

第9回 続高橋セミナー  
最尤法によるポアソン回帰分析入門  
2020年7月11日

第11章 デビアンズ・逸脱度・テコ比・4種の残差

一般化線形モデルによるポアソン回帰分析では、通常の回帰分析とは結果の表記法がかなり異なり、なかなか馴染めないのではないだろうか。デビアンズ・逸脱度は、その代表的な例であろう。逸脱度を理解するためには対数尤度と最尤法の理解も必要となり最小2乗法に慣れ親しんだ人たちは、茫然自失の状態になるかもしれない。さらに、スチューデント化デビアンズ残差に関連して「テコ比」も登場する。「テコ比」は、通常の回帰分析にも登場することもあるが、マイナーな存在である。そこで、通常の回帰分析と対比しつつポアソン回帰で使われている統計用語とその意味づけについてこれまで取り上げてきた事例を用いて関連付けを行う。

第11章 目次

11. デビアンズ・逸脱度・テコ比・4種の残差	359
11.1. デビアンズ	359
11.2 通常の回帰分析におけるスチューデント化残差	361
回帰パラメータの推定, 分散分析表, パラメータの共分散行列, スチューデント化残差, テコ比・ハット行列, テコ比の活用, Excelの「標準残差」に対する使用上の注意	
11.3. ポアソン回帰におけるデビアンズ・逸脱度	368
デビアンズ・カイ2乗, Pearson・カイ2乗, AICc	
11.4. ポアソン回帰における4種の残差	372
デビアンズ残差, スチューデント化デビアンズ残差, スチューデント化Pearson残差, SAS/GENMODによる各種の残差	
11.5. カブトガニの事例における4種の残差	379
JMPによる4種の残差の計算, 4種の残差の比較	
文献索引, 索引, 解析用ファイル一覧	383

## 第9回 続高橋セミナー 最尤法によるポアソン回帰分析入門

第9回 続高橋セミナー「最尤法によるポアソン回帰分析入門」は、ページ数が多いので章ごとに公開する。全体の章立てを次に示す。

### 目 次

はじめに -----	1
1. ポアソン分布に従う各種のカウント・データ-----	7
2. ニュートン・ラフソン法によるポアソン回帰 -----	63
3. 尤度比検定のためのデザイン行列-----	95
4. デザイン行列を用いた回帰分析入門-----	135
5. 反復重み付き最尤法によるポアソン回帰 -----	175
6. 過分散・ゼロ過剰への対応 -----	207
7. 過分散がある場合の探索的ポアソン回帰 -----	237
8. 2本の回帰直線の比較-----	269
9. 花数を共変量とした種子数の分析 -----	293
10. オフセットを含む探索的ポアソン回帰-----	323
<b>11. デビアン스・逸脱度・テコ比・4種の残差 -----</b>	<b>359</b>
12. パラメータの共分散行列の活用 -----	383
13. 最小2乗平均の謎を予測プロファイルで解く -----	421
文献, 文献索引, 索引, (解析用ファイル) 一覧 -----	461

## 11. デビアンズ・逸脱度・テコ比・4種の残差

一般化線形モデルによるポアソン回帰分析では、通常の回帰分析とは結果の表記法がかなり異なり、なかなか馴染めないのではないだろうか。デビアンズ・逸脱度は、その代表的な例であろう。逸脱度を理解するためには対数尤度と最尤法の理解も必要となり最小 2 乗法に慣れ親しんだきた人たちは、茫然自失の状態になるかもしれない。さらに、スチューデント化デビアンズ残差に関連して「テコ比」も登場する。「テコ比」は、通常の回帰分析にも登場することもあるが、マイナーな存在である。そこで、通常の回帰分析と対比しつつポアソン回帰で使われている統計用語とその意味づけについてこれまで取り上げてきた事例を用いて関連付けを行う。

### 11.1. デビアンズ

第 1.9 節では、久保 (2012) で取り上げられている種子数のデータについて表 1.34 に Excel で計算した 3 種の対数尤度 (縮小モデル, 完全モデル, 飽和モデル) を示し、デビアンズを飽和モデルの対数尤度  $\ln L_{\text{飽和}}$  から完全モデルの対数尤度  $\ln L_{\text{完全}}$  の差の 2 倍, 84.9930 として示した。このデビアンズ 84.9930 が、飽和モデルの自由度 100 に対し、完全モデルの自由度 2 との差の自由度 98 のカイ 2 乗分布に従うことから、上側確率が 0.8226 となり、100 個のパラメータで推定した飽和モデルに対し、2 個のパラメータで推定した完全モデルが統計的に遜色ないことを示した。

単にデビアンズと言う場合は、データに対して何らかの仮定をした 2 つのモデルから計算される対数尤度の差の 2 倍した統計量を意味する。ややこしいのは、統計ソフトを使うと「デビアンズ残差」あるいは「スチューデント化デビアンズ残差」なども出てくる。第 1.9 節で JMP によるポアソン回帰で出力される「スチューデント化デビアンズ残差プロット」を図 1.9 に示したが、その計算方法を示さずに「ほとんどが (-2~+2) の範囲に入っていることからあてはめの妥当性が示されている」との解説をしている。なぜ、そのような判断ができるのか、本章で Excel による計算を通して理解を深めたい。

第 9.5 節では、R による負の 2 項回帰 (ガンマ・ポアソン回帰) による結果を文献から引用した際に「Deviance Residuals」と「Residual deviance」について Excel での計算例を示し、「デビアンズ・逸脱度」について断片的な解説をした。しかし、他では意図的に使ってこなかった。それに代えてマイナス 2 倍の対数尤度、対数尤度の差の 2 倍、尤度比カイ 2 乗値を用い

てきた。第 1.9 節では、Pearson のカイ 2 乗値の計算方法を示したが、デビエンスについては深入りせず、Pearson のカイ 2 乗値を多用してきた。これは、デビアンズ残差の計算で用いている「飽和モデルの対数尤度」の概念が通常の最小 2 乗法の世界で、対応する概念がないために、あえて言及しなかった。

本章では、通常の回帰分析で標準的に使われている「分散分析表」の概念と比較しながらデビアンズについて解説する。さらに、通常の回帰分析で行われている通常の残差に加え、スチューデント化残差（標準化残差）と対比して（Pearson 残差，スチューデント化 Pearson 残差，デビアンズ残差，スチューデント化デビアンズ残差）について、統計ソフトの結果と対比しつつ Excel を用いた計算法を提示し、それらの使い分けについて概説する。

Excel での各種の残差を計算する際に、厄介な問題に直面する。それは、テコ比  $h_{ii}$  がハット行列  $H$

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

の対角要素と定義されているため、Excel の行列関数で対角要素を取り出すことが容易ではない。そのために、デザイン行列  $X$  の  $i$  行目のベクトル  $\mathbf{x}_i$  を使った

$$h_{ii} = \mathbf{x}_i (X^T X)^{-1} \mathbf{x}_i^T$$

計算により代替する。この式を使って  $h_{1,1}$ ,  $h_{2,2}, \dots$  を列ベクトルとして得ることができる。このような、Excel による逐次的な計算手順を経験することにより、テコ比  $h_{ii}$  の意味付けと活用方法について理解を深めてもらいたい。

テコ比  $h_{ii}$  を使うことにより、通常の回帰分析でのスチューデント化残差（標準化残差）を計算することができるようになり、ポアソン回帰では、スチューデント化デビアンズ残差の計算が可能となる。

## 11.2. 通常の回帰分析におけるスチューデント化残差

### 回帰パラメータの推定

第 4 章で示した Excel の行列関数を用いた回帰分析について要点をまとめる．用いるデータは，第 1.4 節のポアソン回帰のための人工データである [ドブソン (2008)]．表 11.1 に示すように，まず，説明変数  $X$  の平均からの偏差を計算し，それらの平方から平方和  $S_{XX}$  を計算する．次に，反応変数  $Y$  の偏差と  $X$  の偏差の積和  $S_{XY}$  を計算し，それらから，回帰パラメータ傾き  $\hat{\beta}_1$  を

$$\hat{\beta}_1 = \frac{S_{XY}}{S_{XX}} = \frac{24.000}{4.8889} = 4.9091$$

により計算する．説明変数  $X$  の平均  $\bar{X}$  と反応変数  $Y$  の平均  $\bar{Y}$  と傾き  $\hat{\beta}_1$  を用いて回帰パラメータの切片  $\hat{\beta}_0$  を，

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} = 8.0 - 4.9091 \times 0.1111 = 7.4545$$

として計算する．なお，表中に  $S_T = S_{YY}$  があるが，回帰パラメータの推定には使われない．

表 11.1 回帰パラメータの Excel シート上での推定 (表 4.1 再掲)

$i$	$X$	$Y$	$X$ 偏差	$X$ 偏差 <sup>2</sup>	$Y$ 偏差	$Y$ 偏差 <sup>2</sup>	$XY$ 偏差		
1	-1	2	-1.1111	1.2346	-6.0000	36.0000	6.6667	$\hat{\beta}_1 =$	4.9091
2	-1	3	-1.1111	1.2346	-5.0000	25.0000	5.5556	$\hat{\beta}_0 =$	7.4545
3	0	6	-0.1111	0.0123	-2.0000	4.0000	0.2222		
4	0	7	-0.1111	0.0123	-1.0000	1.0000	0.1111		
5	0	8	-0.1111	0.0123	0.0000	0.0000	0.0000		
6	0	9	-0.1111	0.0123	1.0000	1.0000	-0.1111		
7	1	10	0.8889	0.7901	2.0000	4.0000	1.7778		
8	1	12	0.8889	0.7901	4.0000	16.0000	3.5556		
9	1	15	0.8889	0.7901	7.0000	49.0000	6.2222		
	0.1111	8.0000	0.0000	4.8889	0.0000	136.0000	24.0000		
	平均	平均	合計	平方和	合計	平方和	平方和		
	$\bar{X}$	$\bar{Y}$		$S_{XX}$		$S_T = S_{YY}$	$S_{XY}$		

Excel での一般的な計算方法は，セルに対し計算式を与え，そのセルをプルダウンして参照セルを変化させつつ計算式をコピーする．これは，Excel の画期的な計算機能であるが，操作ミスがあっても発見しづらい．セルごとの計算に代え，範囲を使った計算を使うのが確実である．たとえば，8 行 1 列の「 $X$  偏差」は， $=(X \text{ の範囲} - X \text{ の平均})$  のように行列計算の要領で一括計算ができる．「 $X$  偏差<sup>2</sup>」は， $=(X \text{ 偏差の範囲})^2$  によって 8 行 1 列分の計算を一括して行っている．「 $XY$  偏差」は， $=(X \text{ 偏差} * Y \text{ 偏差})$  として一括計算している．範囲を指定した場合の四則演算は，対応するセル同士の計算となり，片方がスカラーのような場合は，相手のサイズに合わせてくれる．

## 分散分析表

表 11.2 に示すように分散分析表に必要な平方和の計算を逐次的に行った結果を示す. なお, 各種の平方和の計算は, `SumSq()` 関数を使用するのが効率的である.

$$\text{総平方和} \quad S_T = \sum_{i=1}^9 (Y_i - \bar{Y})^2 = \text{SumSq}(\mathbf{Y} \text{の範囲} - \bar{Y}) = 136.0000 \quad df = 9 - 1 = 8$$

$$\text{推定値} \quad \hat{\mathbf{Y}} = \mathbf{X} \hat{\boldsymbol{\beta}} = \text{Mmult}(\mathbf{X} \text{の範囲}, \hat{\boldsymbol{\beta}} \text{の範囲}) \quad df = 2$$

$$\text{回帰の平方和} \quad S_R = \sum_{i=1}^9 (\hat{Y}_i - \bar{Y})^2 = \text{SumSq}(\hat{\mathbf{Y}} \text{の範囲} - \bar{Y}) = 117.8182 \quad df = 2 - 1 = 1$$

$$\text{誤差平方和} \quad S_e = \sum_{i=1}^9 (Y_i - \hat{Y}_i)^2 = \text{SumSq}(\mathbf{Y} \text{の範囲} - \hat{\mathbf{Y}} \text{の範囲}) = 18.1818 \quad df = 9 - 2 = 7$$

回帰の平方和  $S_R$  は,

$$S_R = S_T - S_e = 136.0000 - 18.1818 = 117.8182$$

としても計算できる.

表 11.2 回帰の平方和の計算 (表 4.3 再掲)

$i$	$\mathbf{X}$		$Y_i$	$Y^-$	$Y_i - Y^-$	$Y_i^\wedge$	$Y_i^\wedge - Y^-$	$Y_i - Y_i^\wedge$		$\boldsymbol{\beta}^\wedge$
1	1	-1	2	8.00	-6.00	2.55	-5.45	-0.55	$\beta_0^\wedge =$	7.4545
2	1	-1	3	8.00	-5.00	2.55	-5.45	0.45	$\beta_1^\wedge =$	4.9091
3	1	0	6	8.00	-2.00	7.45	-0.55	-1.45		
4	1	0	7	8.00	-1.00	7.45	-0.55	-0.45		
5	1	0	8	8.00	0.00	7.45	-0.55	0.55		
6	1	0	9	8.00	1.00	7.45	-0.55	1.55		
7	1	1	10	8.00	2.00	12.36	4.36	-2.36		
8	1	1	12	8.00	4.00	12.36	4.36	-0.36		
9	1	1	15	8.00	7.00	12.36	4.36	2.64		
		5.00		8.0000	136.0000		117.8182	18.1818		
		$\Sigma X^2$		$Y^-$	$S_T$		$S_R$	$S_e$		
自由度 $df$			9	1	9-1=8	2	2-1=1	9-2=7		2

計算結果を, 表 11.3 の分散分析表にまとめる. 分散分析表の自由度については, 各種の便宜的な説明が行なわれているが, 自由度の本質が把握しづらいので, ここに示したように偏差平方和をベースにした計算式に対応する自由度の計算法を示す.

表 11.3 回帰に対する分散分析表 (表 4.4 再掲)

要因	平方和	自由度	平均平方	$F$ 値	$p$ 値
回帰 $S_R$	117.8182	2-1=1	117.8182	45.3600	0.0003
誤差 $S_e$	18.1818	9-2=7	2.5974		
全体 $S_T$	136.0000	9-1=8			

## パラメータの共分散行列

分散分析表の誤差の平均平方（誤差分散  $\hat{\sigma}^2$ ）を用いて、回帰パラメータ  $\hat{\beta}_1$  の分散  $Var(\hat{\beta}_1)$  は、正規方程式の解を用いて、式 (4.23) および式 (4.22) から

$$\begin{aligned} Var(\hat{\beta}_0) &= \frac{\hat{\sigma}^2 \sum X_i^2}{n \sum (X_i - \bar{X})^2} \\ &= \frac{2.5974 \times 5.00}{9 \times 4.8889} = 0.2952 \end{aligned}$$

$$\begin{aligned} Var(\hat{\beta}_1) &= \frac{\hat{\sigma}^2}{\sum (X_i - \bar{X})^2} \\ &= \frac{2.5974}{4.8889} = 0.5313 \end{aligned}$$

が得られる。また、式 (4.25) から、それぞれの分散は、

$$\Sigma(\hat{\beta}) = (X^T X)^{-1} \hat{\sigma}^2 = \begin{bmatrix} Var(\hat{\beta}_0) & Cov(\hat{\beta}_0, \hat{\beta}_1) \\ Cov(\hat{\beta}_0, \hat{\beta}_1) & Var(\hat{\beta}_1) \end{bmatrix} = \begin{bmatrix} \frac{\sum X_i^2}{n \sum (X_i - \bar{X})^2} & \frac{-\bar{X}}{\sum (X_i - \bar{X})^2} \\ \frac{-\bar{X}}{\sum (X_i - \bar{X})^2} & \frac{1}{\sum (X_i - \bar{X})^2} \end{bmatrix} \hat{\sigma}^2$$

パラメータの共分散行列の対角要素であることも示した。実際の計算は、

$$\begin{aligned} \Sigma(\hat{\beta}) &= \begin{bmatrix} 0.1136 & -0.0227 \\ -0.0227 & 0.2045 \end{bmatrix} \begin{bmatrix} 2.5974 \\ \sigma^2 \end{bmatrix} \\ &= \begin{bmatrix} 0.2952 & -0.0590 \\ -0.0590 & 0.5313 \end{bmatrix} \\ &= \Sigma(\hat{\beta}^{\wedge}) = (X^T X)^{-1} \sigma^{\wedge 2} \end{aligned}$$

となり、 $Var(\hat{\beta}_0) = 0.2952$ 、 $Var(\hat{\beta}_1) = 0.5313$  が得られる。それらの平方根から  $SE$  を計算しパラメータに関する  $t$  検定が行なう。この行列計算の方法は、変数の数が増えても同じであり、偏差平方和をベースにした計算手順よりも簡潔である。詳しくは第 12.3 節で示す。

表 11.4 回帰パラメータの推定値（表 4.2 再掲）

項	推定値	分散	$SE$	$t$ 値	$p$ 値
$\beta_0^{\wedge}$	7.4545	0.2952	0.5433	13.72	0.0000
$\beta_1^{\wedge}$	4.9091	0.5313	0.7289	6.73	0.0003

## スチューデント化残差

通常の回帰分析において、残差の検討の重要性は常に強調されている。代表的なのは、図 11.1 に示す予測値  $\hat{y}_i$  に対する残差  $\hat{\varepsilon}_i = y_i - \hat{y}_i$  のプロットおよびスチューデント化残差プロットである。この残差プロットは、予測値  $\hat{y}_i$  が大きくなるにつれ扇型に広がっている。したがって、推定された回帰直線に対して誤差が均一とはみなせないため、変数変換などで、誤差が均一になるような変換を検討することになる。

スチューデント化残差  $\hat{\varepsilon}'_i$  は、残差  $\hat{\varepsilon}_i$  の標準誤差で割って基準化したものであり、残差自体について統計的な考察ができる。図 11.1 右に示すように、扇型に広がってはいないが残差  $\hat{\varepsilon}'_i$  は、 $(-2 \sim +2)$  の範囲に入っており、通常の回帰分析を適用しても差し支えないと判断される。

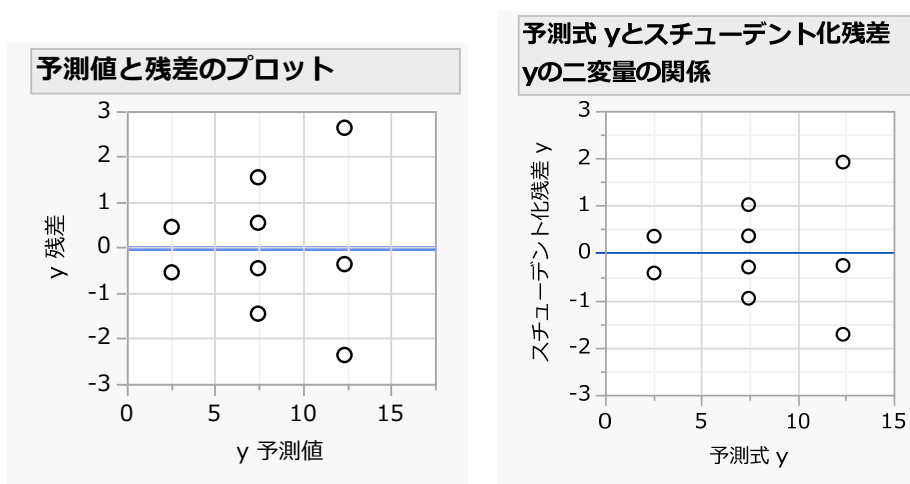


図 11.1 通常の回帰分析における残差およびスチューデント化残差プロット

## テコ比・ハット行列

スチューデント化残差  $\hat{\varepsilon}'_i$  は、残差  $\hat{\varepsilon}_i$  の分散  $Var(\hat{\varepsilon}_i) = \hat{\sigma}^2(1 - h_{ii})$  を考慮したものである。ここで  $h_{ii}$  は、通称テコ比ともいわれており、スチューデント化残差  $\hat{\varepsilon}'_i$  は、ハット行列の対角要素  $h_{ii}$  を用いて基準化したもので

$$\begin{aligned}\hat{\varepsilon}'_i &= \frac{\hat{\varepsilon}_i}{\sqrt{Var(\hat{\varepsilon}_i)}} \\ &= \frac{\hat{\varepsilon}_i}{\sqrt{\hat{\sigma}^2(1 - h_{ii})}}\end{aligned}$$

として計算されている。ハット行列  $\mathbf{H}$  は、 $\mathbf{Y}$  の推定値を求める式  $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$  に  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$  を代入し、

$$\begin{aligned}\hat{\mathbf{Y}} &= \mathbf{X}\hat{\boldsymbol{\beta}} \\ &= \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}\end{aligned}$$



最後の  $\mathbf{Y}$  を除いた行列

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

で定義されている。実際に計算すると次のような結果となる。まず、 $(\mathbf{X}^T \mathbf{X})^{-1}$  の結果を示す。次に、 $(\mathbf{X}^T \mathbf{X})^{-1}$  を挟んで  $\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$  の行列計算を行う。このようにハット行列  $\mathbf{H}$  の計算は、ストレスなく行えるのであるが、その対角要素を列ベクトル化する関数が Excel がないために、この式は使えない。

$\mathbf{X}^T$									$\mathbf{X}$		$\mathbf{X}^T \mathbf{X}$		$(\mathbf{X}^T \mathbf{X})^{-1}$	
1	1	1	1	1	1	1	1	1	1	-1	9	1	0.1136	-0.0227
-1	-1	0	0	0	0	1	1	1	1	-1	1	5	-0.0227	0.2045
=Transpose( $\mathbf{X}$ の範囲)									1	0			=Minverse( $\mathbf{X}^T \mathbf{X}$ の範囲)	
									1	0				
									1	0	=Mmult( $\mathbf{X}^T$ の範囲, $\mathbf{X}$ の範囲)			
									1	0				
									1	1				
									1	1				
									1	1				
(2×9)									(9×2)		(2×2)		(2×2)	

$X$		$(X^T X)^{-1}$		$X^T$									$X(X^T X)^{-1} X^T$											
1	-1	0.1136	-0.0227	1	1	1	1	1	1	1	1	1	1	0.36	0.36	0.14	0.14	0.14	0.14	-0.1	-0.1	-0.1		
1	-1	-0.0227	0.2045	-1	-1	0	0	0	0	1	1	1	1	0.36	0.36	0.14	0.14	0.14	0.14	-0.1	-0.1	-0.1		
1	0													0.14	0.14	0.11	0.11	0.11	0.11	0.09	0.09	0.09		
1	0													0.14	0.14	0.11	0.11	0.11	0.11	0.09	0.09	0.09		
1	0													0.14	0.14	0.11	0.11	0.11	0.11	0.09	0.09	0.09		
1	0													0.14	0.14	0.11	0.11	0.11	0.11	0.09	0.09	0.09		
1	1													-0.1	-0.1	0.09	0.09	0.09	0.09	0.27	0.27	0.27		
1	1													-0.1	-0.1	0.09	0.09	0.09	0.09	0.27	0.27	0.27		
1	1													-0.1	-0.1	0.09	0.09	0.09	0.09	0.27	0.27	0.27		
(9×2)		(2×2)		(2×9)									(9×9)											
=Mmult(Mmult( $X$ の範囲, $(X^T X)^{-1}$ の範囲), $X^T$ の範囲)																								

ハット行列  $\mathbf{H}$  対角要素を  $h_{ii}$  としたとき、デザイン行列  $\mathbf{X}$  の  $i$  行目のベクトルを  $\mathbf{x}_i$  を用いることにより、テコ比  $h_{ii}$  は、

$$h_{ii} = \mathbf{x}_i (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i^T$$

として求めることができる。この式で、それぞれの  $h_{ii}$  を計算することにより、ハット行列  $\mathbf{H}$  を計算しなくとも次式で示すように、対角要素をベクトル化することができる。

$\mathbf{x}_1$		$(\mathbf{X}^T \mathbf{X})^{-1}$	$\mathbf{x}_1^T$	$h_{11}$
1	-1	0.1136 -0.0227	1	0.3636
		-0.0227 0.2045	-1	
:				
$\mathbf{x}_8$		$(\mathbf{X}^T \mathbf{X})^{-1}$	$\mathbf{x}_8^T$	$h_{88}$
1	1	0.1136 -0.0227	1	0.2727
		-0.0227 0.2045	1	

## テコ比の活用

単純な誤差のプロットに加えて、テコ比を考慮したスチューデント化残差（標準化残差）による検討も有益である。なお、テコ比とハット行列の意味付けについては、野沢（1992）、「テコ比とハット行列」が詳しい。

テコ比  $h_{ii}$  は、回帰の 95%信頼区間を求めるための分散  $Var(\hat{y}_i)$  の計算にも関係している。式 (4.36) から、

$$\begin{aligned}
 Var(\hat{y}_i) &= \mathbf{x}_i \boldsymbol{\Sigma}(\hat{\boldsymbol{\beta}}) \mathbf{x}_i^T \\
 &= \mathbf{x}_i [(\mathbf{X}^T \mathbf{X})^{-1} \hat{\sigma}^2] \mathbf{x}_i^T \\
 &= h_{ii} \hat{\sigma}^2
 \end{aligned}$$

が導かれる。個々のデータの分散  $\hat{\sigma}^2$  に対しテコ比  $h_{ii}$  は、回帰の推定値の分散  $Var(\hat{y}_i)$  を求めるための割引係数として解される。表 11.5 に Excel でテコ比  $h_{ii}$  を計算し、スチューデント化残差  $\hat{\varepsilon}'_i$  および分散  $Var(\hat{y}_i)$  を計算した結果を示す。

表 11.5 テコ比を用いたスチューデント化残差

$i$	$\mathbf{X}$		$\mathbf{Y}$	$\mathbf{Y}^\wedge$	残差 $\boldsymbol{\varepsilon} = \mathbf{Y} - \mathbf{Y}^\wedge$	テコ比 $h_{ii}$	ス化残差 $\boldsymbol{\varepsilon}'$	残差比 $\boldsymbol{\varepsilon}' / \boldsymbol{\varepsilon}$	分散 $Var(\mathbf{y}^\wedge)$
1	1	-1	2	2.5455	-0.5455	0.3636	-0.4243	0.7778	0.9445
2	1	-1	3	2.5455	0.4545	0.3636	0.3536	0.7778	0.9445
3	1	0	6	7.4545	-1.4545	0.1136	-0.9586	0.6591	0.2952
4	1	0	7	7.4545	-0.4545	0.1136	-0.2996	0.6591	0.2952
5	1	0	8	7.4545	0.5455	0.1136	0.3595	0.6591	0.2952
6	1	0	9	7.4545	1.5455	0.1136	1.0185	0.6591	0.2952
7	1	1	10	12.3636	-2.3636	0.2727	-1.7197	0.7276	0.7084
8	1	1	12	12.3636	-0.3636	0.2727	-0.2646	0.7276	0.7084
9	1	1	15	12.3636	2.6364	0.2727	1.9182	0.7276	0.7084
	9.00	1.00	0.1136	-0.0227	72.0000	$\beta_0 =$	7.4545	$\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon} =$	18.1818
	1.00	5.00	-0.0227	0.2045	32.0000	$\beta_1 =$	4.9091	$\hat{\sigma}^2 =$	2.5974
	$\mathbf{X}^T \mathbf{X}$		$(\mathbf{X}^T \mathbf{X})^{-1}$		$\mathbf{X}^T \mathbf{Y}$	$(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$			

実際の計算過程を  $i=1$  の場合について示す．まずテコ比  $h_{11}$  を計算し，

$$h_{11} = \begin{array}{|c|c|} \hline 1 & -1 \\ \hline \end{array} \begin{array}{|c|c|} \hline 0.1136 & -0.0227 \\ \hline \end{array} \begin{array}{|c|c|} \hline 1 & -1 \\ \hline \end{array} = \begin{array}{|c|} \hline 0.3636 \\ \hline \end{array}$$

$$\begin{array}{|c|} \hline \mathbf{x}_1 \\ \hline \end{array} \begin{array}{|c|c|} \hline -0.0227 & 0.2045 \\ \hline \end{array} \begin{array}{|c|} \hline (X^T X)^{-1} \\ \hline \end{array} \begin{array}{|c|} \hline \mathbf{x}_1^T \\ \hline \end{array} \begin{array}{|c|} \hline h_{11} \\ \hline \end{array}$$

スチューデント化残差  $\hat{\varepsilon}_1'$  を求める．

$$\begin{aligned} \hat{\varepsilon}_1' &= \frac{\hat{\varepsilon}_1}{\sqrt{\hat{\sigma}^2(1-h_{11})}} \\ &= \frac{-0.5455}{\sqrt{2.5974 \times (1-0.3636)}} = -0.4243 \end{aligned}$$

テコ比  $h_{11}$  を用いて回帰の推定値  $\hat{y}_1$  の分散  $Var(\hat{y}_1)$  の計算もできる．

$$\begin{aligned} Var(\hat{y}_1) &= h_{11}\hat{\sigma}^2 \\ &= 0.3636 \times 2.5974 = 0.9445 \end{aligned}$$

なお，図 11.1 右に示したスチューデント化残差プロットは，JMP のスチューデント化残差プロットに X 軸に予測値を指定できないので，スチューデント化残差をファイルに出力して，別途「二変量の関係」を使って作図したものである．

テコ比は，回帰直線の推定値の分散を誤差分散  $\hat{\sigma}^2$  に対して割引係数としても理解される．回帰直線の中心部は小さく，外側に向かって大きくなり，推定値の分散が大きくなり 95%信頼区間の幅の変化を示す統計量とも解される．

## Excel の「標準残差」に対する使用上の注意

Excel の回帰分析の活用は，分散分析表および回帰パラメータの推定値に関してデザイン行列の計算の煩わしさを軽減するために有益であることを示したきた．さらに，Excel の回帰分析で「残差」に加えて「標準残差」を Excel シートに出力することができる．「標準残差」はスチューデント化残差（標準化残差）と紛らわしいが，テコ比を含まない計算であり，別物である．表 11.5 の場合では，全ての残差  $\hat{\varepsilon}_i$  に 0.6633 を掛けている．これは，分散  $\hat{\sigma}^2$  の平方根の逆数 0.6205 に近い値であるが，どのような計算なのかは不明である．いずれにしても，スチューデント化残差（標準化残差）ではないことに注意が必要である．

Excel の統計解析については，正確性が欠けるとの指摘があることは十分に承知しており，「標準残差」もその一例であろう．もっとひどい事例は，第 10.3 節で例示した，折れ線グラフの誤差範囲の設定のいい加減さであり，それを回避する方法を知る必要がある．なお，行列計算などの精度で問題になった経験はないが，計算過程の脆弱性を常に認識し，他のソフトでの検証を怠ってはならない．

### 11.3. ポアソン回帰におけるデビアンズ・逸脱度

対数尤度を用いた恒等リンクのポアソン回帰は、回帰式を

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad \varepsilon_i \sim \text{ポアソン分布}$$

としたときに、対数尤度  $\ln L_{\text{回帰}}$  を最大にするような回帰パラメータ  $\hat{\beta}_0$  と  $\hat{\beta}_1$

$$\ln L_{\text{回帰}} = \sum_i \ln[\text{Poisson.dist}(y_i, \hat{\beta}_0 + \hat{\beta}_1 x_i, \text{false})]$$

を推定する． $\ln L_{\text{回帰}}$  を最大化するために、第 1.4 節および第 5 章で反復重み付き回帰を用いる方法、第 2 章で対数尤度関数の 2 階の偏微分式行列を用いる方法、第 1.9 節で Excel のソルバーにより  $\ln L_{\text{回帰}}$  が最大になるような  $\hat{\beta}_0$  と  $\hat{\beta}_1$  を直接求める方法を示してきた．ここでは、簡便な Excel のソルバーを用いる方法を使う．

#### デビアンズ・カイ 2 乗

ポアソン回帰では、第 1.9 節で示したように飽和モデル、完全（最大）モデル、縮小モデル、切片のみの場合には（null モデル）など幾つかの「モデル」が登場する．そして、逸脱度/デビアンズは、飽和モデルと各モデルとのマイナス 2 倍の対数尤度の差で定義されている．

「飽和モデル」の概念が通常の回帰分析にはないので、理解に苦しむことになる．表 11.6 に示すように、飽和モデルの対数尤度  $\ln L_{\text{飽和}}$  は、 $y_i$  のポアソン分布の確率を求めるための推定値として  $\hat{y}_i = y_i$  のように自分自身  $y_i$  を用いて、

$$\text{飽和モデル： } \hat{y}_i = y_i, \quad \left\{ \begin{array}{l} \ln L_{\text{飽和}} = \sum_i \ln(\text{Poisson.dist}(y_i, y_i, \text{false})) \\ = \ln(0.2707) + \ln(0.2240) + \cdots + \ln(0.1024) \\ = -1.3069 - 1.4959 - \cdots - 2.2785 \\ = -17.0566 \end{array} \right.$$

としてポアソン分布の確率を計算している．通常の回帰分析は、偏差平方和の計算で組み立てられているので、無理に飽和モデルを考えても、次のようにゼロなので

$$S_{\text{飽和}} = \sum_i (y_i - y_i)^2 = 0$$

何の役にも立たない．したがって、飽和モデルの概念がない．（完全 or 最大 or 回帰）モデルは、 $\hat{\beta}_0$  と  $\hat{\beta}_1$  を使った回帰モデルで、 $\ln L_{\text{回帰}}$  に変えて  $\ln L_{\text{完全}}$  とし、

$$\text{完全モデル， } \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i : \quad \left\{ \begin{array}{l} \ln L_{\text{完全}} = \sum_i \ln(\text{Poisson.dist}(y_i, \hat{\beta}_0 + \hat{\beta}_1 x_i, \text{false})) \\ = \ln(0.2557) + \ln(0.2144) + \cdots + \ln(0.0791) \\ = -1.3639 - 1.5397 - \cdots - 2.5366 \\ = -18.0039 \end{array} \right.$$

となる．これは、表 11.7 に示すように通常の回帰分析の「誤差平方和  $S_e$ 」に該当する．

表 11.6 最尤法によるポアソン回帰

			縮小 (null) モデル			完全 (最大) モデル			飽和モデル		
			$y^{\wedge} = \hat{\beta}_0$			$y^{\wedge} = \hat{\beta}_0 + \hat{\beta}_1 x$			$y^{\wedge} = y$		
			$\hat{\beta}_0 = 8.0000$			$\hat{\beta}_0 = 7.4516$ $\hat{\beta}_1 = 4.9353$					
			$\ln L_{\text{縮小}} = -26.2669$			$\ln L_{\text{完全}} = -18.0039$			$\ln L_{\text{飽和}} = -17.0566$		
			$\hat{\beta}_0$	確率	対数尤度	回帰	確率	対数尤度	$y$	確率	対数尤度
$i$	$x$	$y$	$y^{\wedge}$	$P$	$\ln L_i$	$y^{\wedge}$	$P$	$\ln L_i$	$y^{\wedge}$	$P$	$\ln L_i$
1	-1	2	8.00	0.0107	-4.5343	2.52	0.2557	-1.3639	2	0.2707	-1.3069
2	-1	3	8.00	0.0286	-3.5534	2.52	0.2144	-1.5397	3	0.2240	-1.4959
3	0	6	8.00	0.1221	-2.1026	7.45	0.1380	-1.9803	6	0.1606	-1.8287
4	0	7	8.00	0.1396	-1.9691	7.45	0.1469	-1.9178	7	0.1490	-1.9038
5	0	8	8.00	0.1396	-1.9691	7.45	0.1369	-1.9888	8	0.1396	-1.9691
6	0	9	8.00	0.1241	-2.0869	7.45	0.1133	-2.1776	9	0.1318	-2.0268
7	1	10	8.00	0.0993	-2.3100	12.39	0.0978	-2.3249	10	0.1251	-2.0786
8	1	12	8.00	0.0481	-3.0339	12.39	0.1137	-2.1744	12	0.1144	-2.1683
9	1	15	8.00	0.0090	-4.7076	12.39	0.0791	-2.5366	15	0.1024	-2.2785

(縮小 or 切片 or Null) モデルは,

$$\text{縮小モデル, } \hat{y}_i = \hat{\beta}_0 : \begin{cases} \ln L_{\text{縮小}} = \sum_i \ln(\text{Poisson.dist}(y_i, \hat{\beta}_0, \text{false})) \\ = \ln(0.0107) + \ln(0.0286) + \dots + \ln(0.0090) \\ = -4.5343 - 3.5534 - \dots - 4.7076 \\ = -26.2669 \end{cases}$$

として計算される. 通常の回帰分析の「総平方和  $S_T$ 」に該当する.

回帰の平方和  $S_R$  は,  $S_R = S_T - S_e$  で求められると同様に, 傾き  $\hat{\beta}_1$  に対する 2 倍の対数尤度は, 差分

$$\text{差分} : \begin{cases} 2 \ln L_R = 2(\ln L_{\text{完全}} - \ln L_{\text{縮小}}) \\ = 2 \times [-18.0039 - (-26.2669)] \\ = 16.5260 \end{cases}$$

として求められる. この 2 倍の対数尤度に対する検定統計量は, それぞれの自由度の差のカイ 2 乗分布に従うことにより有意差検定が行なえる. 完全モデル ( $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ ) と縮小モデル ( $\hat{y}_i = \hat{\beta}_0$ ) の対数尤度の差 2 倍なので,  $\hat{\beta}_1$  に対する尤度比検定統計量となり, カイ 2 乗検定が行なえる.

通常の回帰分析では,  $S_R$  をその自由度  $df_R$  で割った平均平方を,  $S_e$  をその自由度  $df_e$  で割った平均平方 (誤差分散) との分散比

$$F = \frac{S_R / df_R}{S_e / df_e}$$

が、分母の自由度  $df_R$ 、分子の自由度  $df_e$  の  $F$  分布に従うことから有意差検定を行なっている。

JMP によるポアソン回帰の結果を表 11.7 に示す。通常の回帰分析の「分散分析表」と対比して、(差分, 完全, 縮小) の意味を理解してもらいたい。その後に続く「適合度統計量」は通常の「分散分析表」にはない概念である。しいて言えば、仮定した正規分布からの乖離度の統計量であろうか。「適合度統計量」の欄の「デビアンス」の行のカイ 2 乗値=1.8947 となっているのは、飽和モデルの対数尤度  $\ln L_{\text{飽和}} = -17.0566$  と  $\ln L_{\text{完全}} = -18.0039$  との差の 2 倍で

$$\begin{aligned}\text{デビアンス} \cdot \text{カイ2乗} &= 2(\ln L_{\text{飽和}} - \ln L_{\text{完全}}) \\ &= 2 \times [-17.0566 - (-18.0039)] \\ &= 1.8947\end{aligned}$$

と、デビアンスが計算されている。自由度は、 $df_{\text{飽和}} = 9$  と  $df_{\text{完全}} = 2$  の差から 7 となっている。このデビアンスが、第 9.5 節で示した R 言語の一般化線形モデルで出力されている Residual deviance に対応する。

表 11.7 通常の回帰とポアソン回帰の対比

通常の回帰分析		
要因		自由度
回帰	$S_R$	2-1=1
誤差	$S_e$	9-2=7
全体	$S_T$	9-1=8

モデル全体の検定				
モデル	(-1)*対数尤度	尤度比カイ2乗	自由度	p値
差分	8.2630	16.5260	1	<.0001*
完全	18.0039			
縮小	26.2669			
適合度統計量	カイ2乗	自由度	p値	
Pearson	1.8944	7	0.9655	
デビアンس	1.8947	7	0.9654	
AICc				
42.0078				

### Pearson・カイ 2 乗

Pearson の適合度統計量は、反応変数  $y_i$  と推定値  $\hat{y}_i$  の差の 2 乗をその分散  $Var(y_i) = \hat{y}_i$  で割って加えたもので

$$\begin{aligned}\text{Pearson} \cdot \text{カイ2乗} &= \sum_i \frac{(y_i - \hat{y}_i)^2}{Var(\hat{y}_i)} = \sum_i \frac{(y_i - \hat{y}_i)^2}{\hat{y}_i} \\ &= \frac{(2 - 2.52)^2}{2.52} + \frac{(3 - 2.52)^2}{2.52} + \dots + \frac{(15 - 12.39)^2}{12.39} \\ &= 1.8944\end{aligned}$$

として計算されている。これらのデビアンスおよび Pearson のカイ 2 乗値を自由度で除したのが過分散の調整パラメータである。この例では、明らかにカイ 2 乗値が自由度 7 より小さいので、過分散が起きていないと判断される。

通常の回帰分析の分散分析表の「要因」の欄の（回帰，誤差，全体）との表記は，ポアソン回帰の「モデル」の表記（差分，完全，縮小）と全く異なるので，それらの対応関係を関連付けることによりによって理解を深めてもらいたい．ただし，「回帰」が「差分」に対応し，「誤差」が「完全」に，「全体」が「縮小」にそれぞれ対応している．このような用語の対応付けを行なったとしても，同義語としては全く認識されないで，それらの対数尤度の定義を理解した上での説明を加えながら注意深く使うことが必要と思われる．

## AICc

AICc は，修正済み赤池の情報量基準で， $k$  をパラメータ数 2 とし

$$\begin{aligned} \text{AICc} &= -2 \ln L_{\text{完全}} + 2k + \frac{2k(k+1)}{n-k-1} \\ &= -2 \times (-18.0039) + 2 \times 2 + \frac{2 \times 2 \times (2+1)}{9-2-1} \\ &= 42.0078 \end{aligned}$$

と計算されている．AICc は，パラメータ数が異なるモデルを比較する際に役に立つ．

表 11.8 のパラメータに関する尤度比検定の「項」の欄の「x」に対する尤度比カイ 2 乗は，16.5260 であり，モデル全体の検定の差分の尤度比カイ 2 乗 16.5260 に一致する．なお，標準誤差は，第 4.2 節で詳しく述べたように対数尤度関数の 2 階の偏微分行列  $\mathbf{H}$  の負の逆行列  $(-\mathbf{H})^{-1}$  がパラメータの共分散行列なので，その対角要素の平方根である．

表 11.8 パラメータに関する尤度比検定

パラメータ推定値				
項	推定値	標準誤差	尤度比カイ2乗	p値
切片	7.4516	0.8842	71.0299	<.0001*
x	4.9353	1.0915	16.5260	<.0001*

## 11.4. ポアソン回帰における 4 種の残差

ポアソン回帰の場合は、恒等リンクの場合でも対数リンクの場合でも予測値の大きさに比例して分散が大きくなるので、単純な残差プロットによる残差の検討は、不適切である。そのために、ピアソン残差、スチューデント化ピアソン残差、デビアンズ残差、あるいは、スチューデント化デビアンズ残差などが JMP の一般化線形モデル、SAS の GENMOD プロシジャで提供されている。R では、デビアンズ残差が使われている。

### デビアンズ残差

表 11.7 に示した JMP の「モデル全体の検定」での「適合度統計量」としてのデビアンズは、飽和モデルの  $\ln L_{\text{飽和}} = -17.0566$  と完全モデル  $\ln L_{\text{完全}} = -18.0039$  との差の 2 倍で

$$\begin{aligned}\text{デビアンズ} &= 2(\ln L_{\text{飽和}} - \ln L_{\text{完全}}) \\ &= 2 \times [-17.0566 - (-18.0039)], \quad df = 7 \\ &= 1.8947\end{aligned}$$

であることを示した。これが有意であれば、ポアソン回帰で取り上げた説明変数では説明しきれない誤差変動が残っていることを意味している。他に追加できる変数がなければ、誤差分布がポアソン分布に対し過分散となっていることを意味する。過分散となっている場合は、パラメータの推定値の標準誤差が小さくなり、結果を過大評価することになる。そのために、過分散パラメータで標準誤差を大きくする修正、あるいは、負の 2 項分布（ガンマ・ポアソン分布）などを仮定した解析が必要となる。

デビアンズは、飽和モデルと完全モデルの対数尤度の差の 2 倍であり、それぞれの対数尤度は、個々のデータの対数尤度  $\ln L_i$  の和であることを表 11.6 で示した。デビアンズ残差 (Deviance Residuals) は、飽和モデルと完全モデルの個々の対数尤度の差の 2 倍の平方根として定義されている。ただし、符号が全てプラスなので残差とは言えないので、完全モデルの推定値  $\hat{y}_{\text{完全},i}$  と飽和モデルの推定値  $\hat{y}_{\text{飽和},i} = y_i$  の差の符号を付けてデビアンズ残差  $\varepsilon_i^{(D)}$  とする。

$$d_i = 2[\ln(L_{\text{飽和},i}) - \ln(L_{\text{完全},i})] \quad (11.1)$$

$$\varepsilon_i^{(D)} = \text{Sign}(\hat{y}_{\text{飽和},i} - \hat{y}_{\text{完全},i})\sqrt{d_i} \quad (11.2)$$

なお、個々のデビアンズ  $d_i$  は、簡略化した計算公式 (11.3)

$$\begin{aligned}d_i &= 2 \left[ \ln \left( \frac{\hat{y}_{\text{飽和},i}^{y_i} e^{-\hat{y}_{\text{飽和},i}}}{y_i!} \right) - \ln \left( \frac{\hat{y}_{\text{完全},i}^{y_i} e^{-\hat{y}_{\text{完全},i}}}{y_i!} \right) \right] \\ &= 2 \left[ y_i \ln(y_i) - y_i - \ln(y_i!) - y_i \ln(\hat{y}_{\text{完全},i}) + \hat{y}_{\text{完全},i} + \ln(y_i!) \right] \\ &= 2 \left[ y_i \ln \left( \frac{y_i}{\hat{y}_{\text{完全},i}} \right) - (y_i - \hat{y}_{\text{完全},i}) \right]\end{aligned} \quad (11.3)$$



が、一般的に計算公式として用いられている．ただし、この計算式からでは何を意味しているのか、推測しがたい．元の対数尤度での定義式による理解を勧める．

実際に  $i=1$  の場合の計算を次に示し、全ての  $i$  についての結果を表 11.9 に示す．

$$\begin{aligned} d_1 &= 2[(-1.3069) - (-1.3639)] \\ &= 2 \times 0.0570 = 0.1140 \end{aligned}$$

$$\text{Sign}(\hat{y}_{\text{飽和},1} - \hat{y}_{\text{完全},1}) = \text{Sign}(2.0 - 2.5163) = \text{マイナス}$$

$$\begin{aligned} \varepsilon_1^{(D)} &= \text{Sign}(\hat{y}_{\text{飽和},1} - \hat{y}_{\text{完全},1}) \sqrt{d_1} \\ &= -\sqrt{0.1140} = -0.3377 \end{aligned}$$

表 11.9 デビアンズ残差

				完全(最大)モデル			飽和モデル			デビアンズ残差		
					$\hat{\beta}_0 =$	7.4516						
					$\hat{\beta}_1 =$	4.9353					平方和 =	1.8947
					$\ln L_{\text{完全}} =$	-18.0039		$\ln L_{\text{飽和}} =$	-17.0566	尤度の差	平方根	残差
$i$	$X$		$y$	$y_{\text{完全}}^{\wedge}$	$P_{\text{完全}}$	$\ln L_{\text{完全},i}$	$y_{\text{飽和}}^{\wedge}$	$P_{\text{飽和}}$	$\ln L_{\text{飽和},i}$	$d$	$\sqrt{(d)}$	$\varepsilon^{(D)}$
1	1	-1	2	2.5163	0.2557	-1.3639	2	0.2707	-1.3069	0.1140	0.3377	-0.3377
2	1	-1	3	2.5163	0.2144	-1.5397	3	0.2240	-1.4959	0.0875	0.2958	0.2958
3	1	0	6	7.4516	0.1380	-1.9803	6	0.1606	-1.8287	0.3032	0.5506	-0.5506
4	1	0	7	7.4516	0.1469	-1.9178	7	0.1490	-1.9038	0.0279	0.1672	-0.1672
5	1	0	8	7.4516	0.1369	-1.9888	8	0.1396	-1.9691	0.0394	0.1985	0.1985
6	1	0	9	7.4516	0.1133	-2.1776	9	0.1318	-2.0268	0.3015	0.5491	0.5491
7	1	1	10	12.3869	0.0978	-2.3249	10	0.1251	-2.0786	0.4927	0.7019	-0.7019
8	1	1	12	12.3869	0.1137	-2.1744	12	0.1144	-2.1683	0.0122	0.1105	-0.1105
9	1	1	15	12.3869	0.0791	-2.5366	15	0.1024	-2.2785	0.5161	0.7184	0.7184

デビアンズ残差  $\varepsilon_i^{(D)}$  の平方和は、

$$\begin{aligned} \text{デビアンズ} &= \sum_{i=1}^9 \left[ \varepsilon_i^{(D)} \right]^2 \\ &= (-0.3377)^2 + 0.2958^2 + \dots + 0.7184^2 = 1.8947 \end{aligned}$$

となり、デビアンズ・カイ 2 乗値に

$$\begin{aligned} \text{デビアンズ・カイ2乗} &= 2(\ln L_{\text{飽和}} - \ln L_{\text{完全}}) \\ &= 2 \times [-17.0566 - (-18.0039)] = 1.8947 \end{aligned}$$

一致することが確認される．もちろん、計算公式によっても、次のように

$$\begin{aligned} d_1 &= 2 \left[ y_1 \ln \left( \frac{y_1}{\hat{y}_{\text{完全},1}} \right) - (y_1 - \hat{y}_{\text{完全},1}) \right] \\ &= 2 \times \left[ 2 \times \ln \left( \frac{2}{2.5163} \right) - (2 - 2.5163) \right] = 0.1140 \end{aligned}$$

求め、一致することが確認できる．

## スチューデント化デビアンズ残差

スチューデント化デビアンズ残差は、反復重み付き回帰での重みを加味したハット行列  $\mathbf{H}'$  の対角要素テコ比  $h_{ii}'$  を用いる。恒等リンクの場合の重みは、推定値  $\hat{y}_i$  の逆数であるので、対角要素にそれぞれの重を持つ行列を  $\hat{\mathbf{W}}$  とする。推定値  $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$  のパラメータ推定値  $\hat{\boldsymbol{\beta}}$  を重み付き回帰の計算式  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \hat{\mathbf{W}} \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{W}} \mathbf{Y}$  を代入すると、

$$\begin{aligned}\hat{\mathbf{Y}} &= \mathbf{X}\hat{\boldsymbol{\beta}} \\ &= \mathbf{X}(\mathbf{X}^T \hat{\mathbf{W}} \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{W}} \mathbf{Y}\end{aligned}\quad (11.4)$$

となる。ハット行列  $\mathbf{H}'$  は、回帰の推定値を求める際に  $\hat{\mathbf{Y}}$  の式から最後の  $\mathbf{Y}$  を除いた行列

$$\mathbf{H}' = \mathbf{X}(\mathbf{X}^T \hat{\mathbf{W}} \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{W}} \quad (11.5)$$

で定義されている。重み行列  $\hat{\mathbf{W}}$  の対角要素からなるベクトルを  $\hat{\mathbf{w}}$  とし、デザイン行列  $\mathbf{X}$  の  $i$  行目のベクトルを  $\mathbf{x}_i$  とすれば、テコ比  $h_{ii}'$  は、

$$h_{ii}' = \mathbf{x}_i[(\mathbf{X} * \hat{\mathbf{w}})^T \mathbf{X}]^{-1} \mathbf{x}_i^T \hat{\mathbf{w}} \quad (11.6)$$

として求めることができる。スチューデント化デビアンズ残差  $\varepsilon_i^{(D)}$  は、

$$\hat{\varepsilon}_i^{(D)} = \frac{\hat{\varepsilon}_i^{(D)}}{\sqrt{1 - h_{ii}'}} \quad (11.7)$$

として求められる。表 11.10 にスチューデント化デビアンズ残差の計算結果を示す。デビアンズ残差は、表 11.9 に示した結果を用いる。重み  $\hat{w}_i$  を推定値  $\hat{y}_i$  の逆数とし、テコ比を求め、スチューデント化デビアンズ残差  $\hat{\varepsilon}_i^{(D)}$  が計算されている。

表 11.10 スチューデント化デビアンズ残差

			完全モデル		飽和モデル							
			$\hat{\beta}_0 =$	7.4516					デビアンズ残差		スチューデント化	
			$\hat{\beta}_1 =$	4.9353					平方和 = 1.8947		デビアンズ残差	
			$\ln L_{\text{完全}} =$	-18.0039	$\ln L_{\text{飽和}} =$	-17.0566	尤度の差	残差	重み $w$	テコ比	残差	
$i$	$X$		$y$	$y_{\text{完全}}^{\wedge}$	$\ln L_{\text{完全},i}$	$y_{\text{飽和}}^{\wedge}$	$\ln L_{\text{飽和},i}$	$d$	$\varepsilon^{(D)}$	$1/y_{\text{完全}}^{\wedge}$	$h_{ii}'$	$\varepsilon'^{(D)}$
1	1	-1	2	2.5163	-1.3639	2.0	-1.3069	0.1140	-0.3377	0.3974	0.4510	-0.4558
2	1	-1	3	2.5163	-1.5397	3.0	-1.4959	0.0875	0.2958	0.3974	0.4510	0.3993
3	1	0	6	7.4516	-1.9803	6.0	-1.8287	0.3032	-0.5506	0.1342	0.1049	-0.5820
4	1	0	7	7.4516	-1.9178	7.0	-1.9038	0.0279	-0.1672	0.1342	0.1049	-0.1767
5	1	0	8	7.4516	-1.9888	8.0	-1.9691	0.0394	0.1985	0.1342	0.1049	0.2098
6	1	0	9	7.4516	-2.1776	9.0	-2.0268	0.3015	0.5491	0.1342	0.1049	0.5804
7	1	1	10	12.3869	-2.3249	10.0	-2.0786	0.4927	-0.7019	0.0807	0.2261	-0.7979
8	1	1	12	12.3869	-2.1744	12.0	-2.1683	0.0122	-0.1105	0.0807	0.2261	-0.1256
9	1	1	15	12.3869	-2.5366	15.0	-2.2785	0.5161	0.7184	0.0807	0.2261	0.8167
										0.7817	0.4166	
										0.4166	1.1863	
										[(X * w)' X] <sup>-1</sup>		

実際の計算を  $i=1$  の場合について示す．スチューデント化デビアンズ残差は，

$$\begin{aligned}\hat{w}_1 &= 1 / \hat{y}_{\text{完全},1} \\ &= 1 / 2.5163 = 0.3974\end{aligned}$$

$$\begin{aligned}\Sigma(\hat{\beta}) &= [(X * \hat{w})^T X]^{-1} \\ &= \text{Minverse}(\text{Mmult}(\text{Transpose}(X \text{ の範囲} * \hat{w} \text{ の範囲}), X \text{ の範囲}))\end{aligned}$$

$$h_{11}' = \begin{array}{|c|c|} \hline 1 & -1 \\ \hline \end{array} \begin{array}{|c|c|} \hline 0.7817 & 0.4166 \\ \hline 0.4166 & 1.1863 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline -1 \\ \hline \end{array} \times \begin{array}{|c|} \hline 0.3974 \\ \hline \end{array} = \begin{array}{|c|} \hline 0.4510 \\ \hline \end{array}$$

$\begin{array}{|c|} \hline x_1 \\ \hline \end{array} \quad \begin{array}{|c|} \hline [(X * \hat{w})^T X]^{-1} \\ \hline \end{array} \quad \begin{array}{|c|} \hline x_1^T \\ \hline \end{array} \quad \begin{array}{|c|} \hline w_1 \\ \hline \end{array} \quad \begin{array}{|c|} \hline h_{ii}' \\ \hline \end{array}$

$$\begin{aligned}\hat{\varepsilon}_i^{(D)} &= \frac{\hat{\varepsilon}_i^{(D)}}{\sqrt{1 - h_{ii}'}} \\ &= \frac{-0.3377}{\sqrt{1 - 0.4510}} = -0.4558\end{aligned}$$

で求められる．

デビアンズ残差およびスチューデント化デビアンズ残差について，JMP のポアソン回帰で作成した残差プロットを図 11.2 に示す．デビアンズ残差に対してスチューデント化デビアンズ残差の方が，残差の絶対値が大きめになっていることが確認できる．ただし，実際の解析で複数の残差を用いることは非現実的であり，第 11.5 節を参考にして，選択してほしい．

自ら計算する場合には，手軽に計算できるデビアンズ残差であるが，統計的には，スチューデント化デビアンズ残差が望ましいと思われる．JMP のデフォルトの残差プロットは，スチューデント化デビアンズ残差が使用されているが，SAS/GENMOD プロシジャでは，ユーザの選択に任されている．

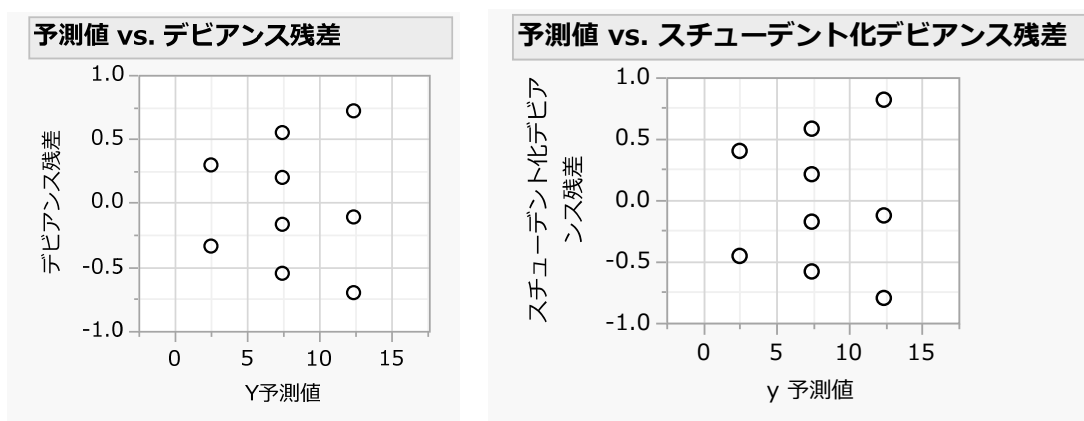


図 11.2 デビアンズ残差とスチューデント化デビアンズ残差の比較

## スチューデント化 Pearson 残差

第 11.2 節で通常の回帰分析の場合に残差  $\hat{\varepsilon}_i = y_i - \hat{y}_i$  に対してテコ比  $h_{ii}$  を用いてスチューデント化残差  $\hat{\varepsilon}_i'$  を求めた。ポアソン回帰の場合、単純な残差  $\hat{\varepsilon}_i$  は、推定値  $\hat{y}_i$  に比例して大きくなるので残差の検討に使えない、そこで、 $\hat{y}_i$  の標準誤差で基準化した Pearson 残差  $\hat{\varepsilon}_i^{(P)}$  が使われている。Pearson 残差  $\hat{\varepsilon}_i^{(P)}$  は、

$$\begin{aligned}\hat{\varepsilon}_i^{(P)} &= \frac{y_i - \hat{y}_i}{\sqrt{\text{Var}(\hat{y}_i)}} \\ &= \frac{y_i - \hat{y}_i}{\sqrt{\hat{y}_i}}\end{aligned}\quad (11.8)$$

として計算される。Pearson 残差  $\hat{\varepsilon}_i^{(P)}$  は過小評価になりがちなので、テコ比  $h_{ii}$  により調整したスチューデント化 Pearson 残差も使われている。重みは、 $\hat{w}_i = 1/\hat{y}_i$  であり、テコ比は、スチューデント化デビエンス残差で計算した表 11.10 と同じ式であり、スチューデント化 Pearson 残差  $\hat{\varepsilon}_i'^{(P)}$  は、スチューデント化デビエンス残差  $\varepsilon_i'^{(D)}$  の場合と同様に

$$\hat{\varepsilon}_i'^{(P)} = \frac{\hat{\varepsilon}_i^{(P)}}{\sqrt{1 - h_{ii}}}\quad (11.9)$$

となる。表 11.11 に Pearson 残差およびスチューデント化 Pearson 残差についての計算結果を示す。

表 11.11 Pearson 残差およびスチューデント化 Pearson 残差

					Pearson	スチューデント化 Pearson 残差		
				推定値	残差	重み	テコ比	残差
$i$	$X$		$y$	$\hat{y}$	$\varepsilon^{(P)}$	$w = 1/y$	$h_{ii}'$	$\varepsilon'^{(P)}$
1	1	-1	2	2.5163	-0.3255	0.3974	0.4510	-0.4393
2	1	-1	3	2.5163	0.3049	0.3974	0.4510	0.4115
3	1	0	6	7.4516	-0.5318	0.1342	0.1049	-0.5621
4	1	0	7	7.4516	-0.1654	0.1342	0.1049	-0.1749
5	1	0	8	7.4516	0.2009	0.1342	0.1049	0.2123
6	1	0	9	7.4516	0.5672	0.1342	0.1049	0.5995
7	1	1	10	12.3869	-0.6782	0.0807	0.2261	-0.7709
8	1	1	12	12.3869	-0.1099	0.0807	0.2261	-0.1250
9	1	1	15	12.3869	0.7425	0.0807	0.2261	0.8440

推定値は、表 11.10 の完全モデルの推定値、  
テコ比は、スチューデント化デビエンス残差でのテコ比に等しい。

実際の計算を  $i=1$  の場合について示す。Pearson 残差は、

$$\begin{aligned}\varepsilon_1^{(P)} &= \frac{y_1 - \hat{y}_1}{\sqrt{\hat{y}_1}} \\ &= \frac{2 - 2.5163}{\sqrt{2.5163}} \\ &= -0.3255\end{aligned}$$

スチューデント化デ Pearson 残差は,

$$\begin{aligned}\hat{w}_1 &= 1 / \hat{y}_1 \\ &= 1 / 2.5163 = 0.3974\end{aligned}$$

$$h_{11}' = \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & -1 & 0.7817 & 0.4166 & 1 & \times & 0.3974 = 0.4510 \\ \hline & & 0.4166 & 1.1863 & -1 & & \\ \hline \mathbf{x}_1 & & [(X^*w)^T X]^{-1} & & \mathbf{x}_1^T & & w_1 & h_{ii}' \\ \hline \end{array}$$

$$\begin{aligned}\hat{\varepsilon}_1^{(P')} &= \frac{\hat{\varepsilon}_i^{(P)}}{\sqrt{1-h_{11}'}} \\ &= \frac{-0.3255}{\sqrt{1-0.4510}} = -0.4393\end{aligned}$$

で求められる.

### SAS/GENMOD による各種の残差

SAS/GENMOD を用いて, スチューデント化 Pearson 残差を含め, これまでに示したポアソン回帰での 4 種の残差, (デビアンズ残差, スチューデント化デビアンズ残差, Pearson 残差, スチューデント化 Pearson 残差) を計算し, これまでの結果を検証する. SAS/GENMOD のオプションで「residual」が各種の残差の一括出力指示で, 「plots=stdresdev(xbeta)」が, スチューデント化デビアンズ残差を推定値  $\hat{y}_i$  に対する残差プロットの作成を指示している.

#### <<SAS/GENMOD によるポアソン回帰>>

```
Title "デビアンズ残差_a01.sas 2020/01/20 Y.Takahashi" ;
data d01 ;
  input x y @@ ;
datalines ;
-1 2 -1 3 0 6 0 7 0 8 0 9 1 10 1 12 1 15
;
proc genmod data=d01 plots=reschi(xbeta) plots=stdreschi(xbeta)
               plots=resdev(xbeta) plots=stdresdev(xbeta) ;
  model y = X / dist=poisson link= identity residual ;
run;
```

表 11.12 に示した SAS の残差の出力で, 「標準化」となっているのが「スチューデント化」の意味である「未加工残差」は,  $\hat{\varepsilon}_i = y_i - \hat{y}_i$  であるが, 右端の「尤度残差」は, スチューデント化デビアンズ残差とスチューデント化 Pearson 残差の両方を使って, それらの中間的な残差である [SAS Institute (2016), The GENMOD Procedure, The GENMOD Procedure :3164-3165]. SAS の残差の出力結果と, Excel で計算した 4 種の残差を照合し, 一致することが確認される.

このように沢山の残差があり，どれを使うか選択に窮するのである．最も簡便なのが Pearson 残差であるが，スチューデント化 Pearson 残差に比べて小さめに出るので，残差の大きさを少々過小評価となる．同様にデビアンズ残差もスチューデント化デビアンズ残差に比べて小さめになっている．

スチューデント化 Pearson 残差か，スチューデント化デビアンズ残差かの選択については，この事例ではどちらとも言い切れない．元々の計算過程で，対数尤度を使っているのだから，飽和モデルと完全モデルの対数尤度の差を用いるデビアンズ残差，あるいは，スチューデント化デビアンズ残差を使うのが自然の流れのように思われる．

表 11.12 SAS/GENMOD の出力：ポアソン回帰の各種の残差統計量

観測値の統計量						
オブザベーション	未加工残差	Pearson 残差	デビアンズ 残差	標準化デビアンズ残差	標準化 Pearson 残差	尤度残差
1	-0.5163	-0.3255	-0.3377	-0.4558	-0.4393	-0.4484
2	0.4837	0.3049	0.2958	0.3993	0.4115	0.4048
3	-1.4516	-0.5318	-0.5506	-0.5820	-0.5621	-0.5799
:						
9	2.6131	0.7425	0.7184	0.8167	0.8440	0.8229

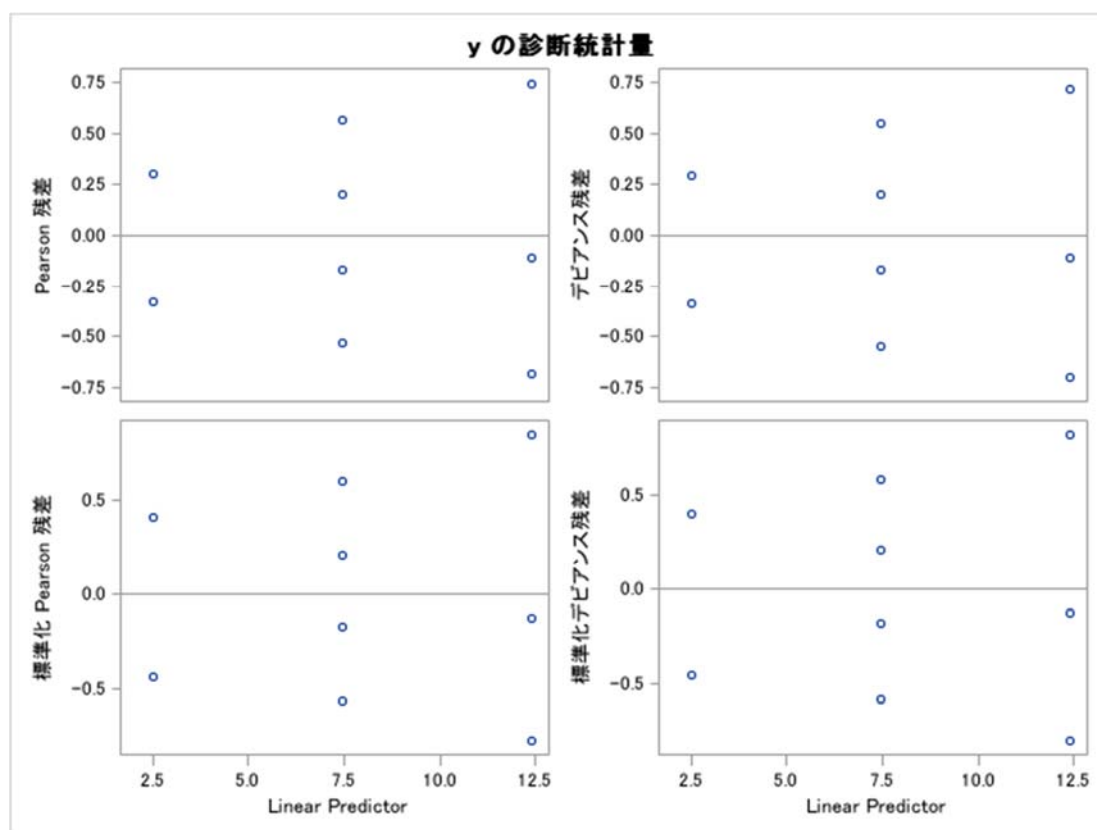


図 11.3 SAS/GENMOD による 4 種の残差プロット

## 11.5. カブトガニの事例における 4 種の残差

アグレスティ (2003) のカブトガニのデータについて第 1.13 節で概要を示し、第 7.2 節で探索的な解析の事例として用いた。ここでは、ポアソン回帰で用いられている 4 種の残差を比較検討するために用いる。このデータは、雌のカブトガニに連結する雄のサテライト数 (Satellite 数) について 173 匹について、名義尺度 (甲羅の色, 後体部の棘の状態) の 2 変数, 連続尺度 (甲羅の幅, 体重) の 2 変数, 反応変数としてサテライト数が含まれている。

### JMP による 4 種の残差の計算

第 1.13 節では、雌のカブトガニに連結する雄のサテライト数を反応変数とし、雌の甲羅の幅を説明変数とした対数リンクでのポアソン回帰を行い、散布図に回帰曲線および 95%信頼区間と予測区間 (個別データの 95%信頼区間) を示し、多くのデータが 95%予測区間の外にあることを示した。残差プロットには、JMP の過分散の調整に使われている Pearson 残差について、「予測値 vs. Pearson 残差」を例示した。

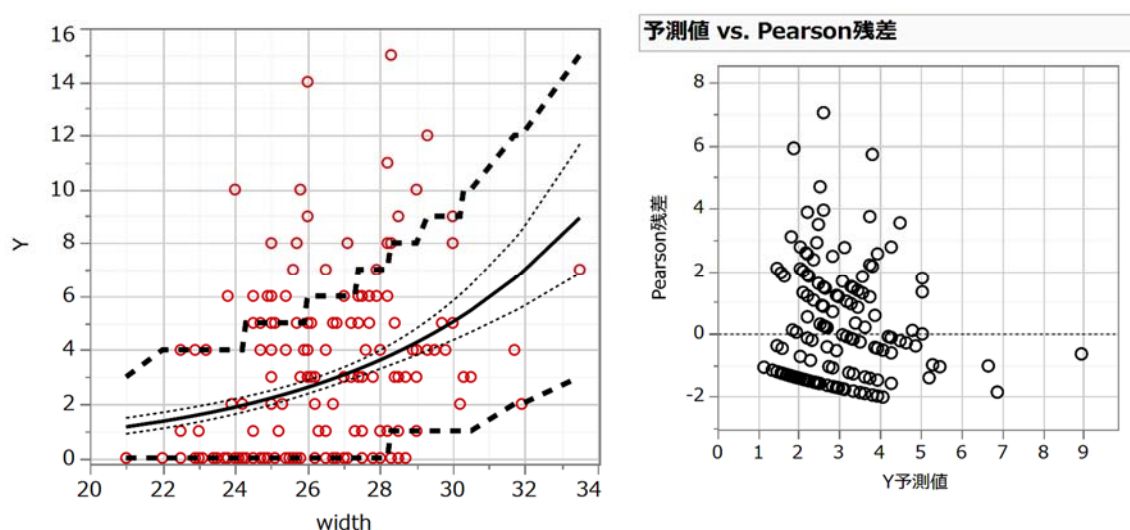


図 11.4 ポアソン回帰に対する 95%信頼区間および pearson 残差 (図 1.15 再掲)

Pearson 残差プロットは、プラス側に大きく歪んでいる。これは、推定値に対してポアソン分布がプラス側に裾を引くことによる必然的な現象であること、さらに、ゼロを含む場合にプラス側に大きく裾を引くことも影響している。このような必然的に起きるバイアスを少しでも解消するためにデビエンス残差の使用が望ましい。

JMP で過分散なしの対数リンクでのポアソン回帰分析を行い、4 種の診断プロットを選択し、さらに、4 種の残差を JMP ファイルとして出力して、相互の比較を行う。

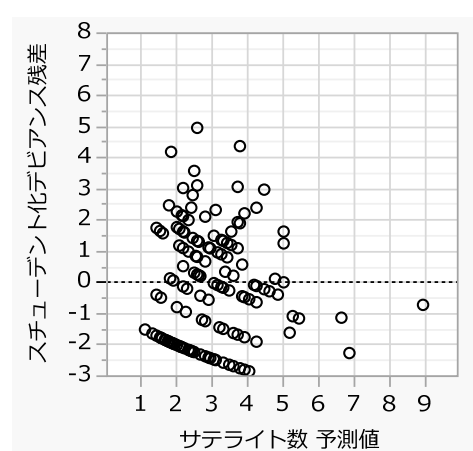
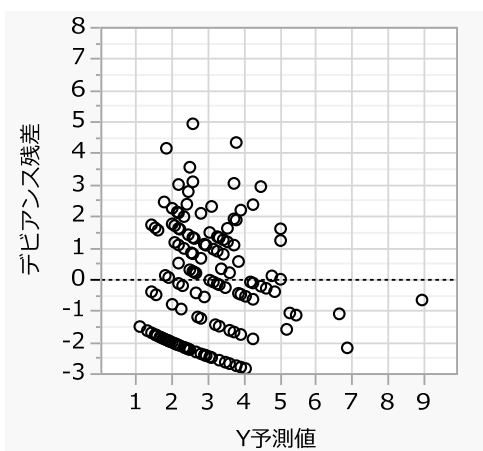
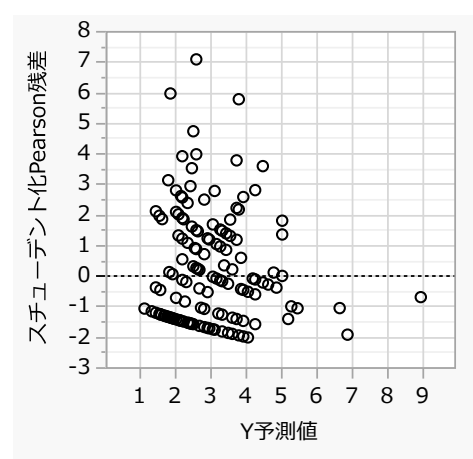
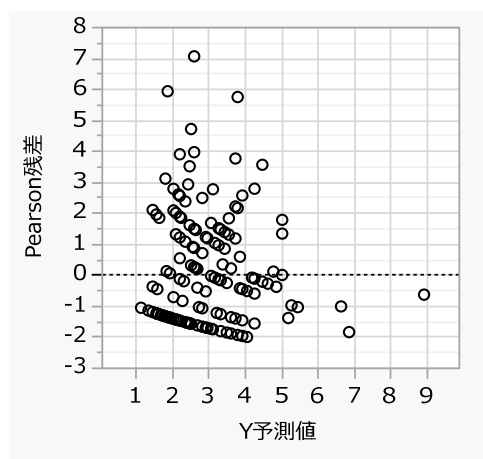
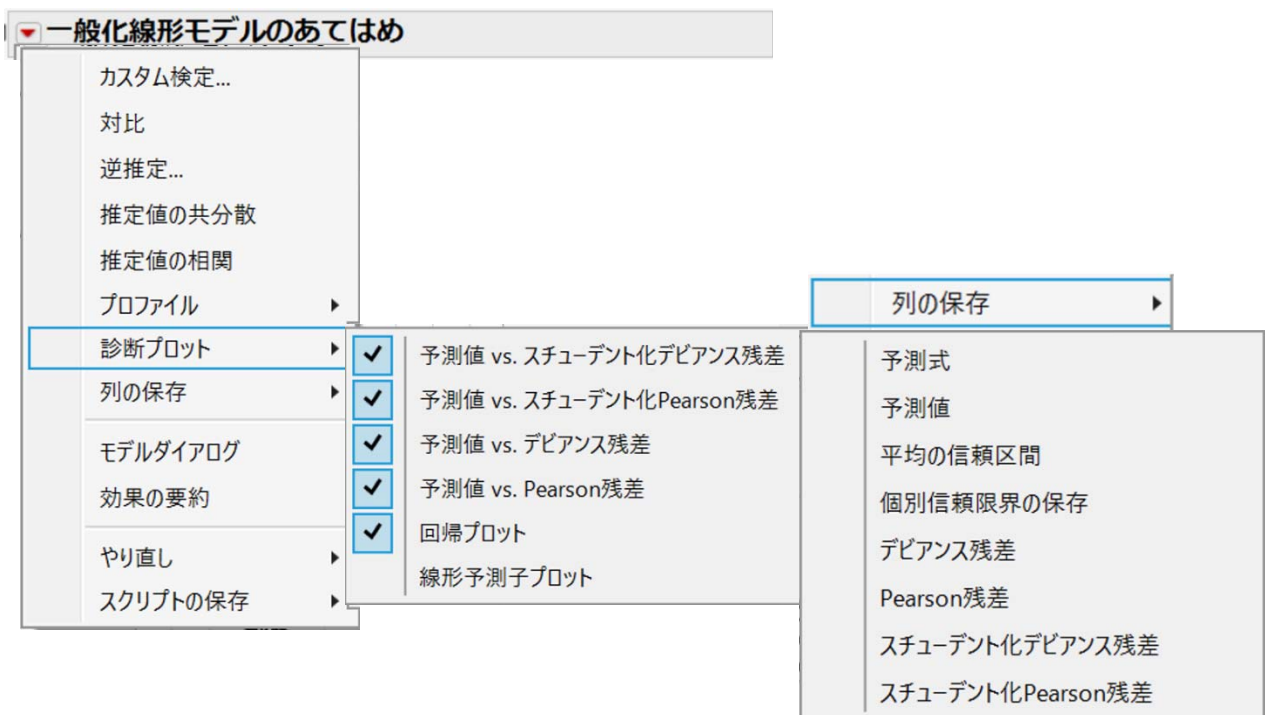


図 11.5 JMP/ポアソン回帰による 4 種の残差プロット



#### 4 種の残差の比較

(Pearson 残差 vs. デビアンズ残差) では、正の残差については、デビアンズ残差を全体的に小さい方向へシフトし、負の残差の場合は、負の大きい方向にシフトしている。(スチューデント化 Pearson 残差 vs. スチューデント化デビアンズ残差) も同様である。

(Pearson 残差 vs. スチューデント化 Pearson 残差) および (デビアンズ残差 vs. スチューデント化デビアンズ残差) については、式 (11.7) および式 (11.9) を用いたスチューデント化に用いているテコ比  $h_{ii}$  が 1 以下であることから、必然的に残差ゼロを基準に残差の絶対値を大きい方に広げることになる。言い換えれば、プラスの残差はプラス方向に引き伸ばし、マイナスの残差は、マイナス方向に引き伸ばす。

どの程度の差が実際に起きるのかを実感するために、表 11.13 に JMP ファイルに出力された 4 種の残差を Excel に取り込みデビアンズ残差の大きい順に並べて抜粋した結果を示す。スチューデント化した場合の残差の差について計算した結果は、大きいもので小数点以下 2 桁目での差であり、目立った違いではない。

Pearson 残差とデビアンズ残差を比較すると、最初の行の No.15 では、Pearson 残差が 7.0448 であるのに対し、デビアンズ残差は、4.9221 と明らかにデビアンズ残差に縮小効果が表れている。また最後の No.94 の場合は、-2.0171→-2.8526 とマイナス方向への明らかな引き延ばしを確認される。

表 11.13 ポアソン回帰の各種の残差統計量

No	甲羅の色	後体部の棘	体重	甲羅の幅	サテライト数	予測式サテライト数	Pearson 残差	スチューデント化 Pearson 残差	Pearson 残差の差	デビアンズ残差	スチューデント化デビアンズ残差	デビアンズ残差の差
15	2	1	2.30	26.0	14	2.6128	7.0448	7.0673	<b>0.0225</b>	4.9221	4.9379	<b>0.0157</b>
56	2	3	3.00	28.3	15	3.8103	5.7324	5.7608	<b>0.0284</b>	4.3279	4.3494	<b>0.0215</b>
13	2	3	3.05	28.2	11	3.7483	3.7456	3.7632	<b>0.0176</b>	3.0301	3.0443	<b>0.0142</b>
165	2	3	2.75	26.5	7	2.8361	2.4725	2.4799	<b>0.0074</b>	2.0787	2.0849	<b>0.0062</b>
28	2	1	2.70	26.8	5	2.9792	1.1708	1.1743	<b>0.0035</b>	1.0659	1.0692	<b>0.0032</b>
91	2	1	3.85	29.7	5	4.7941	0.0940	0.0951	<b>0.0011</b>	0.0934	0.0945	<b>0.0011</b>
124	2	3	1.65	24.2	2	1.9448	0.0396	0.0398	<b>0.0002</b>	0.0394	0.0396	<b>0.0002</b>
44	2	1	3.30	30.0	5	5.0359	-0.0160	-0.0162	<b>-0.0002</b>	-0.0160	-0.0162	<b>-0.0002</b>
63	3	1	2.45	27.0	3	3.0786	-0.0448	-0.0449	<b>-0.0001</b>	-0.0450	-0.0451	<b>-0.0001</b>
50	2	1	3.60	30.3	3	5.2900	-0.9956	-1.0122	<b>-0.0166</b>	-1.0848	-1.1028	<b>-0.0181</b>
81	3	2	2.25	24.5	0	2.0429	-1.4293	-1.4361	<b>-0.0068</b>	-2.0213	-2.0309	<b>-0.0096</b>
94	2	1	3.20	28.7	0	4.0688	-2.0171	-2.0297	<b>-0.0126</b>	-2.8526	-2.8705	<b>-0.0178</b>

全体で 173 サンプルをデビアンズ残差の大きい順に並べ、最大値と最小値を残し、デビアンズ残差がおおむね等間隔になるように 12 サンプルを抽出した。

図 11.6 に全 173 匹についてスチューデント化した場合の残差の変化について JMP の「対応のあるペア」によって比較した結果を示す. スチューデント化した場合, 残差がゼロを起点にプラス方向とマイナス方向に引き伸ばされていることが確認される.

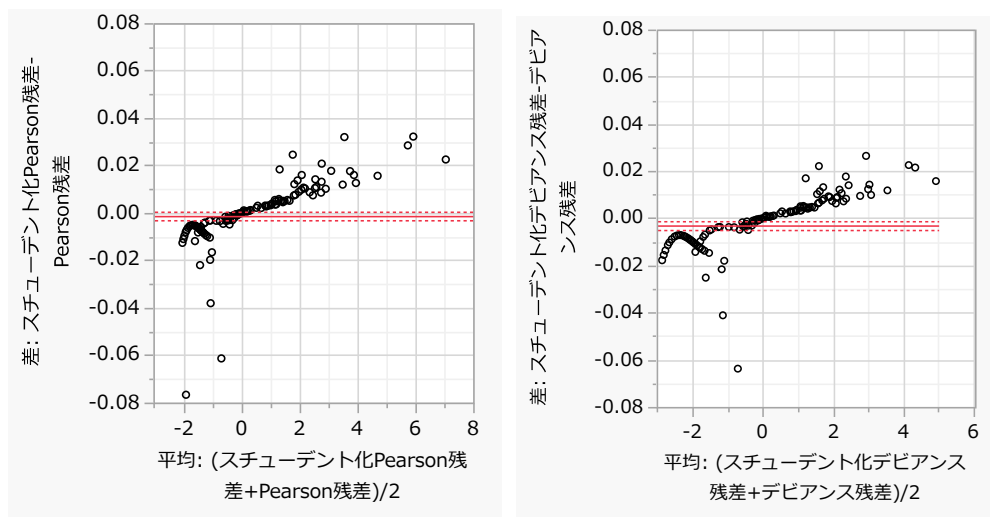


図 11.6 スチューデント化による残差の大きさの比効

図 11.7 に, スチューデント化 Pearson 残差とスチューデント化デビアンس残差の比較を JMP の「二変量の関係」によって作図した結果を示す. デビアンスにした場合に, 残差がプラスの場合は圧縮され, マイナスの場合は, 引き伸ばされていることが確認される. 対数リンクの場合には, 観測値が大いの方に裾を引くので, スチューデント化 Pearson 残差は, 大きい方に引っ張られる. ビアンス残差あるいはスチューデント化デビアンス残差による補正が行われていることが確認される. したがって, 対数リンクの場合には, ビアンス残差あるいはスチューデント化デビアンス残差の使用が望ましい.

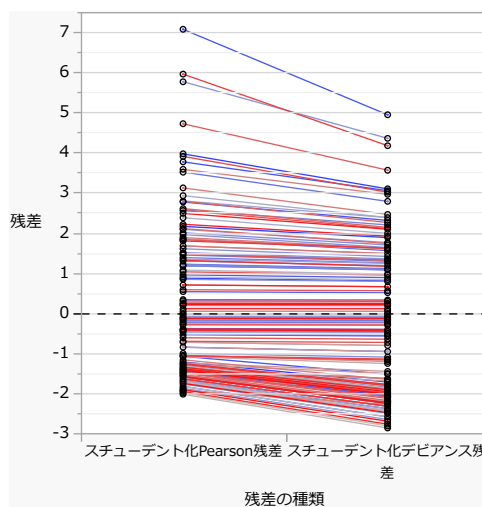


図 11.7 スチューデント化 Pearson 残差とスチューデント化デビアンス残差の比較

## 第11章 文献索引









アグレスティ著, 渡邊・菅波・吉田・角野・寒水・松永 訳(2003) - カテゴリカルデー解析入門	379
SAS Institute(2016) - SAS/STAT® 14.2 User's Guide, The GENMOD Procedure	377
ドブソン著, 田中・森川・山中・富田 訳(2008) - 一般化線形モデル入門, 原著 第2版	361
野沢昌弘(1992) - テコ比とハット行列	366

## 第11章 索引

あ 赤池の情報量基準 - AICc	371	残差プロット - JMP	375
- 修正済み	371	GENMODプロシジャ - 4種の残差	377
アグレスティ(2003) - カブトガニの事例	379	JMP - 一般化線形モデル	380
逸脱度 - デビアン	359	- 残差プロット	375
一般化線形モデル - JMP	380	修正済み - 赤池の情報量基準	371
- 診断プロット	380	縮小モデル - 切片	369
- 列の保存	380	- 総平方和	369
AICc - 赤池の情報量基準	371	種子数 - 久保(2012)	359
Excel - 標準残差	367	使用上の注意 - 標準残差	367
Excel 回帰分析 - 標準化残差	367	診断プロット - 一般化線形モデル	380
Excelの行列関数 - 回帰分析	361	- ポアソン回帰	379
重み行列 - 対角要素	374	人工データ - ドブソン(2008)	361
- ハット行列	374	スチューデント化 - 残差の比較	382
か カイ2乗値 - デビアン	370	- テコ比	372
- Pearson	370	- デビアン	359
回帰の平方和 - 誤差平方和	362	- デビアン残差	372, 381
- 差分	369	- Pearson残差	376, 381
回帰パラメータ - 平方和	361	- 標準化	377
回帰分析 - Excelの行列関数	361	スチューデント化デビアン残差 - テコ比	374
- 行列関数	361	スチューデント化残差 - 残差の分散	364
- 通常の	361	- 通常の回帰分析	364
各種の残差 - SAS/GENMOD	377	- テコ比	367
各種の残差統計量 - SAS/GENMOD	378	切片 - 縮小モデル	369
カブトガニの事例 - アグレスティ(2003)	379	総平方和 - 回帰の平方和	362
- 4種の残差の比較	379	- 縮小モデル	369
完全モデル - 誤差平方和	368	た 対角要素 - 重み行列	374
- 最大モデル	368	- 共分散行列	363
共分散行列 - 対角要素	363	- 分散	363
- パラメータ	363	対数尤度 - 3種	359
行列関数 - 回帰分析	361	- 飽和モデル	360, 368
久保(2012) - 種子数	359	対数尤度の差の2倍 - デビアン	372
恒等リンク - ポアソン回帰	368	通常の - 回帰分析	361
誤差平方和 - 完全モデル	368	通常の回帰分析 - スチューデント化残差	364
- 誤差平方和	362	通常の残差 - Pearson残差	360
個別データの95%信頼区間 - 予測区間	379	適合度統計量 - デビアン	370
さ 最大モデル - 完全モデル	368	- Pearson	370
SAS Institute(2016) - 尤度残差	377	テコ比 - スチューデント化デビアン残差	374
SAS/GENMOD - 各種の残差	377	- ハット行列	374
- 各種の残差統計量	378	- 残差	359
- ポアソン回帰	377	- スチューデント化	372
- 未加工残差	377	- スチューデント化残差	367
- 尤度残差	377	- 野沢(1992)	366
- 4種の残差プロット	378	- ハット行列H	360
差分 - 回帰の平方和	369	- ハット行列の対角要素	364
3種 - 対数尤度	359	- 分散	367
残差 - デビアン	359	- 割引係数	366
- バイアスの補正	382	Deviance Residuals - デビアン残差	372
残差の比較 - スチューデント化	382	デビアン - 逸脱度	359
残差の分散 - スチューデント化残差	364	- カイ2乗値	370

- 残差	359
- スチューデント化	359
- 対数尤度の差の2倍	372
- 適合度統計量	370
デビアンズ残差 - スチューデント化	372, 381
- Deviance Residuals	372
- Pearson残差	381
- 平方根	372
デビアンズ残差 $\epsilon_i$ - 平方和	373
ドブソン(2008) - 人工データ	361
な 野沢(1992) - テコ比	366
- ハット行列	366
は バイアスの補正 - 残差	382
ハット行列 - 重み行列	374
- テコ比	374
- 野沢(1992)	366
ハット行列H - テコ比	360
ハット行列の対角要素 - テコ比	364
パラメータ - 共分散行列	363
Pearson - カイ2乗値	370
- 適合度統計量	370
Pearson残差 - スチューデント化	376, 381
- 通常残差	360
- デビアンズ残差	381
- 標準誤差で基準化	376
標準化残差 - Excel 回帰分析	367
標準化 - スチューデント化	377
標準誤差で基準化 - Pearson残差	376
標準残差 - Excel	367
- 使用上の注意	367
分散 - 対角要素	363
- テコ比	367
分散分析表 - 平方和	362
平方根 - デビアンズ残差	372
平方和 - 回帰パラメータ	361
- デビアンズ残差 $\epsilon_i$	373
- 分散分析表	362
ポアソン回帰 - 恒等リンク	368
- SAS/GENMOD	377
- 診断プロット	379
- 4種の残差	372
飽和モデル - 対数尤度	360, 368
ま 未加工残差 - SAS/GENMOD	377
や 尤度残差 - SAS Institute (2016)	377
- SAS/GENMOD	377
予測区間 - 個別データの95%信頼区間	379
4種の残差 - GENMODプロシジャ	377
- ポアソン回帰	372
4種の残差の比較 - カプトガニの事例	379
4種の残差プロット - SAS/GENMOD	378
ら 列の保存 - 一般化線形モデル	380
わ 割引係数 - テコ比	366

## 第 11 章 Excel, JMP, SAS ファイル一覧

	4 KB	第11章02_人工データ_線形	JMP Data Table
	37 KB	第11章02_人工データ_線形	Microsoft Excel ワークシート
	57 KB	第11章03_ポアソン回帰_デビアンズ残差	Microsoft Excel ワークシート
	49 KB	第11章04_人工データ_4種の残差	Microsoft Excel ワークシート
	1 KB	第11章04_人工データ_4種の残差_SAS	テキストドキュメント
	27 KB	第11章05_カブトガニ_4種の残差	JMP Data Table
	28 KB	第11章05_カブトガニ_4種残差	Microsoft Excel ワークシート
	72 KB	第11章05_カブトガニ_転置_グラフ作成	JMP Data Table

非売品, 無断複製を禁ずる

第 9 回 続高橋セミナー

最尤法によるポアソン回帰分析入門<第 11 章>

## 第 11 章 デビアン스・逸脱度・テコ比・4 種の残差

BioStat 研究所(株)

〒105-0014 東京都 港区 芝 1-12-3 の 1005

2020 年 7 月 11 日 高橋 行雄