.18
Y 33.10 32.10 33.50 32.50 31.40 32.40 30.80
```

```
> cor(x, y)
[1] -0.4459788
> mycor1(f1$pred, y)
[1] 0.4459788
> cor(x, y)^2
[1] 0.1988971
> mycor1(f1$pred, y)^2
[1] 0.1988971
> f1$fsummary
Mean Sum Sq R Squared F-value Df 1 Df 2 Pr(>F)
Y 1.304014 0.1988971 11.17256 1 45 0.001678815
> ## 数値例4 (2次式)
> mycor1(f2$pred, y)
[1] 0.686967
> mycor1(f2$pred, y)^2
[1] 0.4719236
> f2$fsummary
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
```

あてはまりのよさ

重相関係数

y_i と \hat{y}_i の相関を重相関係数 R と呼ぶ。
重相関係数の二乗 R^2 は決定係数と呼ばれる。

$$R = \frac{S_{\hat{y}y}}{\sqrt{S_{yy}S_{\hat{y}\hat{y}}}}$$

$S_{\hat{y}y} = \sum (y_i - \bar{y})(\hat{y}_i - \bar{y}), S_{yy} = \sum (y_i - \bar{y})^2, S_{\hat{y}\hat{y}} = \sum (\hat{y}_i - \bar{y})^2$
すべての要素が \hat{y} のベクトル $\hat{y} = \hat{y}1_n$ を用いて
中心化 $\hat{y} \leftarrow \hat{y} - \bar{y}, \hat{y} \leftarrow \hat{y} - \bar{y}$ をおこなうと

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

あてはまりの良さを R または R^2 で判断する。

とくに単回帰の場合の R は r_{xy} と y_i の相関係数である。

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 $X'e = X'(I_n - H)y = 0$ 。とくに

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

ピタゴラスの定理

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

```
> ## 数値例4 (3次式)
> mycor1(f3$pred, y)
[1] 0.7374264
> mycor1(f3$pred, y)^2
[1] 0.5437978
> f3$fsummary
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
```

```
> 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] 1.675807e-15
> ## ピタゴラスの定理
> sum(y^2)
[1] 49980.06
> sum(f1$pred^2)
[1] 49903.54
> sum(e^2)
[1] 76.52033
> sum(f1$pred^2) + sum(e^2) - sum(y^2)
[1] -7.275958e-12
```

```
Y 0.23 -0.89 0.51 -0.57 -1.76 -0.68 -2.38
```

重相関係数と残差の関係

まず $y \leftarrow y - \bar{y}$, $\hat{y} \leftarrow \hat{y} - \bar{y}$ と中心化しておく,

$$R = \frac{\hat{y}'y}{\|\hat{y}\| \cdot \|y\|} = \cos \theta = \frac{\|\hat{y}\|}{\|y\|}$$

したがって

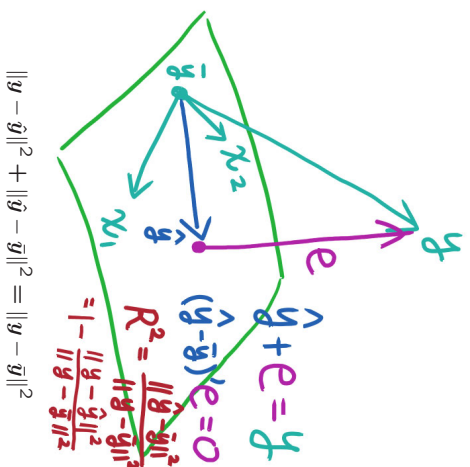
$$R^2 = (\cos \theta)^2 = \frac{\|\hat{y}\|^2}{\|y\|^2}, \quad 1 - R^2 = (\sin \theta)^2 = \frac{\|e\|^2}{\|y\|^2}$$

($e'y = 0$ なので $y'y = (\hat{y} + e)'y = \|\hat{y}\|^2$ をつかって示しても良い)

中心化していない場合は

$$R^2 = \frac{\|\hat{y} - \bar{y}\|^2}{\|y - \bar{y}\|^2} = 1 - \frac{\|y - \hat{y}\|^2}{\|y - \bar{y}\|^2}$$

25



26

```
> ### 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$pred^2)/sum(y^2) # 中心化しないとダメ
[1] 0.998469
> sum((f1$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) # 中心化しないとダメ
[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

QR分解

$$X = [x_0 \dots x_p], \quad Q = [q_0 \dots q_p], \quad Q'Q = I_{p+1}$$

$$X = QR, \quad r_{ij} = 0, i > j$$

- 回帰係数の変数変換: $\beta \leftrightarrow \gamma$

$$\gamma = R\beta$$

- 回帰モデルも変換される

$$X\beta = (QR)\beta = Q(R\beta) = Q\gamma$$

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

27

部分回帰 (Subset Regression)

モデル (k)

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \epsilon$$

すなわち

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

$$X^{(k)} = \begin{bmatrix} 1 & x_{11} & \dots & x_{k1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1n} & \dots & x_{kn} \end{bmatrix}, \quad \beta^{(k)} = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_k \end{bmatrix}$$

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

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

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

28

回帰係数の推定

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

$$\hat{\gamma} = (Q'Q)^{-1} Q'y = Q'y$$

$$\hat{\gamma} = R\hat{\beta}, \quad \hat{\beta} = R^{-1}\hat{\gamma}$$

つまりどちらで推定しても互いに変換できる.

- γ をつかう利点: $\gamma_k = r_{kk}\beta_k + \dots + r_{kp}\beta_p$ に注意すると

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

部分回帰モデル (k)

$$\hat{\gamma}^{(k)} = Q^{(k)' } y = \begin{bmatrix} q_0' \\ \vdots \\ q_k' \end{bmatrix} y, \quad \hat{\gamma}_i^{(k)} = \hat{\gamma}_i$$

29

- $\hat{\gamma}$ を $\hat{\beta}^{(0)}, \dots, \hat{\beta}^{(p)}$ に戻すのも容易.

$$\begin{bmatrix} \hat{\beta}_0^{(k)} \\ \vdots \\ \hat{\beta}_k^{(k)} \\ 0 \end{bmatrix} = R^{-1} \begin{bmatrix} \hat{\gamma}_0 \\ \vdots \\ \hat{\gamma}_k \\ 0 \end{bmatrix}$$

つまり $\hat{\gamma}$ と R^{-1} の計算は一度だけでよい.

- 部分モデルの計算を γ を通して行なうのは, 計算量を減らすという目的というだけでなく, むしろ, 回帰分析の性質を理論的に調べる際に重要.

```

> ## 直交化とQR分解
> xx <- cbind(1, 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 3046265 -0.145865 -0.46469866 0.58635783 0.54555490
Aomori 1 119 14161 1685159 -0.145865 -0.35684288 0.16893012 -0.16646877
Iwate 1 110 12100 1331000 -0.145865 -0.31950791 0.089565370 -0.26620527
Miyagi 1 73 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 -0.0654976 -0.28255814
Fukushima 1 69 4761 328509 -0.145865 -0.14942732 -0.21017241 -0.08803184
Ibaraki 1 10 100 1000 -0.145865 0.09594280 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 0.09389241 -0.07023124
Kanagawa 1 7 49 343 -0.145865 0.10776772 0.09389241 -0.07023124
Niigata 1 74 5476 405224 -0.145865 -0.17016886 -0.19734628 -0.13616848
Toyama 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 -0.08268060 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 0.14995720
Kyoto 1 17 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 0.04554312 -0.05826859 0.12799090
Nara 1 25 625 15625 -0.145865 0.03309820 -0.08268060 0.14536604
Wakayama 1 13 169 2197 -0.145865 0.08287788 0.02700782 0.03316958

```

```

Tottori 1 43 1849 79507 -0.145865 -0.04157133 -0.0.18701133 0.13242784
Shimane 1 38 1444 54872 -0.145865 -0.02082980 -0.16527710 0.15321596
Okayama 1 17 289 4913 -0.145865 0.06629465 -0.0.01312251 0.08367402
Hiroshima 1 21 441 9261 -0.145865 0.04669142 -0.0.04686532 0.12070886
Yamaguchi 1 22 484 10648 -0.145865 0.04554312 -0.0.05826859 0.12799090
Tokushima 1 12 144 1728 -0.145865 0.08702619 0.0.03759783 0.01831088
Kagawa 1 9 81 729 -0.145865 0.09947111 0.0.07070567 -0.03189642
Ehime 1 10 100 1000 -0.145865 0.09532280 0.0.05944675 -0.01420506
Kochi 1 8 64 512 -0.145865 0.10361942 0.0.08218756 -0.05056613
Fukuoka 1 14 196 2744 -0.145865 0.07872957 0.0.01664078 0.04711818
Saga 1 17 289 4913 -0.145865 0.06629465 -0.0.01312251 0.08367402
Nagasaki 1 10 100 1000 -0.145865 0.09532280 0.0.05944675 -0.01420506
Kumamoto 1 10 100 1000 -0.145865 0.09532280 0.0.05944675 -0.01420506
Oita 1 6 36 216 -0.145865 0.11191603 0.0.10582024 -0.09090882
Miyazaki 1 1 1 125 -0.145865 0.13265757 0.0.16880393 -0.21008090
Kagoshima 1 5 25 125 -0.145865 0.11606434 0.0.11797104 -0.11261593
Okinawa 1 0 0 0 -0.145865 0.13680587 0.0.18206958 -0.23719159
> R
1 x x~2 x~3

```

平方和の分解

$$\hat{y} = Q\hat{\gamma} = \tilde{\gamma}_0q_0 + \dots + \tilde{\gamma}_pq_p$$

$$\|\hat{y}\|^2 = \|\tilde{\gamma}_0q_0\|^2 + \dots + \|\tilde{\gamma}_pq_p\|^2 = \tilde{\gamma}_0^2 + \dots + \tilde{\gamma}_p^2$$

$$q_0 = \frac{1}{\sqrt{n}}\mathbf{1}_n, \quad \tilde{\gamma}_0 = \sqrt{n}\bar{y}$$

$$R^2 = \frac{\tilde{\gamma}_1^2 + \dots + \tilde{\gamma}_k^2}{\|\mathbf{y} - \tilde{y}\mathbf{1}_n\|^2}$$

モデル(k)

$$\|Q^{(k)}\hat{\gamma}^{(k)} - \tilde{y}\mathbf{1}_n\|^2 = \tilde{\gamma}_1^2 + \dots + \tilde{\gamma}_k^2$$

$$\|e\|^2 = \|\mathbf{y}\|^2 - \tilde{\gamma}_0^2 - \dots - \tilde{\gamma}_k^2 = \|\mathbf{y} - \tilde{y}\mathbf{1}_n\|^2 - \tilde{\gamma}_1^2 - \dots - \tilde{\gamma}_k^2$$

練習問題

2-1. 直線当てはめ(単回帰)の関数kaiik1を作れ.

```

kaiik1 <- function(x,y) {
# x,yは同じ長さのベクトル
# 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.
return(coef,pred,resid)
}

```

2-2. 重回帰分析の関数kaiik2を作れ.

```

kaiik2 <- function(x,y) {
# xはn * p次元の行列
# yは長さnのベクトル
}

```

```

> ## 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$pred,y)^2 # R^2(1次)
[1]
> cor(f2$pred,y)^2 # R^2(2次)
[1]
> cor(f3$pred,y)^2 # R^2(3次)
[1]
> sum(g3[2]^2)/sum((y-mean(y))^2)
[1] 0.1988971

```

```

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

```

```

> 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

```



```
# 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。
return(coef,pred,resid)
}
```

2-3. kaiki2 の返す値から決定係数 R^2 を計算する関数 ketteikeisu を作れ。

```
ketteikeisu <- function(kou) {
  # kou$pred は予測値, kou$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
  return(rsq)
}
```

2-4. kaiki2, ketteikeisu を使い、 χ^2_{2000} から適当な項目を選んで重回帰分析する。係数 β と決定係数 R^2 を計算する。myfunc20020919.R にある mylsfit 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。

2-5. 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
```