

回帰分析

予測と発展的なモデル

村田 昇

講義の内容

- 第1回: 回帰モデルの考え方と推定
- 第2回: モデルの評価
- 第3回: モデルによる予測と発展的なモデル

回帰分析の復習

線形回帰モデル

- 目的変数 を 説明変数 で説明する関係式を構成
 - 説明変数: x_1, \dots, x_p (p 次元)
 - 目的変数: y (1 次元)
- 回帰係数 $\beta_0, \beta_1, \dots, \beta_p$ を用いた一次式

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

- 誤差項 を含む確率モデルで観測データを表現

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

問題設定

- 確率モデル

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \text{確率分布}$$

- 式の評価: 残差平方和 の最小化による推定

$$S(\boldsymbol{\beta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

解とその一意性

- 解の条件: 正規方程式

$$\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^\top \mathbf{y}$$

- 解の一意性: Gram 行列 $\mathbf{X}^\top \mathbf{X}$ が正則

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

解析の事例

東京の8月の気候の分析

- データの一部

Table 1: 東京の8月の気候

日付	気温	降雨	日射	降雪	風向	風速	気圧	湿度	雲量
2022-08-01	30.6	0	24.53	0	SSE	2.8	1010.1	72	8.8
2022-08-02	31.6	0	24.78	0	SSE	2.5	1008.8	71	9.8
2022-08-03	31.5	0	21.24	0	SSE	2.3	1005.1	75	7.3
2022-08-04	24.6	18	3.46	0	NE	2.7	1006	89	10
2022-08-05	23.8	0	7.65	0	NE	2.9	1006.1	83	9.8
2022-08-06	25.2	0	17.06	0	SSE	2.4	1008.1	73	10
2022-08-07	27.6	0	14.45	0	SSE	2.2	1009.3	80	8.3
2022-08-08	29.8	0	22.52	0	S	4.5	1008.5	75	4.8
2022-08-09	30.9	0	25.5	0	S	5.5	1006.9	69	6.8
2022-08-10	30.5	0	25.99	0	S	5.3	1007.2	70	6
2022-08-11	29.5	0	22.9	0	S	5.4	1007.5	75	6
2022-08-12	28.3	2	15.36	0	S	5.8	1007.5	81	9.8
2022-08-13	25.5	47.5	4.53	0	S	4.8	1005.6	94	10
2022-08-14	28.2	0	16.28	0	SSE	2.6	1003	84	8.8
2022-08-15	29.4	0	18.65	0	S	2.5	1003.4	78	8.8
2022-08-16	31	0	20.5	0	SSW	4.8	1000.6	70	8.3
2022-08-17	27.3	5	8.87	0	NE	2.5	1005.8	77	10
2022-08-18	26.8	13	8.74	0	S	2.8	1001.7	81	6
2022-08-19	27.5	0	23.52	0	SSE	3.4	1001.7	62	3
2022-08-20	26.4	1.5	13.5	0	NW	1.8	1000.6	82	9.8
2022-08-21	26	1	8.96	0	NE	2.1	1002.3	87	10
2022-08-22	26.2	0	9.05	0	NNE	2.5	1005.5	82	10
2022-08-23	28.7	0	17.94	0	S	3.2	1003.2	83	8.3
2022-08-24	27.8	2	12.86	0	NE	2.9	1003.2	79	10
2022-08-25	25.7	0	9.83	0	SE	2	1004.1	77	10
2022-08-26	27	3.5	10.05	0	SSE	2.1	1002.5	89	10
2022-08-27	29	0	19.87	0	SSE	3.3	1002.7	80	5.5
2022-08-28	23.7	5	4.58	0	NE	3	1009.2	87	9.8
2022-08-29	23.3	0.5	15.45	0	NE	2.8	1016.1	69	8
2022-08-30	22.8	5	10.12	0	NNE	1.9	1012.5	88	10
2022-08-31	27.1	1	17.46	0	S	3.2	1007.6	85	8.8

- 気温を説明する5種類の線形回帰モデルを検討
 - モデル1: 気温 = F(気圧)
 - モデル2: 気温 = F(日射)
 - モデル3: 気温 = F(気圧, 日射)
 - モデル4: 気温 = F(気圧, 日射, 湿度)
 - モデル5: 気温 = F(気圧, 日射, 雲量)

分析の視覚化

- 関連するデータの散布図
- 観測値とあてはめ値の比較

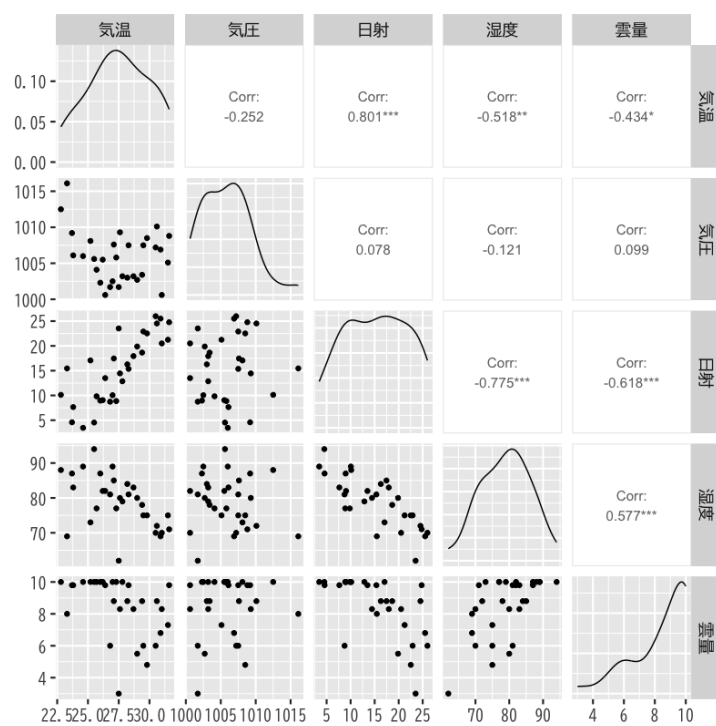


Figure 1: 散布図

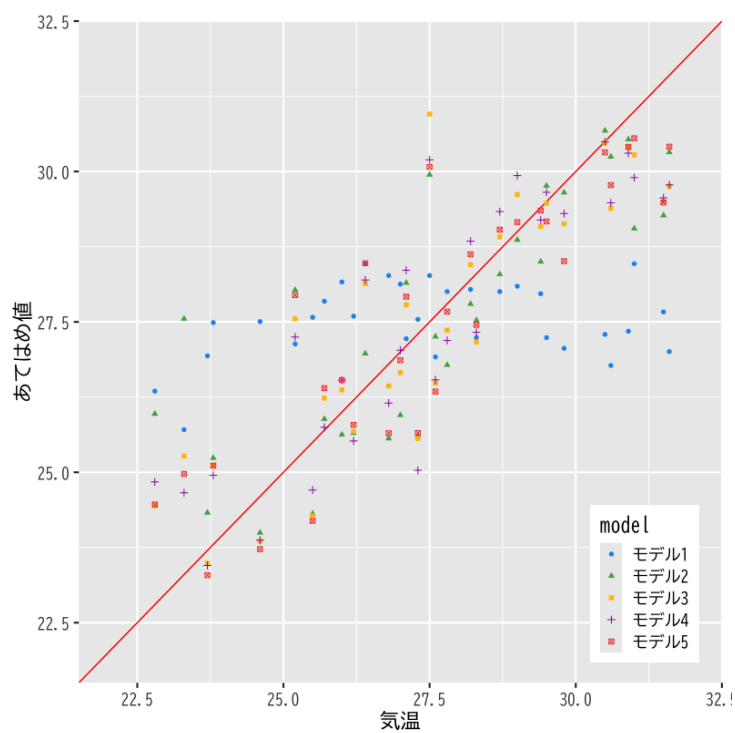


Figure 2: モデルの比較

寄与率

- 決定係数 (R-squared)

$$R^2 = 1 - \frac{\sum_{i=1}^n \hat{\epsilon}_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

- 自由度調整済み決定係数 (adjusted R-squared)

$$\bar{R}^2 = 1 - \frac{\frac{1}{n-p-1} \sum_{i=1}^n \hat{\epsilon}_i^2}{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2}$$

– 不偏分散で補正

モデルの評価

- 決定係数 (R^2 , Adjusted R^2)

Table 2: 寄与率によるモデルの比較

	目的変数				
	モデル 1	モデル 2	モデル 3	モデル 4	モデル 5
気温					
気圧	-0.178 (0.127)		-0.223*** (0.068)	-0.214*** (0.067)	-0.242*** (0.068)
日射		0.297*** (0.041)	0.306*** (0.036)	0.366*** (0.056)	0.348*** (0.045)
湿度				0.071 (0.051)	
雲量					0.238 (0.161)
Constant	206.535 (127.430)	22.969*** (0.690)	247.477*** (68.433)	231.843*** (68.254)	263.717*** (67.941)
R^2	0.064	0.641	0.741	0.758	0.760
Adjusted R^2	0.031	0.628	0.722	0.731	0.733

F 統計量による検定

- 説明変数のうち 1 つでも役に立つか否かを検定する
 - 帰無仮説 $H_0: \beta_1 = \dots = \beta_p = 0$
 - 対立仮説 $H_1: \exists j \beta_j \neq 0$ (少なくとも 1 つは役に立つ)
- F 統計量: 決定係数 (または残差) を用いて計算

$$F = \frac{n-p-1}{p} \frac{R^2}{1-R^2}$$

- p 値: 自由度 $p, n-p-1$ の F 分布で計算

モデルの評価

- F 統計量

Table 3: F 統計量によるモデルの比較

	目的変数				
	モデル 1	モデル 2	気温 モデル 3	モデル 4	モデル 5
気圧	-0.178 (0.127)		-0.223*** (0.068)	-0.214*** (0.067)	-0.242*** (0.068)
日射		0.297*** (0.041)	0.306*** (0.036)	0.366*** (0.056)	0.348*** (0.045)
湿度				0.071 (0.051)	
雲量					0.238 (0.161)
Constant	206.535 (127.430)	22.969*** (0.690)	247.477*** (68.433)	231.843*** (68.254)	263.717*** (67.941)
R ²	0.064	0.641	0.741	0.758	0.760
Adjusted R ²	0.031	0.628	0.722	0.731	0.733
Residual Std. Error	2.463 (df = 29)	1.526 (df = 29)	1.320 (df = 28)	1.298 (df = 27)	1.293 (df = 27)
F Statistic	1.973 (df = 1; 29)	51.743*** (df = 1; 29)	39.964*** (df = 2; 28)	28.174*** (df = 3; 27)	28.484*** (df = 3; 27)

Note:

*p<0.1; **p<0.05; ***p<0.01

t 統計量による検定

- 回帰係数 β_j が回帰式に寄与するか否かを検定する
 - 帰無仮説 $H_0: \beta_j = 0$
 - 対立仮説 $H_1: \beta_j \neq 0$ (β_j は役に立つ)
- t 統計量: 各係数ごと, ζ は $(X^T X)^{-1}$ の対角成分

$$t = \frac{\hat{\beta}_j}{\hat{\sigma} \zeta_j}$$

- p 値: 自由度 $n-p-1$ の t 分布を用いて計算

モデルの評価

- t 統計量

Table 4: t 統計量によるモデルの比較

	目的変数				
	モデル 1	モデル 2	気温 モデル 3	モデル 4	モデル 5
気圧	-0.178 (0.127) t = -1.405 p = 0.171		-0.223*** (0.068) t = -3.281 p = 0.003	-0.214*** (0.067) t = -3.185 p = 0.004	-0.242*** (0.068) t = -3.566 p = 0.002
日射		0.297*** (0.041) t = 7.193 p = 0.00000	0.306*** (0.036) t = 8.547 p = 0.000	0.366*** (0.056) t = 6.582 p = 0.00000	0.348*** (0.045) t = 7.699 p = 0.00000
湿度				0.071 (0.051) t = 1.390 p = 0.176	
雲量					0.238 (0.161) t = 1.474 p = 0.152
Constant	206.535 (127.430) t = 1.621 p = 0.116	22.969*** (0.690) t = 33.277 p = 0.000	247.477*** (68.433) t = 3.616 p = 0.002	231.843*** (68.254) t = 3.397 p = 0.003	263.717*** (67.941) t = 3.882 p = 0.001

診断プロットによる評価

- モデル 4

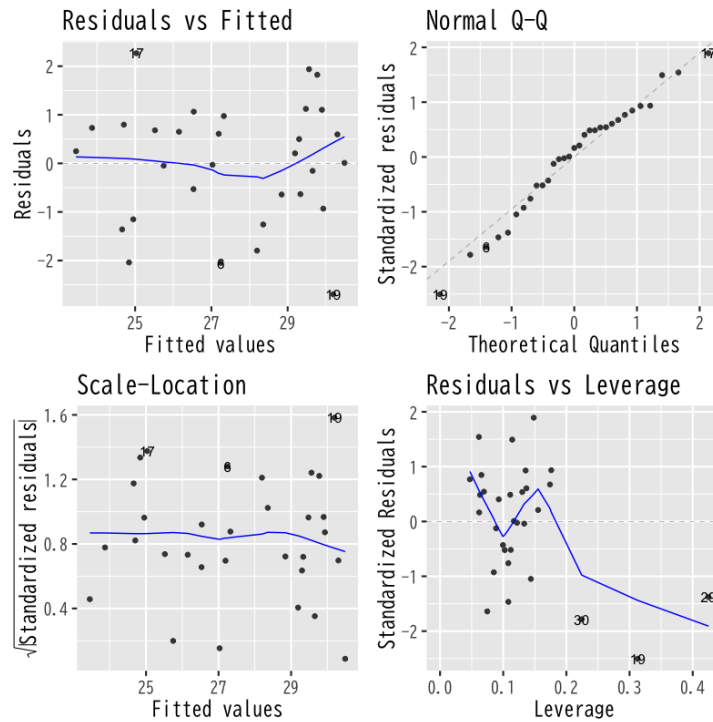


Figure 3: モデル 4 の診断

- モデル 5

回帰モデルによる予測

予測

- 新しいデータ (説明変数) x に対する **予測値**

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

- 予測値は元データの目的変数の重み付け線形和

$$\hat{y} = w(x)^T y, \quad w(x)^T = (1, x^T) (X^T X)^{-1} X^T$$

- 重みは元データと新規データの説明変数で決定

予測値の性質

- 推定量は以下の性質をもつ多変量正規分布

$$\begin{aligned} \mathbb{E}[\hat{\beta}] &= \beta \\ \text{Cov}(\hat{\beta}) &= \sigma^2 (X^T X)^{-1} \end{aligned}$$

- この性質を利用して以下の 3 つの値の違いを評価

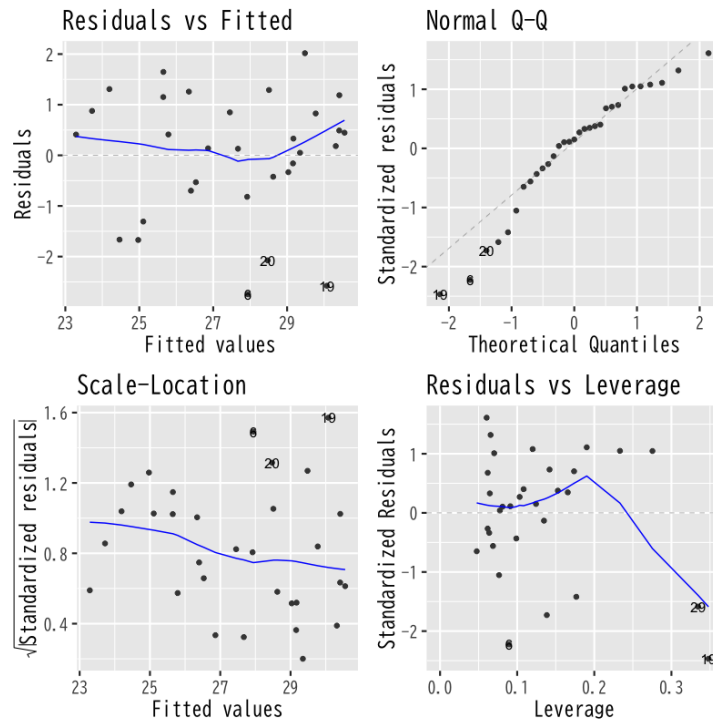


Figure 4: モデル 5 の診断

$$\begin{aligned}\hat{y} &= (1, \mathbf{x}^\top) \hat{\boldsymbol{\beta}} && \text{(回帰式による予測値)} \\ \tilde{y} &= (1, \mathbf{x}^\top) \boldsymbol{\beta} && \text{(最適な予測値)} \\ y &= (1, \mathbf{x}^\top) \boldsymbol{\beta} + \epsilon && \text{(観測値)}\end{aligned}$$

– \hat{y} と y は独立な正規分布に従うことに注意

演習

問題

- 誤差が平均 0 分散 σ^2 の正規分布に従うとき、以下の問に答えなさい
 - 予測値 \hat{y} の平均を求めよ
 - 予測値 \hat{y} の分散を求めよ

解答例

- 定義にもとづいて計算する

$$\begin{aligned}\mathbb{E}[\hat{y}] &= \mathbb{E}[(1, \mathbf{x}^\top) \hat{\boldsymbol{\beta}}] \\ &= (1, \mathbf{x}^\top) \mathbb{E}[\hat{\boldsymbol{\beta}}] \\ &= (1, \mathbf{x}^\top) \boldsymbol{\beta} \\ &= \tilde{y}\end{aligned}$$

– 真の回帰式による最適な予測値

- 定義にもとづいて計算する

$$\begin{aligned}
\text{Var}(\hat{y}) &= \text{Var}((1, \mathbf{x}^\top)(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})) \\
&= (1, \mathbf{x}^\top) \text{Cov}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})(1, \mathbf{x}^\top)^\top \\
&= (1, \mathbf{x}^\top) \text{Cov}(\hat{\boldsymbol{\beta}})(1, \mathbf{x}^\top)^\top \\
&= (1, \mathbf{x}^\top) \sigma^2 (X^\top X)^{-1} (1, \mathbf{x}^\top)^\top \\
&= \sigma^2 (1, \mathbf{x}^\top) (X^\top X)^{-1} (1, \mathbf{x}^\top)^\top
\end{aligned}$$

信頼区間

最適な予測値との差

- 差の分布は以下の平均・分散をもつ正規分布に従う

$$\begin{aligned}
\mathbb{E}[\tilde{y} - \hat{y}] &= (1, \mathbf{x}^\top) \boldsymbol{\beta} - (1, \mathbf{x}^\top) \mathbb{E}[\hat{\boldsymbol{\beta}}] = 0 \\
\text{Var}(\tilde{y} - \hat{y}) &= \underbrace{\sigma^2 (1, \mathbf{x}^\top) (X^\top X)^{-1} (1, \mathbf{x}^\top)^\top}_{\hat{\boldsymbol{\beta}} \text{ の推定誤差による分散}} = \sigma^2 \gamma_c(\mathbf{x})^2
\end{aligned}$$

- 正規化による表現

$$\frac{\tilde{y} - \hat{y}}{\sigma \gamma_c(\mathbf{x})} \sim \mathcal{N}(0, 1)$$

信頼区間

- 未知の分散を不偏分散で推定

$$Z = \frac{\tilde{y} - \hat{y}}{\hat{\sigma} \gamma_c(\mathbf{x})} \sim \mathcal{T}(n-p-1) \quad (t \text{ 分布})$$

- 確率 α の信頼区間

$$I_\alpha^c = (\hat{y} - C_\alpha \hat{\sigma} \gamma_c(\mathbf{x}), \hat{y} + C_\alpha \hat{\sigma} \gamma_c(\mathbf{x}))$$

$$P(|Z| < C_\alpha | Z \sim \mathcal{T}(n-p-1)) = \alpha$$

- 最適な予測値 \tilde{y} が入ることが期待される区間

演習

問題

- 以下の問に答えなさい
 - 信頼区間について以下の式が成り立つことを示せ

$$P(\tilde{y} \in I_\alpha^c) = \alpha$$

- 観測値と予測値の差 $y - \hat{y}$ の平均と分散を求めよ

解答例

- C_α の定義にもとづいて計算すればよい

$$\begin{aligned}\alpha &= P(|Z| < C_\alpha) \\ &= P\left(\left|\frac{\tilde{y} - \hat{y}}{\hat{\sigma}\gamma_c(\mathbf{x})}\right| < C_\alpha\right) \\ &= P(|\tilde{y} - \hat{y}| < C_\alpha \hat{\sigma}\gamma_c(\mathbf{x})) \\ &= P(-C_\alpha \hat{\sigma}\gamma_c(\mathbf{x}) < \tilde{y} - \hat{y} < C_\alpha \hat{\sigma}\gamma_c(\mathbf{x})) \\ &= P(\hat{y} - C_\alpha \hat{\sigma}\gamma_c(\mathbf{x}) < \tilde{y} < \hat{y} + C_\alpha \hat{\sigma}\gamma_c(\mathbf{x}))\end{aligned}$$

- 独立性を利用して計算する

$$\begin{aligned}\mathbb{E}[y - \hat{y}] &= \mathbb{E}[y] - \mathbb{E}[\hat{y}] \\ &= \tilde{y} - \tilde{y} \\ &= 0 \\ \text{Var}(y - \hat{y}) &= \text{Var}(y) + \text{Var}(\hat{y}) \\ &= \sigma^2 + \sigma^2(1, \mathbf{x}^\top)(X^\top X)^{-1}(1, \mathbf{x}^\top)^\top\end{aligned}$$

予測区間

観測値との差

- 差の分布は以下の平均・分散をもつ正規分布に従う

$$\begin{aligned}\mathbb{E}[y - \hat{y}] &= (1, \mathbf{x}^\top)\boldsymbol{\beta} + \mathbb{E}[\epsilon] - (1, \mathbf{x}^\top)\mathbb{E}[\hat{\boldsymbol{\beta}}] = 0 \\ \text{Var}(y - \hat{y}) &= \underbrace{\sigma^2(1, \mathbf{x}^\top)(X^\top X)^{-1}(1, \mathbf{x}^\top)^\top}_{\hat{\boldsymbol{\beta}} \text{ の推定誤差による分散}} + \underbrace{\sigma^2}_{\text{誤差の分散}} = \sigma^2\gamma_p(\mathbf{x})^2\end{aligned}$$

- 正規化による表現

$$\frac{y - \hat{y}}{\sigma\gamma_p(\mathbf{x})} \sim \mathcal{N}(0, 1)$$

予測区間

- 未知の分散を不偏分散で推定

$$Z = \frac{y - \hat{y}}{\hat{\sigma}\gamma_p(\mathbf{x})} \sim \mathcal{T}(n-p-1) \quad (t \text{ 分布})$$

- 確率 α の予測区間

$$I_\alpha^P = (\hat{y} - C_\alpha \hat{\sigma}\gamma_p(\mathbf{x}), \hat{y} + C_\alpha \hat{\sigma}\gamma_p(\mathbf{x}))$$

$$P(|Z| < C_\alpha | Z \sim \mathcal{T}(n-p-1)) = \alpha$$

- 観測値 y が入ることが期待される区間
- $\gamma_p > \gamma_c$ なので信頼区間より広くなる

解析の事例

信頼区間と予測区間

- 東京の気候データを用いて以下を試みる
 - 8月のデータで回帰式を推定する
気温 = $F(\text{気圧}, \text{日射}, \text{湿度})$
 - 上記のモデルで9月のデータを予測する

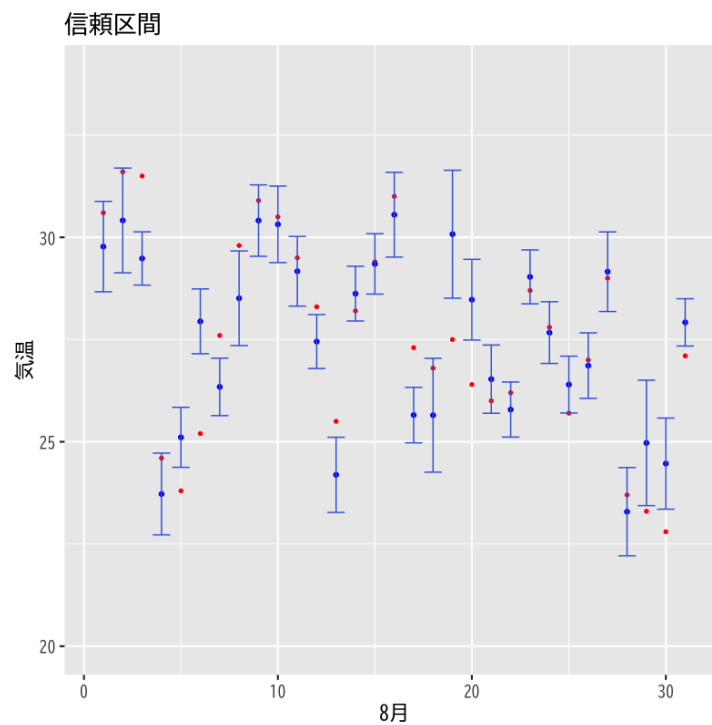


Figure 5: 8月のあてはめ値の信頼区間

発展的なモデル

非線形性を含むモデル

- 目的変数 y
- 説明変数 x_1, \dots, x_p
- 説明変数の追加で対応可能
 - 交互作用 (交差項): $x_i x_j$ のような説明変数の積
 - 非線形変換: $\log(x_k)$ のような関数による変換

カテゴリカル変数を含むモデル

- 数値ではないデータ
 - 悪性良性
 - 血液型
- 適切な方法で数値に変換して対応:

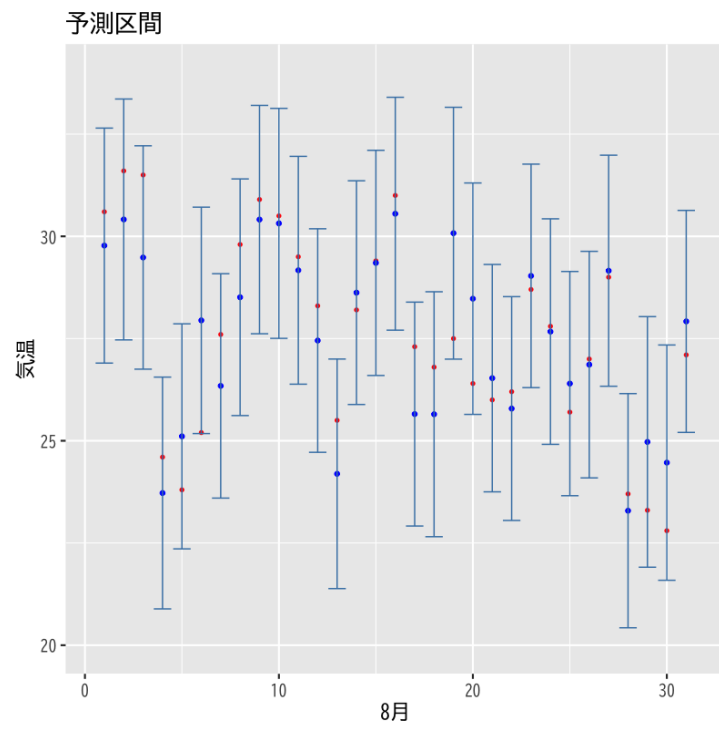


Figure 6: 8月のあてはめ値の予測区間

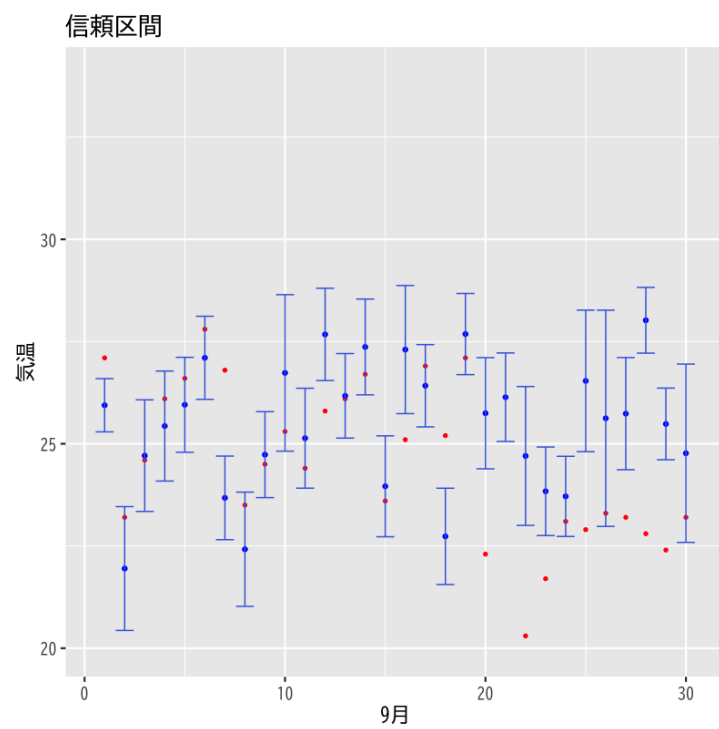


Figure 7: 8月モデルによる9月の予測値の信頼区間

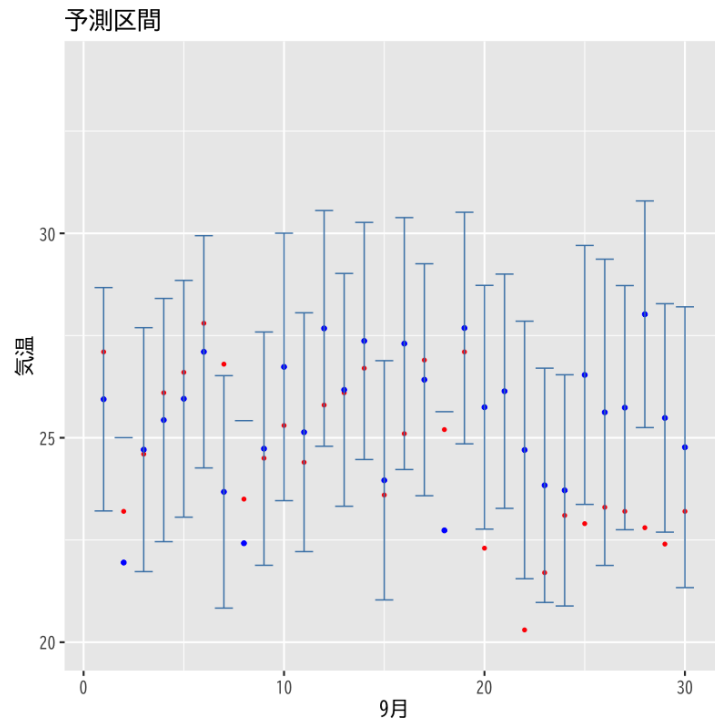


Figure 8: 8月モデルによる9月の予測値の予測区間

- 2 値の場合は 1,0 (真, 偽) を割り当てる
 - * 悪性 : 1
 - * 良性 : 0
- 3 値以上の場合は **ダミー変数** を利用する (カテゴリ数-1 個)
 - * A 型 : (1,0,0)
 - * B 型 : (0,1,0)
 - * O 型 : (0,0,1)
 - * AB 型 : (0,0,0)

解析の事例

非線形変換による線形化

- 様々な動物の体重と脳の重さの関係を調べる
 - 体重は 5 桁程度のばらつき
 - 脳の重さは 4 桁程度のばらつき
- 以下の変換を検討する
 - 変換なし
 - 体重を対数変換
 - 体重および脳の重さを対数変換
- 散布図 (変換なし)
- 散布図 (x 軸を対数変換)

Table 5: 体重と脳の重さ

	body	brain
Mountain beaver	1.350	8.100
Cow	465	423
Grey wolf	36.330	119.500
Goat	27.660	115
Guinea pig	1.040	5.500
Dipliodocus	11,700	50
Asian elephant	2,547	4,603
Donkey	187.100	419
Horse	521	655
Putar monkey	10	115
Cat	3.300	25.600
Giraffe	529	680
Gorilla	207	406
Human	62	1,320
African elephant	6,654	5,712
Triceratops	9,400	70
Rhesus monkey	6.800	179
Kangaroo	35	56
Golden hamster	0.120	1
Mouse	0.023	0.400
Rabbit	2.500	12.100
Sheep	55.500	175
Jaguar	100	157
Chimpanzee	52.160	440
Rat	0.280	1.900
Brachiosaurus	87,000	154.500
Mole	0.122	3
Pig	192	180

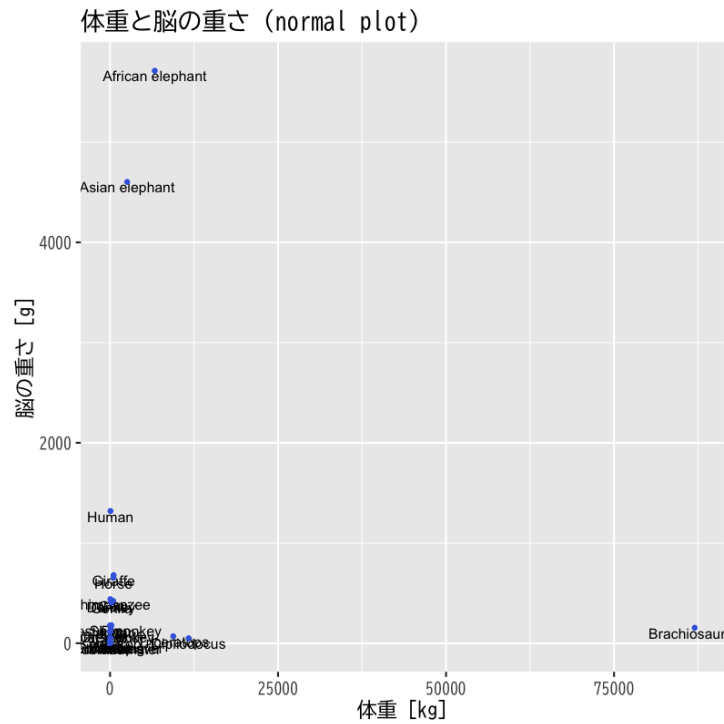


Figure 9: 散布図 (データの変換なし)

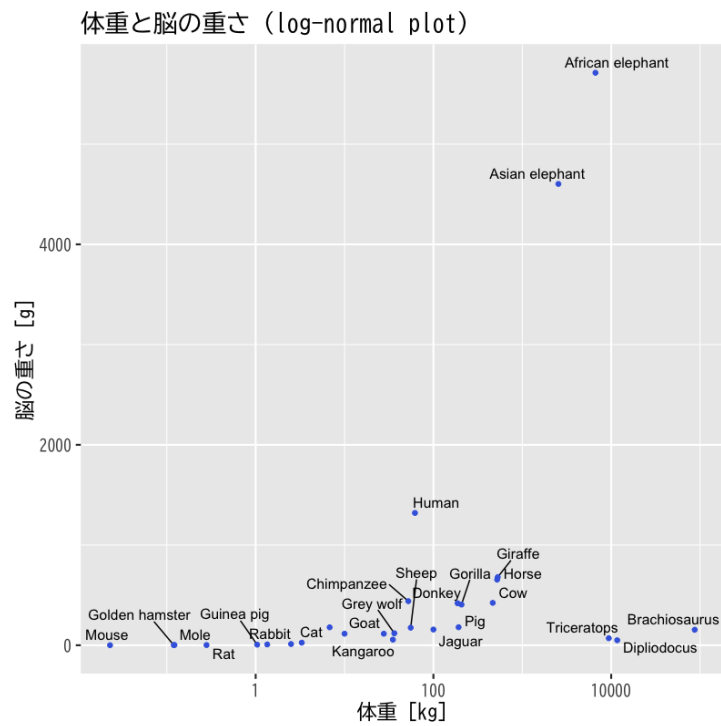


Figure 10: 散布図 (体重を対数変換)

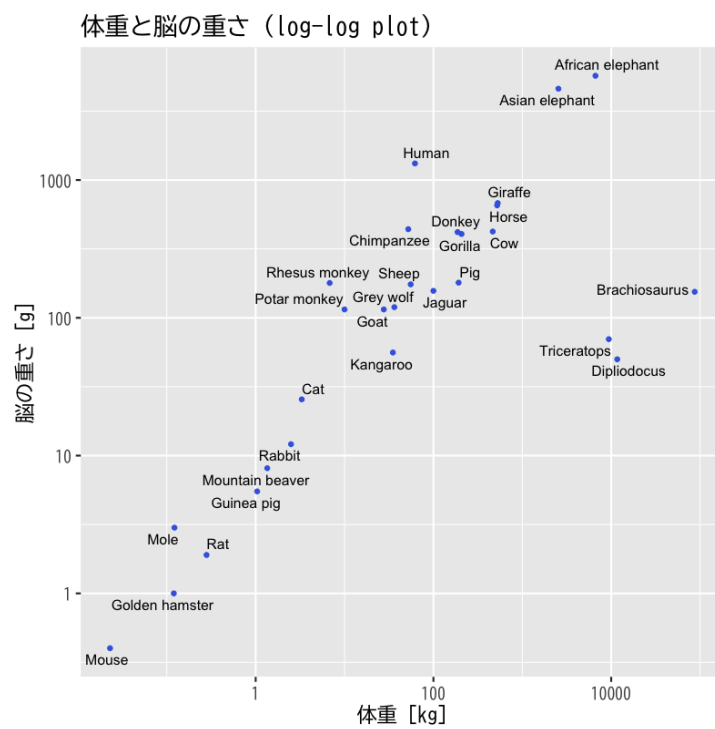


Figure 11: 散布図 (体重と脳の重さを対数変換)

- 散布図 (xy 軸を対数変換)
- 単回帰 (全データ)

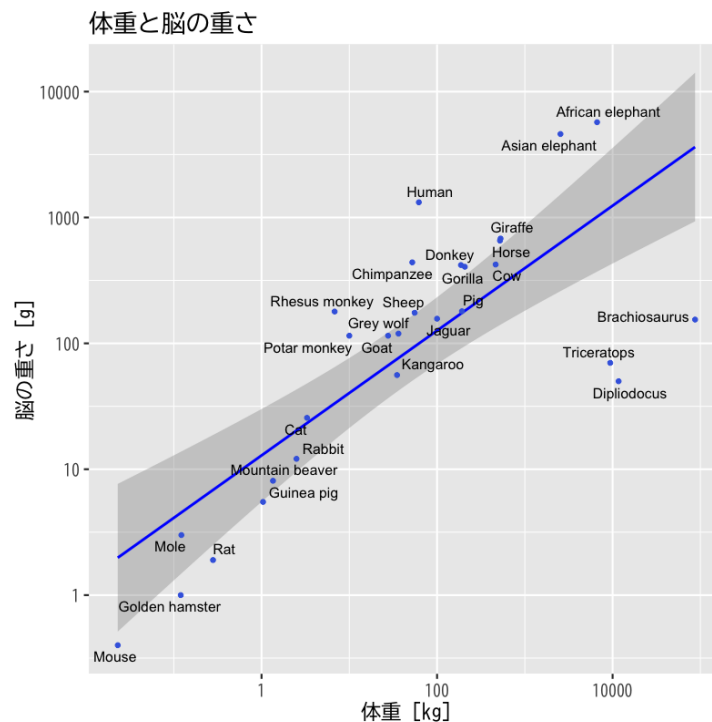


Figure 12: 単回帰

- 単回帰 (外れ値を除去)

非線形な関係の分析

- 東京の気候データを用いて気温に影響する変数の関係を検討する
 - 日射量と気圧の線形回帰モデル
(日射量と気圧が気温にどのように影響するか検討する)
 - これらの交互作用を加えた線形回帰モデル
(日射量と気圧の相互の関係の影響を検討する)
- 関連データの散布図

交互作用の効果

- 気温への寄与
 - 線形モデル
 - * 日射が高くなるほど高
 - * 気圧が低くなるほど高
 - 交互作用を加えたモデル
 - * ある気圧より高い場合には日射量が高くなるほど高
 - * ある日射量より高い場合には気圧が高くなるほど高
 - * 係数の有意性は低いのでより多くのデータでの分析が必要

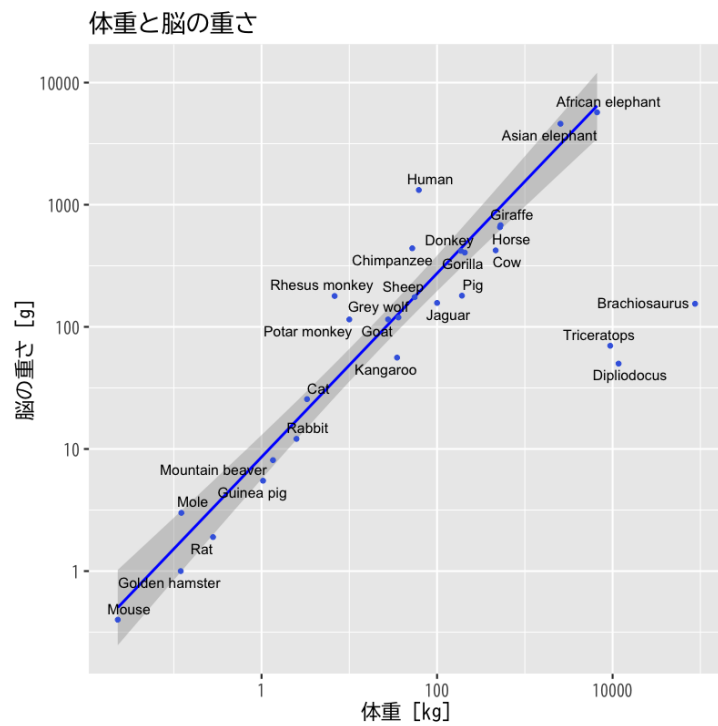


Figure 13: 外れ値を除いた単回帰

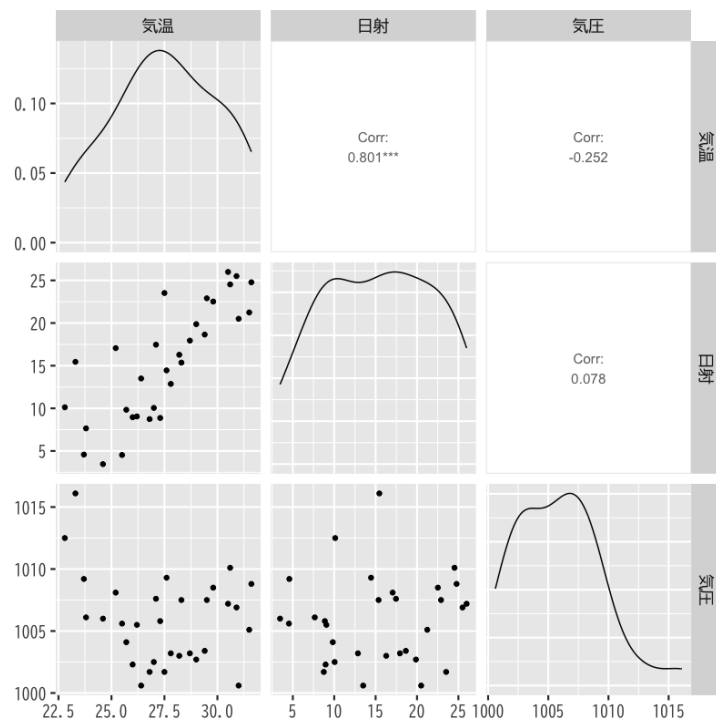


Figure 14: 散布図 (気温・日射・気圧)

Table 6

	目的変数	
	気温	
	交互作用なし	交互作用あり
日射	0.306*** (0.036)	-25.728* (13.058)
気圧	-0.223*** (0.068)	-0.622*** (0.210)
日射 × 気圧		0.026* (0.013)
Constant	247.477*** (68.433)	648.402*** (211.362)
Observations	31	31
R ²	0.741	0.774
Adjusted R ²	0.722	0.749
Residual Std. Error	1.320 (df = 28)	1.255 (df = 27)
F Statistic	39.964*** (df = 2; 28)	30.798*** (df = 3; 27)

Note: *p<0.1; **p<0.05; ***p<0.01

カテゴリカル変数の利用

- 東京の気候データを用いて気温を回帰するモデルを検討する
 - 降水の有無を表すカテゴリカル変数を用いたモデル
(雨が降ると気温が変化することを検証する)
 - 月をカテゴリカル変数として加えたモデル
(月毎の気温の差を考慮する)
- 関連データの散布図

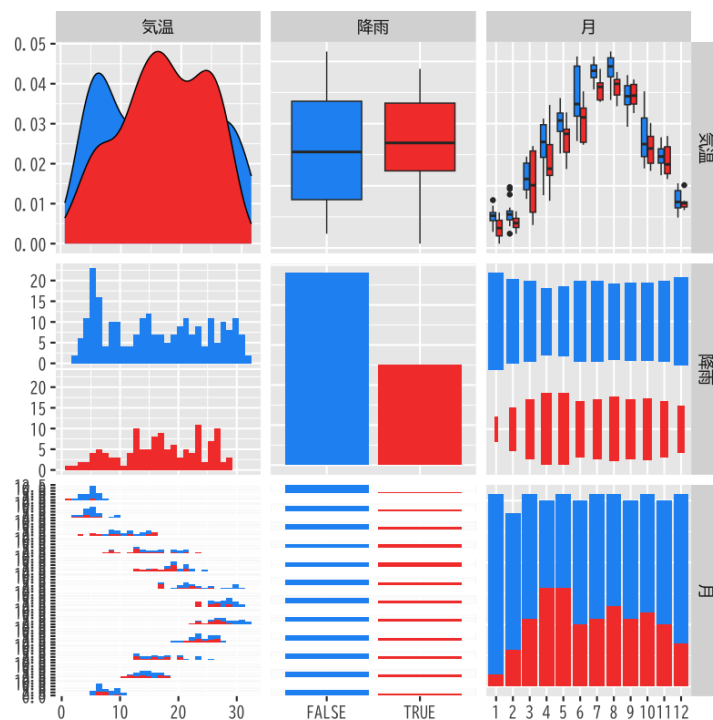


Figure 15: 散布図 (気温・降雨・月)

カテゴリカル変数の効果

- 気温への寄与

Table 7

	目的変数	
	降水	気温 降水+月
降水の有無	0.828 (0.913)	-1.876*** (0.308)
2 月		0.616 (0.705)
3 月		6.580*** (0.691)
4 月		11.243*** (0.706)
5 月		14.735*** (0.700)
6 月		18.598*** (0.696)
7 月		23.035*** (0.691)
8 月		23.295*** (0.694)
9 月		20.040*** (0.698)
10 月		12.915*** (0.693)
11 月		10.098*** (0.696)
12 月		2.945*** (0.687)
Constant	16.161*** (0.534)	5.018*** (0.485)
Observations	365	365
R ²	0.002	0.897
Adjusted R ²	-0.0005	0.894
Residual Std. Error	8.277 (df = 363)	2.700 (df = 352)
F Statistic	0.823 (df = 1; 363)	255.660*** (df = 12; 352)

Note:

*p<0.1; **p<0.05; ***p<0.01

– 降水モデル

- * 降水の有無は気温の予測に無関係ではないと考えられる
- * 決定係数から回帰式としての説明力は極めて低い
- * 通年では雨と気温の関係は積極的に支持されない

– 降水+月モデル

- * 月毎の気温の偏りが月の係数として推定される
- * 雨の日の方が気温が低いことが支持される

次回の予定

- 第 1 回: 主成分分析の考え方
- 第 2 回: 分析の評価と視覚化