

第 8 回 続高橋セミナー

最尤法による探索的ポアソン回帰

2019 年 10 月

高橋 行雄 BioStat 研究所(株)

要約： 成書で取り上げられている「カブトガニのサテライト数の観察データ」について最尤法による探索的ポアソン回帰を行い、2019 年の SAS ユーザ総会で発表をしたところ、予想以上に好評であったので、公表論文およびスライドを元に第 8 回の続・高橋セミナーとする。解析に用いたのは JMP であるが、JMP で行った解析を Excel で追試した結果も追加した。探索的なデータ解析の結果を解釈する際に JMP の予測プロファイルが有用であるが、他の統計ソフトではサポートされていないので、JMP で出力されるパラメータの推定値および共分散行列を用いて、Excel による予測プロファイルの計算方法、予測プロファイルのグラフ表示法について詳細に示した。さらに、Excel を用いて最尤法によるポアソン回帰の計算方法についても詳細に示した。

目 次

1. はじめに-----	1
2. 一般化線形モデル-----	5
3. ポアソン分布のあてはめ-----	8
4. 過分散を調整したポアソン回帰-----	9
5. ポアソン回帰の個別データの 95%信頼区間-----	12
6. ガンマ・ポアソン分布のあてはめ-----	14
7. 層別解析-----	15
8. 甲羅の幅か体重か-----	18
9. 甲羅の色と体重の組み合わせ-----	21
10. 後体部の棘と体重の組み合わせ-----	23
11. 層別散布図行列における回帰の 95%信頼区間-----	25
12. EXCEL による予測プロファイル-----	29
13. 甲羅の幅に対する信頼区間と予測区間-----	39
14. 最尤法による対数リンクのポアソン回帰-----	45
15. まとめ-----	50
文 献-----	54
EXCEL, JMP ファイル 一覧-----	55

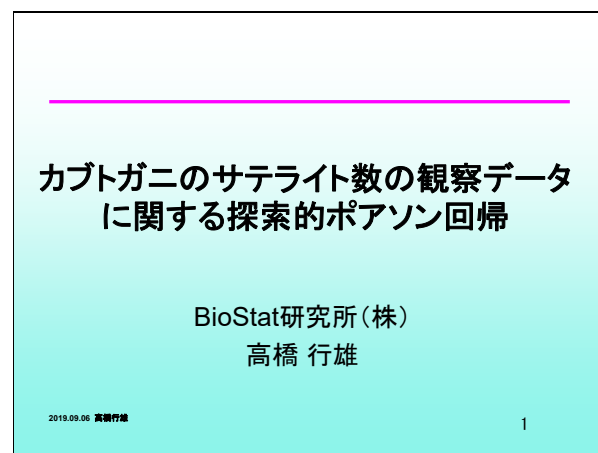
目 次 最尤法による探索的ポアソン回帰

1. はじめに	1
2. 一般化線形モデル	5
3. ポアソン分布のあてはめ	8
4. 過分散を調整したポアソン回帰	9
5. ポアソン回帰の個別データの 95%信頼区間	12
6. ガンマ・ポアソン分布のあてはめ	14
7. 層別解析	15
8. 甲羅の幅か体重か	18
9. 甲羅の色と体重の組み合わせ	21
10. 後体部の棘と体重の組み合わせ	23
11. 散布図行列における回帰の 95%信頼区間	25
12. EXCEL による予測プロファイル	29
13. 甲羅の幅に対する信頼区間と予測区間	39
14. 最尤法による対数リンクのポアソン回帰	45
15. まとめ	50
文献	54
EXCEL, JMP ファイル 一覧	55

1. はじめに

2019 年 5 月 17 日の日本計量生物学会で「ポアソン回帰を用いた勾配比検定」、2019 年 9 月 6 日の SAS ユーザ総会で「カブトガニのサテライト数の観察データに関する探索的ポアソン回帰」と題してポアソン回帰について発表した。これらの講演は、2018 年から取り組んでいる「最尤法によるポアソン回帰分析入門」の応用例として執筆してる事例を深堀するために行った。予想以上の反響があったので、探索的ポアソン回帰について Excel を用いた各種の推定法の解説も含めて第 8 回の続・高橋セミナーとする。

スライド 1



ポアソン回帰については、2004 年の第 17 回 高橋セミナーで「ポアソン回帰分析入門 ー細胞数をカウントしたデータの解析ー」として取り上げた。久保（2012）は、「データ解析のための統計モデリング入門，一般化線形モデル・階層ベイズモデル・MCMC」で、「何でもかんでも正規分布と考えるのはおかしいだろう」というコンセプトで正規分布ではなくポアソン分布を全面的に取り上げて論じている。さらに、第 3 章で植物の種子数を主体した「一般化線形モデル（GLM）ーポアソン回帰ー」を展開し、これまでの正規分布を前提とした統計解析とは異なる切り口を提示している。

回帰分析の基本は、応答変数 Y_i ，説明変数 X_i としたときに回帰直線 $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$ のあてはめである。回帰パラメータの切片 β_0 および傾き β_1 を推定するため，線形最小 2 乗法が標準的に用いられている。ただし，回帰分析を適切に行うためには，誤差 $\varepsilon_i = Y_i - \hat{Y}_i$ に対し，分散が均一の正規分布に従っていることが前提とされている。現実的には， Y_i が互いに独立でなくとも，誤差分散が不均一であっても，正規分布を仮定できなくとも，形式的に線形最小 2 乗法が適用できるので，これらの前提条件は無視されがちである。

説明変数 X_i が増大するにつれて、しばしば誤差分散が増大することが経験的に知られている。応答変数 Y_i が、0, 1, 2, ... のようなカウント・データの場合は、説明変数 X_i の増加に伴い、推定値 $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$ も増加し、誤差分散も同程度に増加することが知られている。このような場合には、誤差に正規分布を仮定する通常の回帰分析ではなく、ポアソン分布を仮定する回帰分析の適用が必要となる。誤差分布にポアソン分布を仮定するので、通常の回帰分析の解析法である最小 2 乗法ではなく、最尤法による解析を行う必要がある。

多くの成書では、最尤法の原理は述べるものの、実際の計算は統計ソフトに丸投げしてしまうので、読者にとっては実際にどのような計算原理なのかは、ブラック・ボックスとなってしまう。このことが、統計ソフト依存症の人達を増やす原因となっている。統計ソフトに丸投げした結果として、統計ソフトからの出力結果の範囲内での理解となり、問題の本質に迫るような応用力を奪われてしまいがちである。

多くの統計の入門書は、その時代の一般的な計算手段を意識しつつ構成されている。手回し計算機の時代、Fortran 言語が使える大型計算機の時代、電卓が普及した時代、パソコン上で Basic 言語が使えるようになった時代、R による統計解析がパソコンで手軽にできるようになった現代は、R の使い方が主体の入門書が増えている。

私も 10 年以上前迄は、統計ソフト SAS および JMP を主体にした統計解析を行ないつつ、SAS の行列計算のための IML、あるいは、JMP のスクリプトで行列計算を用いて、モデルの計算原理を自ら確認してきた。これらを用いて、「高橋セミナー」として成果物を 1999 年 12 月の第 1 回目から電子的に公表してきたが、諸般の事情で 2007 年 1 月の第 26 回で終了とした。その後、Excel の行列計算とソルバーを主体にした「続・高橋セミナー」を 2011 年に再開した。そこでは、Excel での計算結果を主体にし、JMP あるいは SAS で検証するというようなスタイルに変えた。

ポアソン回帰は、一般化線形モデルの枠組みで取り扱える。一般化線形モデルについては、2 値データについて第 6 回目の続・高橋セミナー「一般化線形モデルを Excel で極め活用する ―プロビット法・ロジット法・補 2 重対数法―」で取り上げたが、ポアソン回帰については扱わなかった。ポアソン回帰については、ドブソン著、田中・森川・山中・富田 訳 (2008) 「一般化線形モデル」が詳しいのであるが、計算過程が示されているのは、「恒等リンク」の場合だけで、「対数リンク」の場合、「オフセット」がある場合などについては、計算過程が示されていない。

JMP の一般化線形モデルに含まれているポアソン回帰は、ニュートン・ラフソン法による最尤法が用いられているが、どのような計算をしているのだろうか。SAS の一般化線形モデルの

ための GENMOD プロシジャは、反復重み付き回帰による最尤法を用いているが、どのような計算をしているのであろうか。統計ソフトを使えば結果が得られるのだから、理論を理解すれば実際の計算方法は知らなくていいのだろうか。知ったところで、どうすのか。私も 10 年前まではこの様な葛藤をしていた。救いの神は、Excel の行列関数とソルバーであった。これらを使いこなすことにより、一般化線形モデルの解析のみならず打ち切りがある寿命データ、定量限界があるデータ、経済統計で使われてるトービットモデルなどに対する最尤法による回帰分析も Excel でストレスなく実現できるようになり、その理論的背景をより身近に感ずるようになった。

誤差がポアソン分布に従う場合の 2 群間あるいは他群間の平均値の比較、直線あるいは指数曲線をあてはめるポアソン回帰（対数リンク）、複数のポアソン回帰直線の同時あてはめ、実験計画法で取り上げられている要因配置型の解析などに「ポアソン回帰」でも応用できる。この際に参考になるのは、最小 2 乗法の世界で SAS の一般線形モデル GLM プロシジャが、回帰分析、重回帰分析、共分散分析、繰り返し不揃いの多元配置分散分析、各種の直交表解析、などの従来は別々解析法と見なされていたものを統一的に取り扱えるようにしたことである。一般化線形モデルでのポアソン回帰でも SAS の GLM プロシジャと同様に、ありとあらゆる形式のデータ解析が行えるのであるが、その使い方については適当な教科書がない。

一般化線形モデルを使って、従来の回帰分析と対比して説明しようとしたときに、厄介な問題に直面した。これは、ポアソン回帰直線を最尤法によってあてはめた後、回帰直線の 95% 信頼区間（信頼区間）および個別データ 95% 信頼区間（予測区間）を Excel で計算し図示しようとしたときであった。通常、回帰直線の信頼区間および予測区間のための計算公式は、ほとんどの教科書で画一的に偏差平方和 S_{xx} などを用いた式が掲載されていて、ポアソン回帰では、その考え方が適用できないことにあった。

ポアソン回帰での信頼区間および予測区間の計算では、切片 $\hat{\beta}_0$ の分散 $Var(\hat{\beta}_0)$ 、傾き $\hat{\beta}_1$ の分散 $Var(\hat{\beta}_1)$ 、それらの共分散 $Cov(\hat{\beta}_0, \hat{\beta}_1)$ を用いる。通常、回帰直線の場合に、共分散を用いる方法が一般的になっていれば説明がしやすいのであるが、ほとんど見出すことができない。そもそも、通常、回帰分析の解析手順に共分散 $Cov(\hat{\beta}_0, \hat{\beta}_1)$ を用いる計算が含まれていないためである。

観察データとして得られたカウント・データに対し、ポアソン回帰による探索的な解析を試みようとしたときに（恒等リンク or 対数リンク）、（オフセットの有り or 無し）、（過分散の調整 or ガンマ・ポアソン分布 or ゼロ過剰ポアソン分布）などの選択が迫られる。同じ回帰分析なのに、なれ親しんできた正規分布を仮定した回帰分析とは、全く異次元の世界のごとくである。そこで、成書で取り上げられてい

る「カブトガニのサテライト数の観察データ(付表 A)」に対する探索的な解析方法を示すことにより、ポアソン回帰について理解の向上を図りたい。

このデータには、173 匹のカブトガニについて説明変数として順序尺度データ(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、応答変数としてサテライト数が含まれている。全体としてはゼロ過剰ガンマ・ポアソン分布のあてはめがよいが、探索的な解析での分布としては疑問が残る。また、2 つの順序尺度間の交互作用解析は、セル度数の分布が均一ではないという問題がある。さらに、2 つの連続変数に高い相関があり、変数選択の問題もある。これらの探索的な解析に伴う種々の問題に対し、JMP の「一般化線形モデル」には、各種のプロファイル機能が備わり、交互作用の検討などに威力を発揮する。また、「グラフ・ビルダー」に散布図行列中に回帰直線と 95%信頼区間を上書きできる機能があり、これまで十分とは思えなかったサテライト数におよぼす探索的な解析を行なったので結果を示す。

2. 一般化線形モデル

スライド 2 に示すように、SAS の GENMOD プロシジャにより一般化線形モデルが使えるようになった時に、最初に注目したのは、従来からある LOGISTIC プロシジャとの使い分けであった。その詳細は、高橋(2002)に示したが、ポアソン回帰については、全く扱っていなかった。その後、高橋(2004)で、細胞毒性データについて、GENMOD プロシジャを用いたポアソン回帰による勾配比検定についての検討結果を示した。

久保(2012)は、「何でもかんでも正規分布と考えるのはおかしいだろう」というコンセプトで正規分布ではなくポアソン分布を全面的に取り上げて論じている。さらに、植物の種子数を主体した「一般化線形モデル(GLM)ーポアソン回帰ー」を展開し、これまでの正規分布を前提とした統計解析とは異なる切り口を提示した。

そこで、スライド 3 に示すように「カテゴリーカルデータ解析入門」、アグレスティ著、渡邊・菅波・吉田ら訳(2003)に示されている雌のカブトガニに連結する雄のサテライト数(Satellite 数)の例を取り上げる。

スライド 2

一般化線形モデル

- ◆ SASのGENMODにより一般化線形モデルが使えるようになった時に、最初に注目したのは、従来からあるLOGISTICプロシジャとの使い分けであった。
- ◆ 最近、植物の種子数を主体した「一般化線形モデル(GLM)ーポアソン回帰ー」が注目され、正規分布を前提とした回帰分析とは異なる切り口が提唱されている。

2019.09.06 高橋行雄

2

スライド 3

カウントデータの事例

- ◆ アグレスティ著、渡邊ら訳(2003)、カテゴリーカルデータ解析入門、サイエンティスト社、110-127, 168-179.
- ◆ 雌のカブトガニに連結する雄のサテライト数(Satellite数)の例を取り上げる。
- ◆ 173匹のカブトガニ: 説明変数として甲羅の色、後体部の棘の状態、甲羅の幅、体重の4変数、応答変数としてサテライト数

2019.09.06 高橋行雄

3

このデータには、スライド 4 および(付表 A)に示すように、173 匹のカブトガニについて、説明変数として順序尺度データ(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、応答変数としてサテライト数が含まれている。

スライド 5 に示すように、文献では、甲羅の幅をX軸、サテライト数をY軸とした散布図と共に、対数リンクによるポアソン回帰の結果が示されている。しかし、その後の解析では甲羅の幅を 8 区分とし区分内のカブトガニの数とサテライト数の合計を算出し、カブトガニの数をオフセットとした解析を主体にし

雌のカブトガニに連結する雄のサテライト数

col or	spi ne	width	weight	sat ell	col or	spi ne	width	weight	sat ell	col or	spi ne	width	weight	sat ell	col or	spi ne	width	weight	sat ell
2	3	28.3	3.050	8	3	1	28.5	3.250	9	4	3	23.5	1.900	0	2	1	28.0	2.900	4
3	3	22.5	1.550	0	3	3	28.9	2.800	4	2	2	24.0	1.700	0	4	3	25.8	2.250	10
1	1	26.0	2.300	9	2	3	28.2	2.600	6	2	1	29.7	3.850	5	2	3	27.9	3.050	7
3	3	24.8	2.100	0	2	3	25.0	2.100	4	2	1	26.8	2.550	0	2	3	24.9	2.200	0
3	3	26.0	2.600	4	2	3	28.5	3.000	3	4	3	26.7	2.450	0	2	1	28.4	3.100	5
2	3	27.9	2.800	6	3	1	25.9	2.550	4	4	3	22.5	1.475	4	4	3	27.0	2.625	0
3	3	27.5	3.100	3	2	3	25.8	2.300	0	2	3	26.2	2.025	2	2	2	24.5	2.000	0
1	1	26.1	2.800	5	4	3	27.0	2.250	3	2	1	24.9	2.300	6					
1	1	27.7	2.500	6	2	3	28.5	3.050	0	1	2	24.5	1.950	6					
2	1	30.0	3.300	5	4	1	25.5	2.750	0	2	3	25.1	1.800	0					

注釈: color=色(1=やや明るい, 2=中くらい, 3=やや暗い, 4=暗い);

Spine=後体部の棘の状態(1=いずれも正常, 2=一方が摩耗または破損している, 3=いずれも摩耗または破損している);

width=甲羅の幅(cm); weight=重さ(kg); satell=サテライト数.

出典: <http://lib.stat.cmu.edu/datasets/agresti>. 2019年7月24日 アクセス.

いる. 探索的な解析では, サテライト数が(0, 1 以上)の 2 値データとして, ロジスティック回帰を主体にした解析方法が提示されている.

文献での探索的な解析

- ◆ 甲羅の幅をX軸, サテライト数をY軸とした散布図と共に, 対数リンクによるポアソン回帰
 - その後の解析では甲羅の幅を8区分とし区分内のカブトガニの数とサテライト数の合計を算出し, カブトガニの数をオフセットとした解析を主体
 - 探索的な解析では, サテライト数が(0, 1 以上)の2値データとして, ロジスティック回帰を主体にした解析方法が提示

付表 A 雌のカブトガニに連結する雄のサテライト数

col or	spi ne	width	weight	sat ell	col or	spi ne	width	weight	sat ell	col or	spi ne	width	weight	sat ell	col or	spi ne	width	weight	sat ell
2	3	28.3	3.050	8	3	1	28.5	3.250	9	4	3	23.5	1.900	0	2	1	28.0	2.900	4
3	3	22.5	1.550	0	3	3	28.9	2.800	4	2	2	24.0	1.700	0	4	3	25.8	2.250	10
1	1	26.0	2.300	9	2	3	28.2	2.600	6	2	1	29.7	3.850	5	2	3	27.9	3.050	7
3	3	24.8	2.100	0	2	3	25.0	2.100	4	2	1	26.8	2.550	0	2	3	24.9	2.200	0
3	3	26.0	2.600	4	2	3	28.5	3.000	3	4	3	26.7	2.450	0	2	1	28.4	3.100	5
2	3	23.8	2.100	0	2	1	30.3	3.600	3	2	1	28.7	3.200	0	3	3	27.2	2.400	5
1	1	26.5	2.350	0	4	3	24.7	2.100	5	3	3	23.1	1.550	0	2	2	25.0	2.250	6
3	2	24.7	1.900	0	2	3	27.7	2.900	5	2	1	29.0	2.800	1	2	3	27.5	2.625	6
2	1	23.7	1.950	0	1	1	27.4	2.700	6	3	3	25.5	2.250	0	2	1	33.5	5.200	7
3	3	25.6	2.150	0	2	3	22.9	1.600	4	3	3	26.5	1.967	1	2	3	30.5	3.325	3
3	3	24.3	2.150	0	2	1	25.7	2.000	5	3	3	24.5	2.200	1	3	3	29.0	2.925	3
2	3	25.8	2.650	0	2	3	28.3	3.000	15	3	3	28.5	3.000	1	2	1	24.3	2.000	0
2	3	28.2	3.050	11	2	3	27.2	2.700	3	2	3	28.2	2.867	1	2	3	25.8	2.400	0
4	2	21.0	1.850	0	3	3	26.2	2.300	3	2	3	24.5	1.600	1	4	3	25.0	2.100	8
2	1	26.0	2.300	14	2	1	27.8	2.750	0	2	3	27.5	2.550	1	2	1	31.7	3.725	4
1	1	27.1	2.950	8	4	3	25.5	2.250	0	2	2	24.7	2.550	4	2	3	29.5	3.025	4
2	3	25.2	2.000	1	3	3	27.1	2.550	0	2	1	25.2	2.000	1	3	3	24.0	1.900	10
2	3	29.0	3.000	1	3	3	24.5	2.050	5	3	3	27.3	2.900	1	2	3	30.0	3.000	9
4	3	24.7	2.200	0	3	1	27.0	2.450	3	2	3	26.3	2.400	1	2	3	27.6	2.850	4
2	3	27.4	2.700	5	2	3	26.0	2.150	5	2	3	29.0	3.100	1	2	3	26.2	2.300	0
2	2	23.2	1.950	4	2	3	28.0	2.800	1	2	3	25.3	1.900	2	2	1	23.1	2.000	0
1	2	25.0	2.300	3	2	3	30.0	3.050	8	2	3	26.5	2.300	4	2	1	22.9	1.600	0
2	1	22.5	1.600	1	2	3	29.0	3.200	10	2	3	27.8	3.250	3	4	3	24.5	1.900	0
3	3	26.7	2.600	2	2	3	26.2	2.400	0	2	3	27.0	2.500	6	2	3	24.7	1.950	4
4	3	25.8	2.000	3	2	1	26.5	1.300	0	3	3	25.7	2.100	0	2	3	28.3	3.200	0
4	3	26.2	1.300	0	2	3	26.2	2.400	3	2	3	25.0	2.100	2	2	3	23.9	1.850	2
2	3	28.7	3.150	3	3	3	25.6	2.800	7	2	3	31.9	3.325	2	3	3	23.8	1.800	0
2	1	26.8	2.700	5	3	3	23.0	1.650	1	4	3	23.7	1.800	0	3	2	29.8	3.500	4
4	3	27.5	2.600	0	3	3	23.0	1.800	0	4	3	29.3	3.225	12	2	3	26.5	2.350	4
2	3	24.9	2.100	0	2	3	25.4	2.250	6	3	3	22.0	1.400	0	2	3	26.0	2.275	3
1	1	29.3	3.200	4	3	3	24.2	1.900	0	2	3	25.0	2.400	5	2	3	28.2	3.050	8
1	3	25.8	2.600	0	2	2	22.9	1.600	0	3	3	27.0	2.500	6	4	3	25.7	2.150	0
2	2	25.7	2.000	0	3	2	26.0	2.200	3	3	3	23.8	1.800	6	2	3	26.5	2.750	7
2	1	25.7	2.000	8	2	3	25.4	2.250	4	1	1	30.2	3.275	2	2	3	25.8	2.200	0
2	1	26.7	2.700	5	3	3	25.7	1.200	0	3	3	26.2	2.225	0	3	3	24.1	1.800	0
4	3	23.7	1.850	0	2	3	25.1	2.100	5	2	3	24.2	1.650	2	3	3	26.2	2.175	2
2	3	26.8	2.650	0	3	2	24.5	2.250	0	2	3	27.4	2.900	3	3	3	26.1	2.750	3
2	3	27.5	3.150	6	4	3	27.5	2.900	0	2	2	25.4	2.300	0	3	3	29.0	3.275	4
4	3	23.4	1.900	0	3	3	23.1	1.650	0	3	3	28.4	3.200	3	1	1	28.0	2.625	0
2	3	27.9	2.800	6	3	1	25.9	2.550	4	4	3	22.5	1.475	4	4	3	27.0	2.625	0
3	3	27.5	3.100	3	2	3	25.8	2.300	0	2	3	26.2	2.025	2	2	2	24.5	2.000	0
1	1	26.1	2.800	5	4	3	27.0	2.250	3	2	1	24.9	2.300	6					
1	1	27.7	2.500	6	2	3	28.5	3.050	0	1	2	24.5	1.950	6					
2	1	30.0	3.300	5	4	1	25.5	2.750	0	2	3	25.1	1.800	0					

注釈: color=色 (1=やや明るい, 2=中くらい, 3=やや暗い, 4=暗い);

Spine=後体部の棘の状態 (1=いずれも正常, 2=一方が摩耗または破損している, 3=いずれも摩耗または破損している);

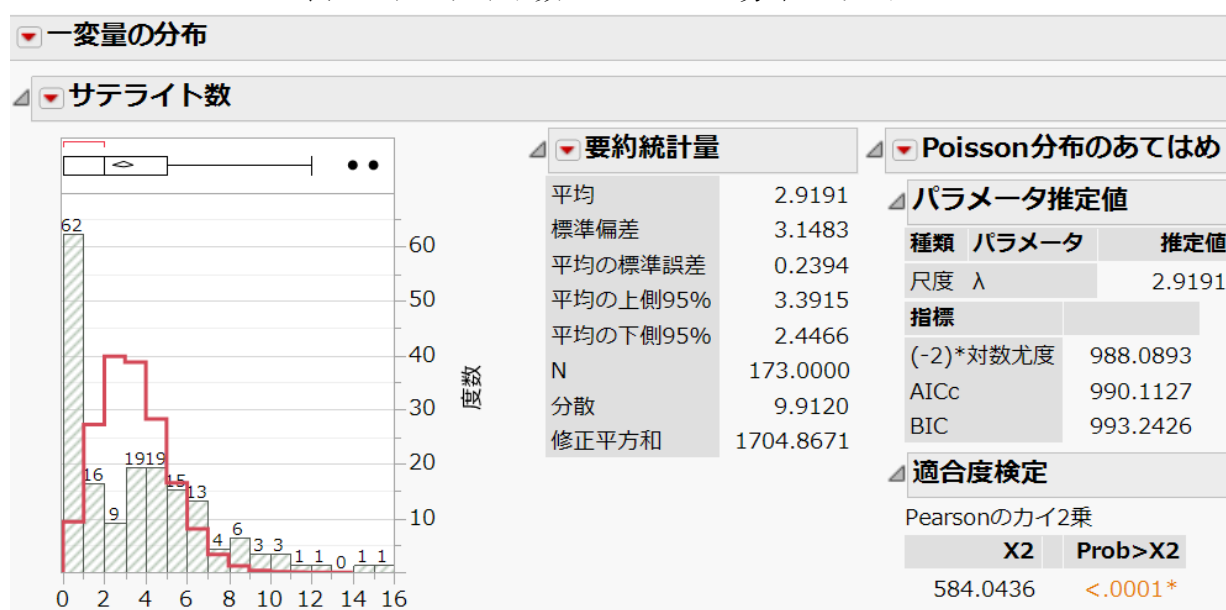
width=甲羅の幅(cm); weight =重さ(kg); satell=サテライト数.

出典: <http://lib.stat.cmu.edu/datasets/agresti>. 2019年7月24日 アクセス.

3. ポアソン分布のあてはめ

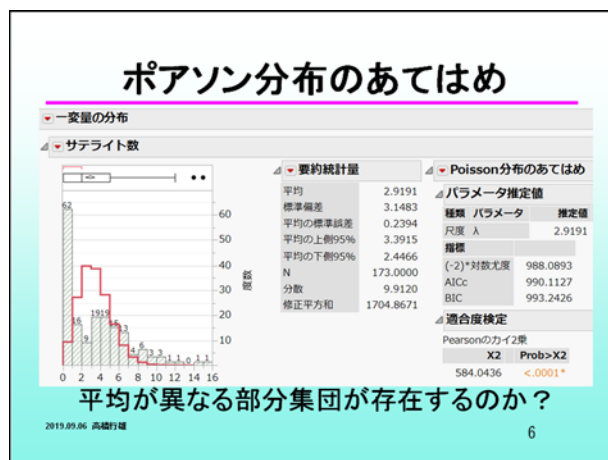
スライド 6 および表 1 に示すように, JMP の「一変量の分布」によりサテライト数の平均は 2.9191, 分散は 9.9120 であり, その比は 3.40 と過分散になっている. ポアソン分布をあてはめ, 棒グラフ上に上書きした結果を見ても, 誤差分布にポアソン分布を仮定することは絶望的とも思われる. もちろん, 適合度の検定でも $\chi^2 = 584.0436$, $p < 0.0001$ でポアソン分布があてはまるとは言えない. このような全データで過分散となる場合では, 何らかの条件によりサテライト数の平均が大きく異なるポアソン分布に従う部分集団の集まりが複数存在する可能性も考えられる.

表 1 サテライト数へのポアソン分布のあてはめ



スライド 6 JMP 「一変量の分布→離散分布のあてはめ→Poisson」

離散分布のあてはめ	Poisson
保存	ガンマPoisson
削除	二項
	ベータ二項



4. 過分散を調整したポアソン回帰

甲羅の幅を説明変数とし、サテライト数を応答変数とした場合に、次式のと対数リンク(両辺に対数を取った時に線形となる)によるポアソン回帰

$$\text{Satellite}_i = \exp(\beta_0 + \beta_1 \cdot \text{width}_i) + \varepsilon_i, \quad \varepsilon_i \sim \text{Poisson}$$

の結果を表2に示す。Pearsonの適合度のカイ2乗値は544.1570と自由度の171に対して3.1822倍と過分散となっている。

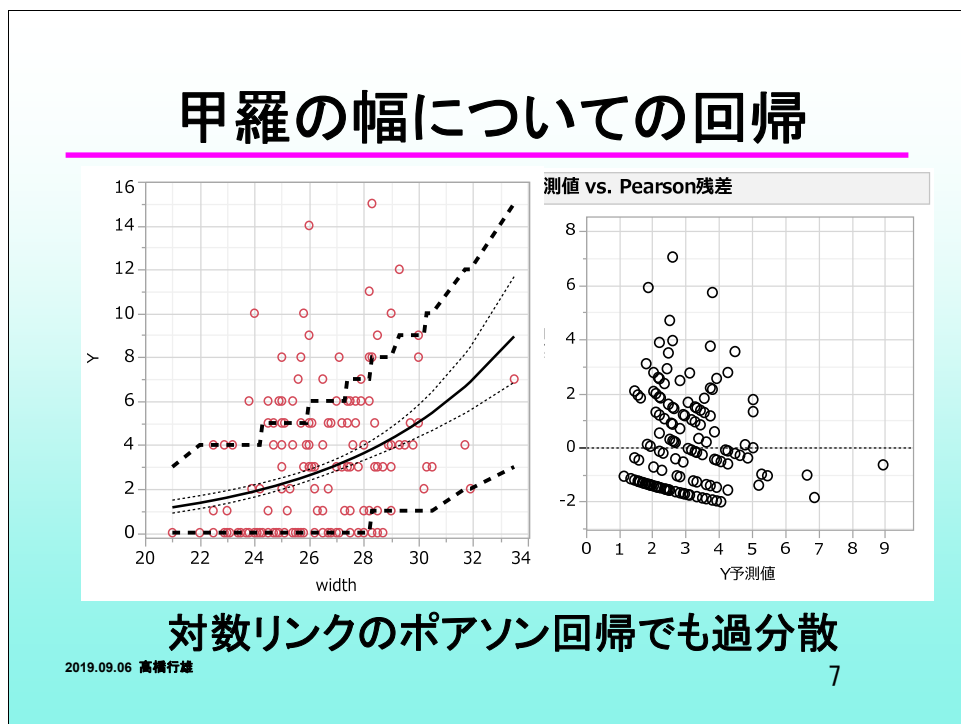
表2 甲羅の幅についての対数リンクの過分散調整なしのポアソン回帰

適合度統計量		カイ2乗	自由度	p値(Prob>ChiSq)
Pearson		544.1570	171	<.0001*
デビアン		567.8786	171	<.0001*

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	-3.3048	0.5422	36.8670	<.0001*
甲羅の幅	0.1640	0.0200	64.9131	<.0001*

スライド7に対数リンクのポアソン回帰から得られた指数曲線、推定値の95%信頼区間、個別データの95%信頼区間、予測値に対するPearson残差で3を超える点が多数あり、これらの結果からポアソン回帰による指数曲線のあてはめは支持されない。

スライド 7



回帰曲線は、 $y = \exp(-3.3048 + 0.1640x)$
 $x = 30 : y = \exp(-3.3048 + 0.1640 \times 30) = 5.0289$

図 1 に示すように、ポアソン回帰から得られた尤度比カイ 2 乗値を自由度で除した過分散パラメータを $\phi = 3.1822$ とし、得られた共分散行列に ϕ 倍して標準誤差を調整する方法が知られていて、JMP のポアソン回帰でもサポートされている。

手法: 一般化線形モデル

分布: Poisson

リンク関数: 対数

☒ 過分散に基づく検定と信頼区間

☐ Firthバイアス調整推定値

図 1 過分散の調整法オプション

過分散を調整したポアソン回帰の結果を表 3 に示す。表 2 に示した甲羅の幅の標準誤差は、 $SE = 0.0200$ であったので、調整後の SE' は、

$$SE' = \sqrt{\phi \cdot SE^2} = \sqrt{3.1822 \times 0.0200^2} = 0.0356$$

と大きくなり、尤度比カイ 2 乗値は、64.9131 から 20.3988 と激減する。

表 3 過分散調整済みのポアソン回帰

項	推定値	標準誤差	尤度比カイ2乗	p値
切片	-3.3048	0.9673	11.5854	0.0007*
甲羅の幅	0.1640	0.0356	20.3988	<.0001*

スライド 8 に過分散の調整パラメータの算出法、スライド 9 に調整前後の尤度比カイ 2 乗統計量と p 値を示す。

スライド 8

過分散の調整

- ◆ 適合度のピアソンのカイ2乗統計量を自由度で除した調整パラメータ ϕ を使う。

適合度統計量	カイ2乗	自由度	p値(Prob>ChiSq)
Pearson	544.1570	171	<.0001*
デビアン	567.8786	171	<.0001*

$\phi = 544.1570 / 171 = 3.1822$

- ◆ ポアソン回帰で得られた分散を ϕ 倍する。

2019.09.06 実験行進 8

スライド 9

過分散パラメータの適用

調整前

項	推定値	標準誤差	尤度比カイ2乗	p値
切片	-3.3048	0.5422	36.8670	<.0001*
甲羅の幅	0.1640	0.0200	64.9131	<.0001*

調整後

項	推定値	標準誤差	尤度比カイ2乗	p値
切片	-3.3048	0.9673	11.5854	0.0007*
甲羅の幅	0.1640	0.0356	20.3988	<.0001*

$SE' = \sqrt{\phi SE^2} = \sqrt{3.1822 \times 0.0200^2} = 0.0356$

2019.09.06 実験行進 9

過分散の係数を用いた方法は、過分散となるカウント・データに対する万能の方法とも思われるかもしれないが、表 1 に示したヒストグラムに重ね書きしたポアソン分布から、このデータにポアソン分布を仮定することは全くできない。もちろん甲羅の幅に対するポアソン回帰で過分散が解消するのであれば嬉しいのであるが、実際にどのような分布になるのか示すことができない。単に SE を割り増ししているだけである。

なお、スライド 28 には、ガンマ・ポアソン回帰による指数曲線、スライド 29 には、ゼロ過剰ガンマ・ポアソン回帰による指数曲線を示す。

5. ポアソン回帰の個別データの 95%信頼区間

ポアソン回帰を行っても過分散が解消していないことを視覚化するために散布図に個別データの 95%信頼区間(予測区間)を重ね書きしてみると、図 2 左に示すように上側に多数の点がはみ出ているのでポアソン回帰(指数曲線)のあてはめには無理があることを実感できる。図 2 右に示すように予測値に対する Pearson 残差をプロットすることにより、Pearson 残差が 3 以上の飛び離れデータが多数存在することからも、ポアソン分布を誤差分布とする回帰分析について否定的な結果となっている。

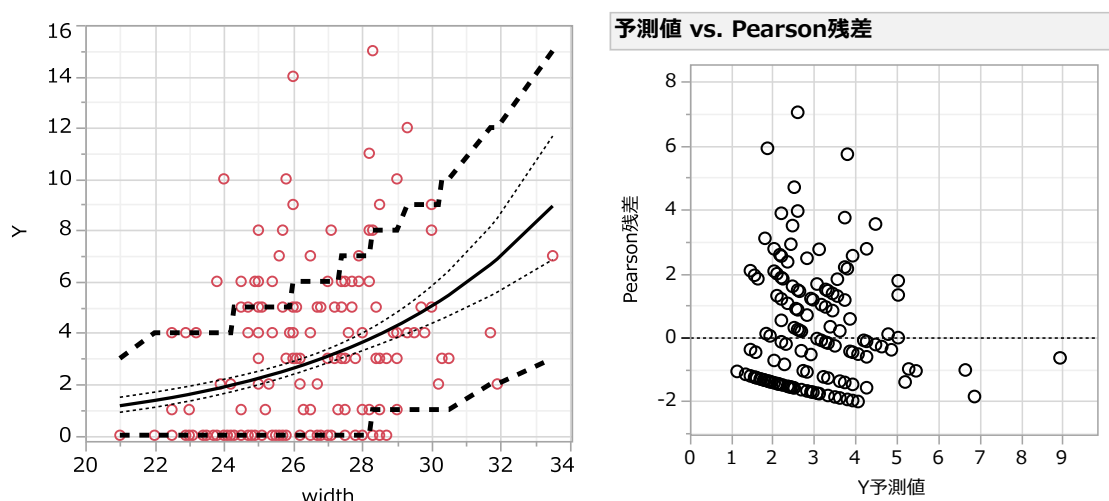


図 2 ポアソン回帰に対する 95%信頼区間および予測値に対する Pearson 残差

他の変数を加えてポアソン回帰を行っても過分散が解消されないのであれば、ポアソン回帰を行う前提がないことになる。主な原因は、173 個体に対してサテライト数ゼロが 62 匹と全体の 35.8%なので、サテライト数を (0, 1) 反応とする解析が望ましいとも考えられる。また、表 1 からサテライト数が 3 と 4 あたりに分布の山があることから、あるいは、3 区分程度の順序データとする解析を行うことが望ましいかも知れない。

ポアソン回帰の予測値の 95%信頼区間および個別データの 95%信頼区間は、JMP で「モデルのあてはめ→一般化線形モデル→分布:Poisson→リンク関数:対数→Y:サテライト数→X:甲羅の幅→実行→列の保存→(予測式, 平均の信頼区間, 個別信頼区間の保存)」のように実行し、データファイルに計算結果を出力し、それらを用いて「重ね合わせプロット→X:甲羅の幅→Y(サテライト数, 平均の下側 95%, 平均の上側 95%, 個別の下側 95%, 個別の上側 95%)→OK」で作図し、編集機能で整形する。

第 13 節では、Excel による対数リンクのポアソン回帰の 95%信頼区間の計算方法および作図法について詳しく解説する。

「モデルのあてはめ→一般化線形モデル→分布:Poisson→リンク関数:対数→Y:サテライト数→X:甲羅の幅→実行→列の保存→(予測式, 平均の信頼区間, 個別信頼区間の保存)」

第1章13_agresti_カブトガニ_探索 - JMP

ファイル(F) 編集(E) テーブル(T) 行(R) 列(C) 実験計画(DOE)(D) 分析(A) グラフ(G) ツール(O) 表示(V) ウィンドウ(W) ヘルプ(H)

第1章13_agres...
widthに...プロット
一般化線形モデル

列(10/0)
color
spine
width
weight
satellite
予測式 satellite
平均の下側95%
平均の上側95%
個別の下側95%
個別の上側95%

行
すべての行 173
選択されている行 0
除外されている行 0
表示しない行 0
ラベルのついた行 0

重ね合わせプロット - JMP

Xが連続に変化する時のYのプロット

列の選択
10列
color
spine
width
weight
satellite
予測式 satellite
平均の下側95%
平均の上側95%
個別の下側95%
個別の上側95%

オプション
☒ X値で並べ替え
☐ X軸を対数にする
☐ 左Y軸を対数にする
☐ 右Y軸を対数にする

選択した列に役割を割り当てる

Y
satellite
予測式 satellite
平均の下側95%
平均の上側95%
個別の下側95%
個別の上側95%

X
width

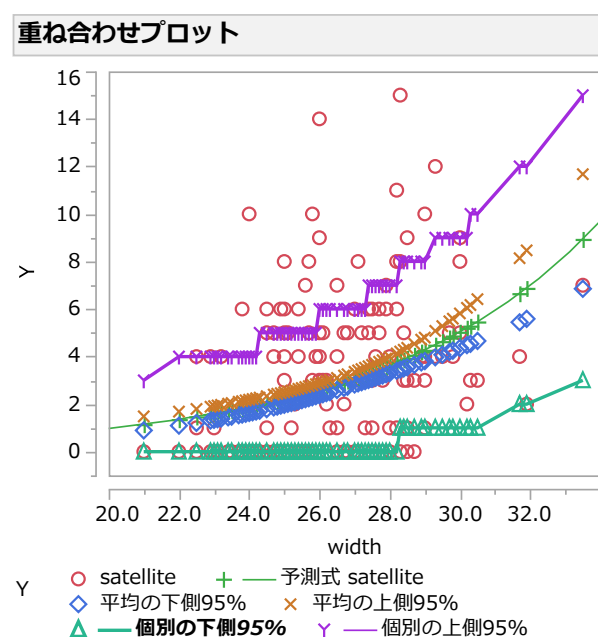
グループ変数
オプション

By
オプション

アクション
OK
キャンセル
削除
前回の設定
ヘルプ

評価が完了しました。

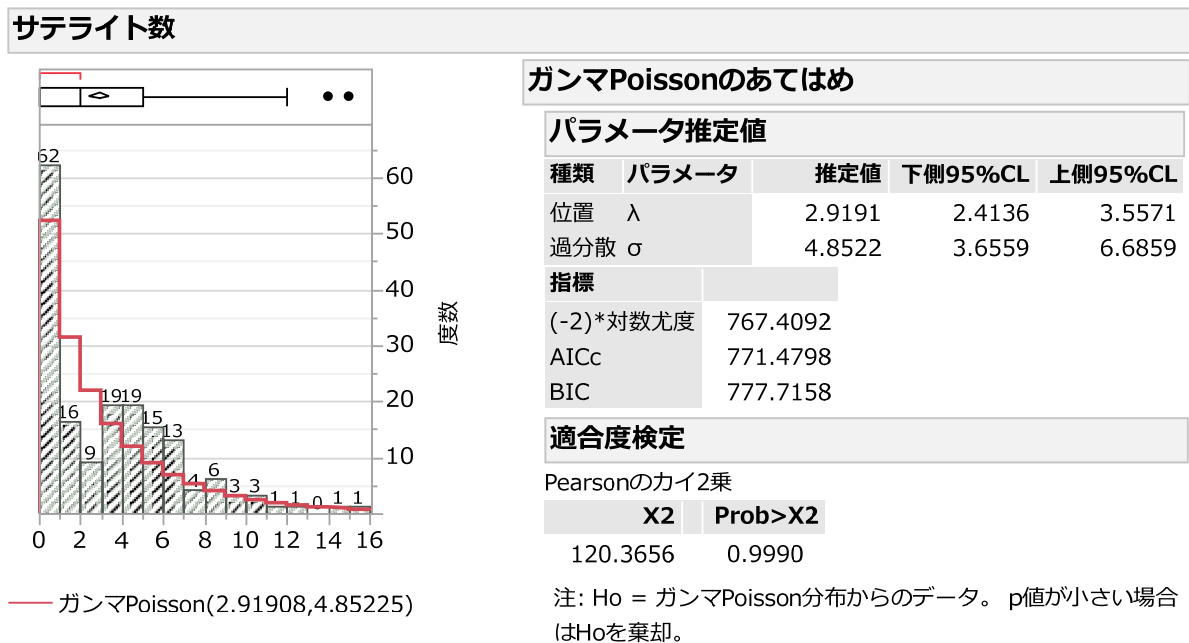
「重ね合わせプロット→X:甲羅の幅→Y(サテライト数, 平均の下側 95%, 平均の上側 95%, 個別の下側 95%, 個別の上側 95%)→OK」



6. ガンマ・ポアソン分布のあてはめ

JMP には、過分散を考慮した負の二項分布から導出されたガンマ・ポアソン分布をあてはめる機能があるので、表 4 およびスライド 10 に結果を示す。結果は、位置 $\lambda = 2.9191$ 、過分散 $\sigma = 4.8522$ となる。表 1 に示したポアソン分布のあてはめでは、サテライト数がゼロの場合について大きな乖離があったが、過分散を考慮したガンマ・ポアソン分布では、まずまずのあてはめが行われているように思われる。

表 4 サテライト数へのガンマ・ポアソン分布のあてはめ



スライド 10

ガンマ・ポアソン分布のあてはめ

サテライト数

— ガンマPoisson(2.91908,4.85225)

ガンマPoissonのあてはめ				
パラメータ推定値				
種類	パラメータ	推定値	下側95%CL	上側95%CL
位置	λ	2.9191	2.4136	3.5571
過分散	σ	4.8522	3.6559	6.6859
指標				
(-2)*対数尤度		767.4092		
AICc		771.4798		
BIC		777.7158		
適合度検定				
Pearsonのカイ2乗				
	X2	Prob>X2		
	120.3656	0.9990		

注: Ho = ガンマPoisson分布からのデータ。p値が小さい場合はHoを棄却。

あてはまりは棄却できない. OKなのか?

2019.09.06 高橋行雄 10

ガンマ・ポアソン分布は、負の 2 項分布から導出される。

7. 層別解析

スライド 4 および付表 A に示したデータには、説明変数として順序尺度(甲羅の色, 後体部の棘の状態)の 2 変数があるので JMP の「二変量の関係」で作成したサテライト数に対する層別分布を図 3 およびスライド 11 および 12 に示す。

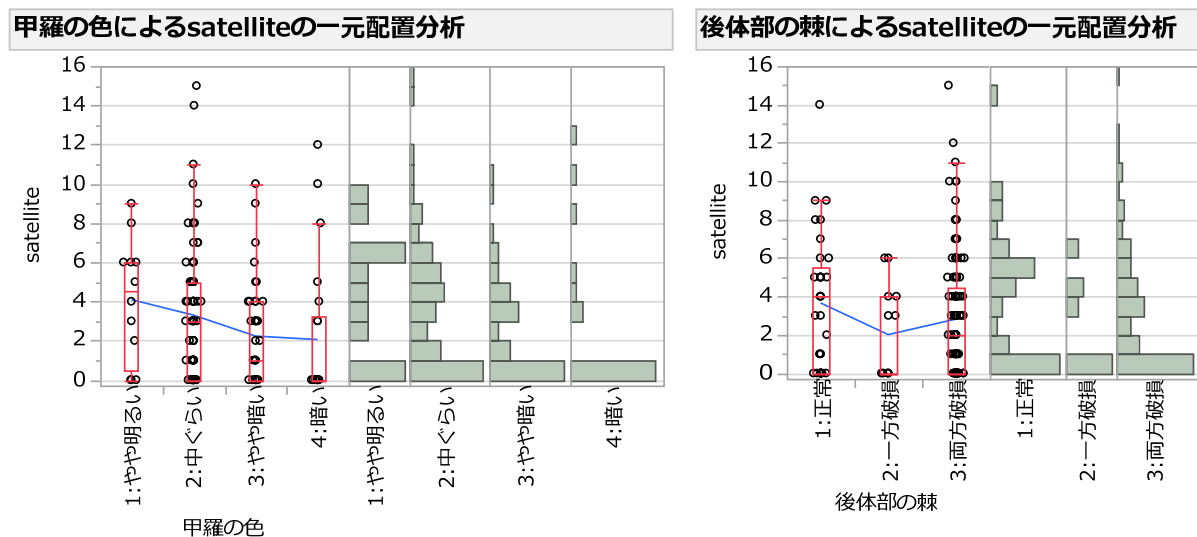


図 3 甲羅の色および後体部の棘の状態とサテライト数の関連

スライド 11

スライド 12

過分散の解消のための層別解析




◆ 甲羅の色が暗くなるとサテライト数が減少

11

層別解析の結果

- ◆ 雌の甲羅の色
 - 暗くなるに従いゼロ・カウントが増える
 - " サテライト数の平均値が減少傾向
- ◆ 雌の後体部の棘の状態
 - 正常の場合には、サテライト数の5匹に山が
 - サテライト数の平均値は同程度
- ◆ 過分散の解消とはならない

12

JMP「二変量の関係→X: (甲羅の色, 後体部の棘)→Y: サテライト数→標示オプション(箱ひげ図, 点をずらす, 平均をつなぐ, ヒストグラム)」

雌の甲羅の色については、暗くなるに従いゼロ・カウントが増えサテライト数の平均値が減少傾向であることが読み取れる。雌の後体部の棘の状態については、正常の場合には、サテライト数の 5 匹に山があり、雄が連結する割合が多いようであるが、サテライト数の平均値は同程度である。

表 5 およびスライド 13 に甲羅の色と後体部の棘の状態を組み合わせた場合のサテライト数 N, サテライト数の平均および分散を示す. 甲羅の色が暗くなるにつれて後体部の棘は, 正常から破損へ移行するが, ある程度のサテライト数がある場合の分散/平均の比は, 2 以上あり過分散が解消する様子はない.

表 5 甲羅の色別 後体部の棘別 のサテライト数の分散/平均の比

甲羅の色	棘の状態	N	平均	分散	分散/平均
1:やや明るい	1:正常	9	4.44	10.53	2.37
	2:一方破損	2	4.50	4.50	1.00
	3:両方破損	1	0.00	-	-
2:中ぐらい	1:正常	24	3.29	12.13	3.68
	2:一方破損	8	1.75	6.21	3.55
	3:両方破損	63	3.49	10.03	2.87
3:やや暗い	1:正常	3	5.33	10.33	1.94
	2:一方破損	4	1.75	4.25	2.43
	3:両方破損	37	2.03	6.25	3.08
4:暗い	1:正常	1	0.00	-	-
	2:一方破損	1	0.00	-	-
	3:両方破損	20	2.25	13.99	6.22
	全体	173	2.92	9.91	3.40

スライド 13

2変数の組合せ

甲羅の色	棘の状態	N	平均	分散	分散/平均
1:やや明るい	1:正常	9	4.44	10.53	2.37
	2:一方破損	2	4.50	4.50	1.00
	3:両方破損	1	0.00	-	-
2:中ぐらい	1:正常	24	3.29	12.13	3.68
	2:一方破損	8	1.75	6.21	3.55
	3:両方破損	63	3.49	10.03	2.87
3:やや暗い	1:正常	3	5.33	10.33	1.94
	2:一方破損	4	1.75	4.25	2.43
	3:両方破損	37	2.03	6.25	3.08
4:暗い	1:正常	1	0.00	-	-
	2:一方破損	1	0.00	-	-
	3:両方破損	20	2.25	13.99	6.22
	全体	173	2.92	9.91	3.40

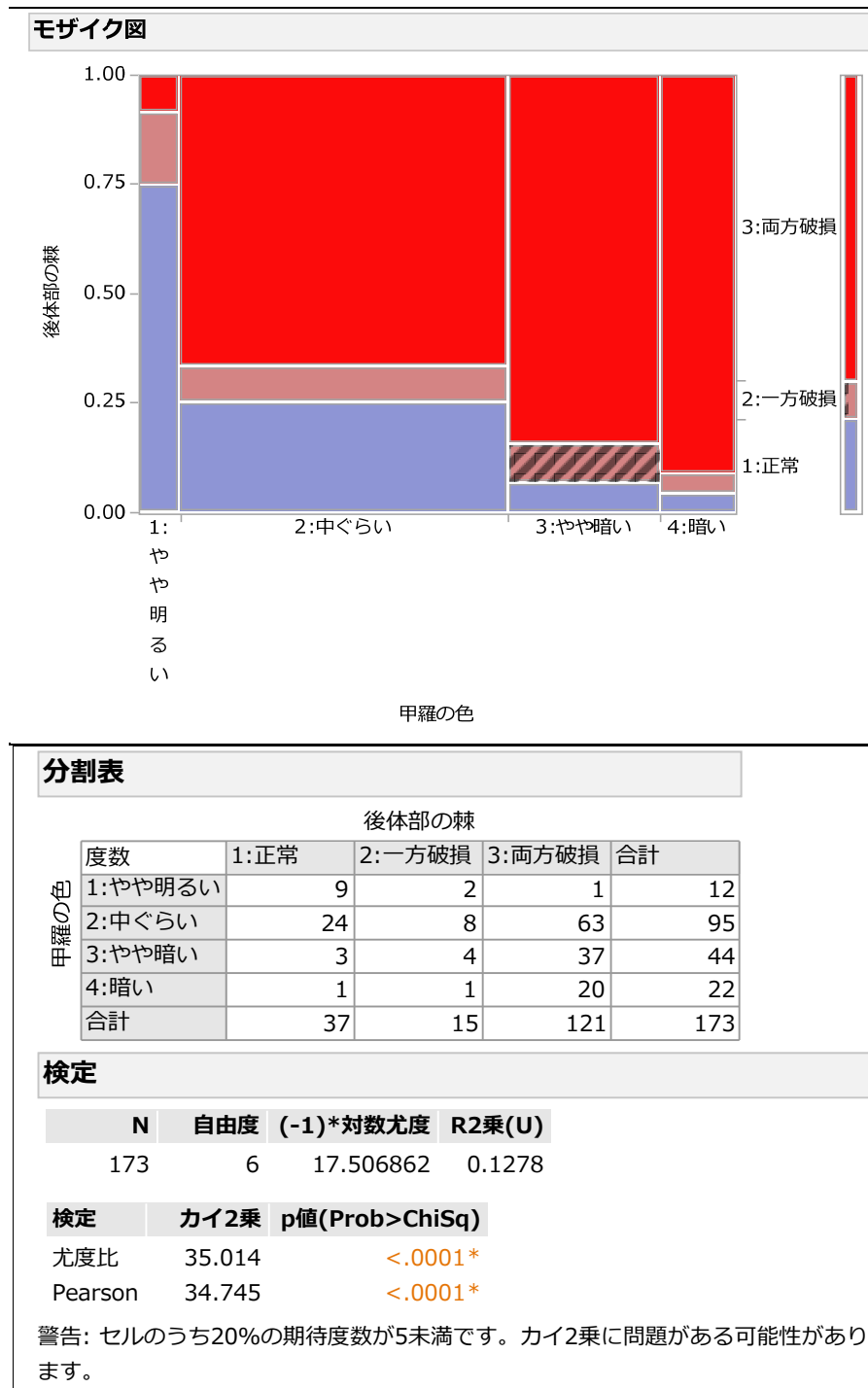
◆ 過分散は解消しない.

◆ 甲羅の色と棘の状態に関連がある.

2019.09.06 高橋行雄
13

甲羅の色と後体部の棘の分布は、明らかに均一ではないが、念のために 4×3 の分割表の検定を行う。甲羅の色が、やや明るいから暗いにかわるに従い後体部の棘の破損が増えることが読み取れる。

4×3 の分割表に対する Pearson の検定



JMP「二変量の関係→X: 甲羅の色→Y: 後体部の棘→OK」

8. 甲羅の幅か体重か

過分散を承知で、対数リンクによる 2 変数のポアソン回帰を行い、幅か体重か、どちらがサテライト数との関連が高いか検討する。表 6 およびスライド 14 に示すように、甲羅の幅の推定値は、0.0461、体重の推定値は、0.4470 であり、尤度比検定の結果は、体重のみが有意な差であった。

表 6 対数リンクによるポアソン 2 変量回帰

項	推定値	標準誤差	尤度比カイ2乗	p値
切片	-1.2952	0.8989	2.0691	0.1503
甲羅の幅	0.0461	0.0467	0.9658	0.3257
体重	0.4470	0.1586	7.9780	0.0047

スライド 15 に示すように、単独では共にサテライト数との正の関連が見いだせるのであるが、2 変量回帰で他方が関連なし、あるいは、負の場合の解釈は注意を要する。内部構造の可視化するために、JMP の「予測プロファイル」が優れている。

スライド 14

甲羅の幅と体重

- ◆ サテライト数に対する2変数の影響を対数リンクによるポアソン2変量回帰

項	推定値	標準誤差	尤度比カイ2乗	p値
切片	-1.2952	0.8989	2.0691	0.1503
甲羅の幅	0.0461	0.0467	0.9658	0.3257
体重	0.4470	0.1586	7.9780	0.0047

- ◆ 甲羅の幅の推定値は、0.0461
- ◆ 体重の推定値は、0.4470
- ◆ 尤度比検定の結果は、体重のみが有意過分散の調整なし

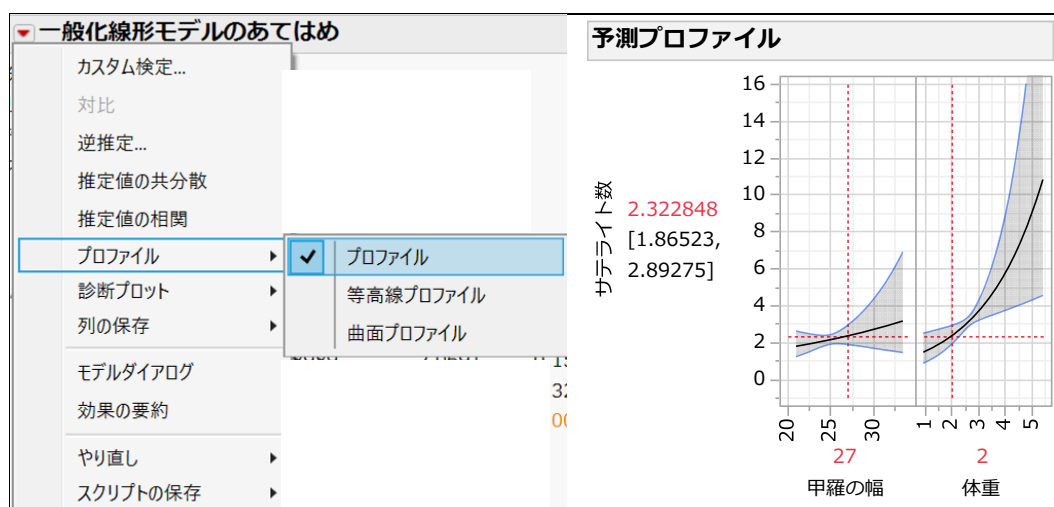
2019.09.06 高橋行雄 14

スライド 15

内部構造の可視化

- ◆ 2変量回帰での係数の解釈
 - 単独では共に関連があるが、2変量で他方が関連なし、あるいは、負の場合の解釈。
- ◆ JMPの「予測プロファイル」で内部構造の可視化する。
 - 体重を(2, 3, 4 kg)と変化させた場合の甲羅の幅がサテライト数に及ぼす影響を図示
 - 甲羅の幅は、体重を固定したときにサテライト数に影響を与えていない。

2019.09.06 高橋行雄 15



予測プロファイルの X 軸の甲羅の幅を 27 にセットし、体重を(2→3→ 4)と変化しつつコピー & ペーストを繰り返す。

図 4 およびスライド 16 は、JMP による対数リンクでの 2 変量ポアソン回帰に引き続き「予測プロファイル」の機能を用い、体重を(2, 3, 4 kg)と変化させた場合の甲羅の幅がサテライト数に及ぼす影響を図示したものである。甲羅の幅は体重の増加に伴いサテライト数も増加しているが、95%信頼区間の表示から、傾きがマイナスになる可能性があることが読み取れ、このことが表 6 の p 値が大きいことに対応する。

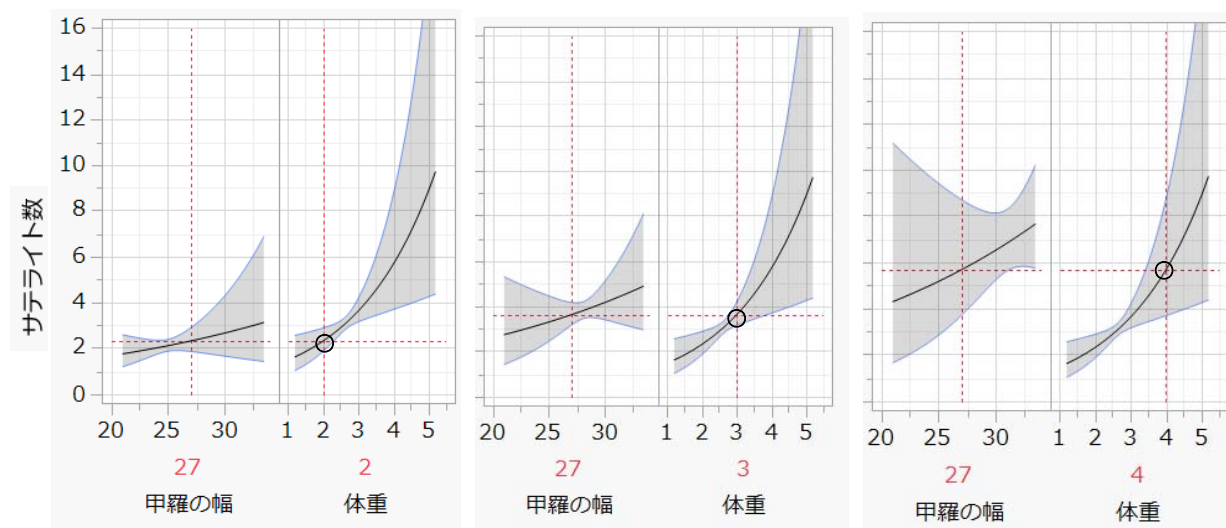
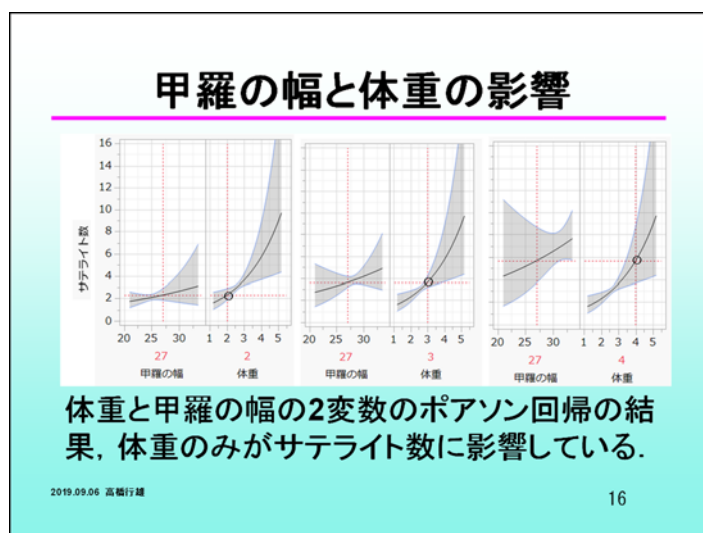


図 4 体重を変化させた場合の甲羅の幅とサテライト数との関連

スライド 16

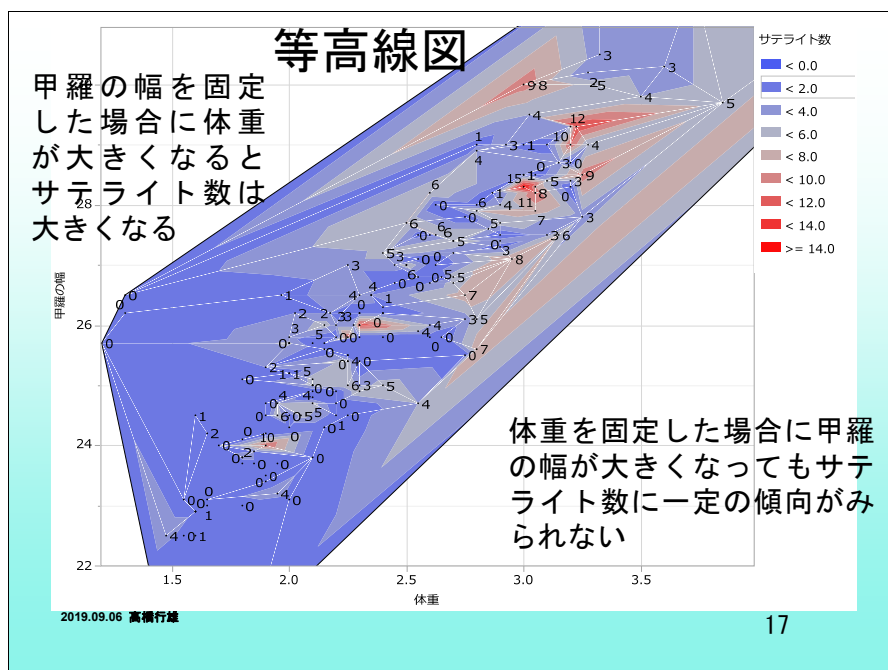


第 12 節で Excel を使った作図方法を示す。

「予測プロファイル」は、JMP ではお馴染みであり、標示されている変数のグラフの X 軸の下の変数に自由に変更でき、他の変数のプロファイルもダイナミックに変化するので、多変量の回帰モデルの結果の解釈になくてはならない。ただし、「予測プロファイル」は、一般的な統計の教科書には乱されないもので、第 12 節で Excel を使って、ここで示されている「予測プロファイル」を再現する。

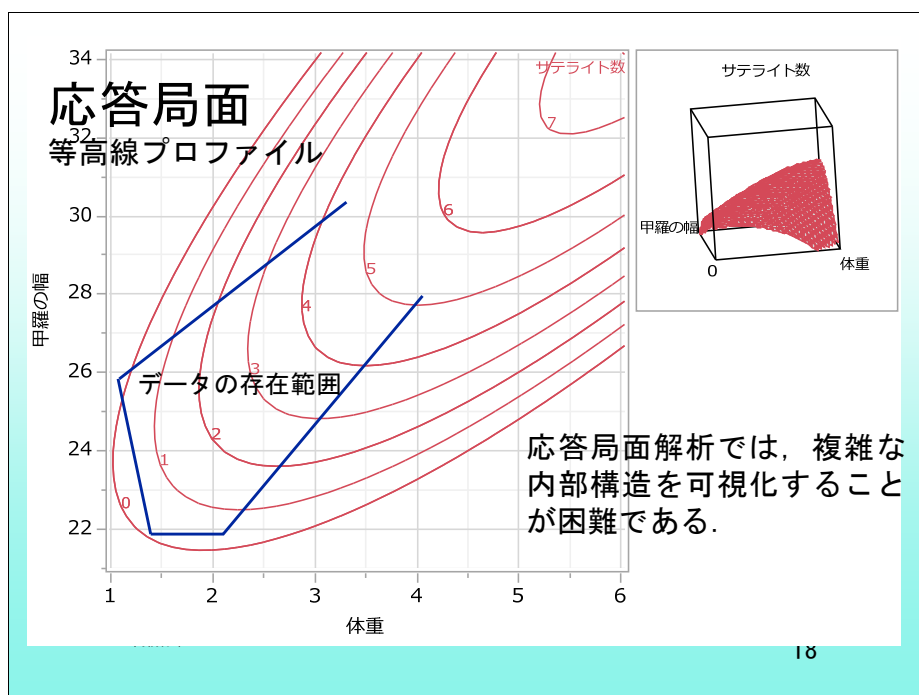
スライド 17 に内部構造を可視化するために等高線図を示す。図の中の数字がサテライト数であるが、変動が大きいため明確な等高線を描くことができない。

スライド 17



スライド 18 に応答局面解析の結果を示すが、等高線図と同様に、サテライト数での変動が大きいためか、甲羅の幅と体重のサテライト数与える影響についての解釈することは難しい。

スライド 18



9. 甲羅の色と体重の組み合わせ

さて、甲羅の色が暗くなるにつれて棘の破損が多くなり、サテライト数が減ることを表 5 で示した。では、甲羅の色と体重を組み合わせた場合に、何らかの関連が見出されるのであろうか。このような関連を、ポアソン回帰で見い出すためには、甲羅の色について何らかの数値を与えてデザイン行列化し、体重との交互作用を含めたポアソン回帰を行う必要がある。スライド 19 に示すように JMP の一般化線形モデルでは、名義尺度に対しては対比型のデザイン行列を自動生成するので、「モデル効果の構成」で(甲羅の色, 体重, 甲羅の色*体重)を設定すればよい。

スライド 19

甲羅の色と体重

- ◆ 甲羅の色を対比型のダミー変数とする
- ◆ 交互作用モデル
 - 甲羅の色, 体重, (甲羅の色 × 体重)
 - 回帰パラメータによる解釈は難解
- ◆ 「予測プロファイル」の機能
 - 甲羅の色ごとの体重の増加によるサテライト数との関連を概観

2019.09.06 高橋行雄
19

パラメータの推定結果を表 7 に示すが、このままでは、結果の解釈は困難を極めるので、「予測プロファイル」の機能を用いて図 5 およびスライド 20 に示すように甲羅の色ごとの体重の増加によるサテライト数との関連を概観する。

表 7 甲羅の色と体重の交互作用を含めた対数リンクでのポアソン重回帰

項	推定値	標準誤差	尤度比カイ2乗	p値
切片	-0.2778	0.3450	0.6530	0.4191
甲羅の色[1:やや明るい]	2.2221	0.7978	7.5086	0.0061*
甲羅の色[2:中ぐらい]	0.2010	0.3797	0.2812	0.5959
甲羅の色[3:やや暗い]	-1.1855	0.4865	6.0352	0.0140*
体重	0.5463	0.1344	16.0804	<.0001*
甲羅の色[1:やや明るい]*体重	-0.7518	0.3050	6.1530	0.0131*
甲羅の色[2:中ぐらい]*体重	-0.0646	0.1456	0.1967	0.6574
甲羅の色[3:やや暗い]*体重	0.3820	0.1870	4.2010	0.0404*

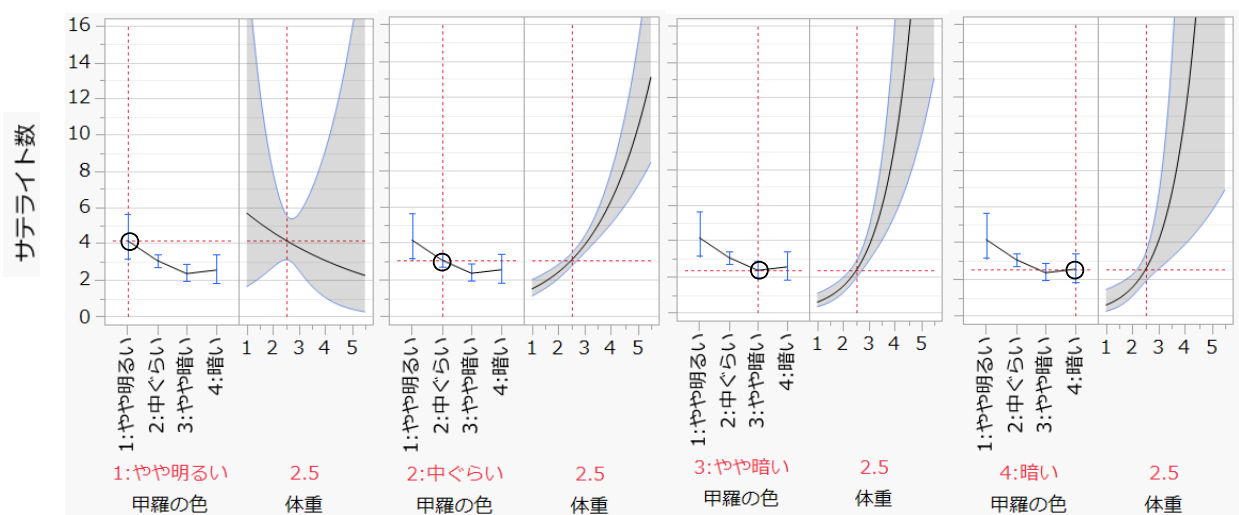
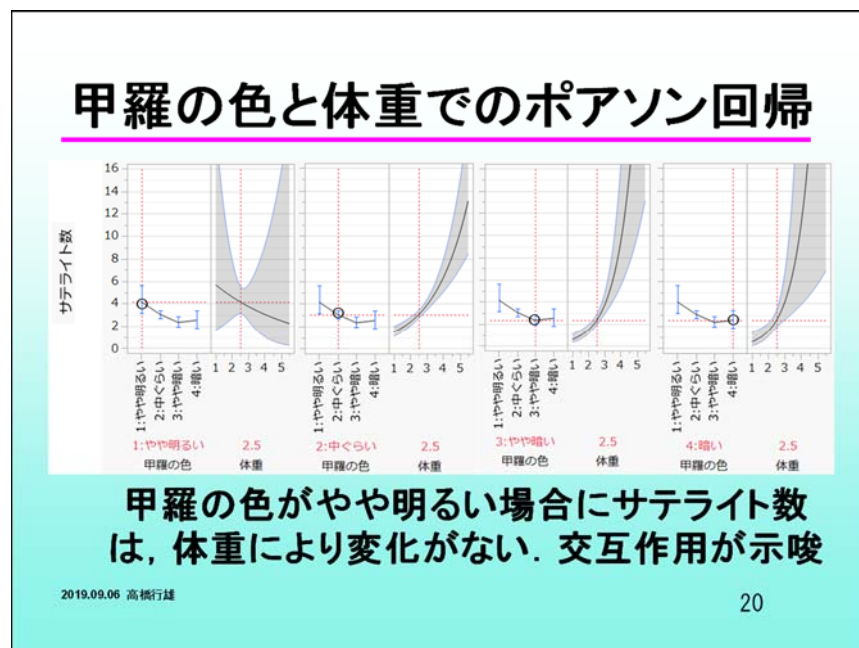


図5 甲羅の色を変化させた場合の体重とサテライト数の関連

スライド 20



第12節で Excel を使った作図方法を示す。

10. 後体部の棘と体重の組み合わせ

予測プロファイルから、甲羅の色が「やや明るい」場合は、体重とサテライト数の関連は、マイナスの傾きも起こりえる信頼区間となっており、関連はみいだせないことが他の色と明らかに異なる。「中ぐらい」以上では、体重が増えればサテライト数も増大する。

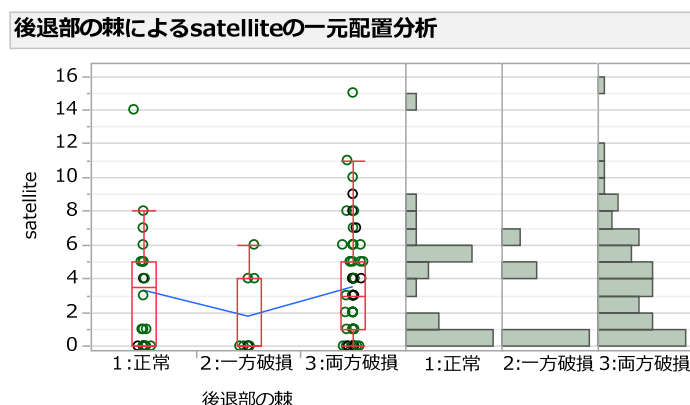


図 6 補足 甲羅の色が「中ぐらい」の場合の後退部の棘とサテライト数の関連

後体部の棘の状態は、甲羅の色によって破損が進行することを表 5 で明らかにした。甲羅の色が「中ぐらい」の場合には、後部の棘が「正常」と「両方破損」に分かれているので、サテライト数との関連を甲羅の色が「中ぐらい」の 95 ツガイに限定して関連を調べた結果を図 6 に示す。

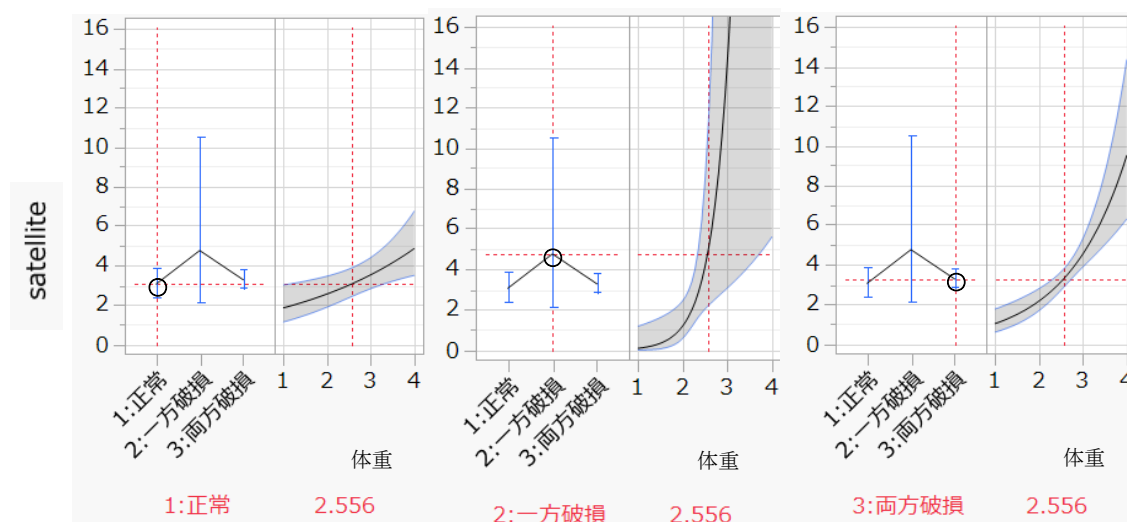


図 6 甲羅の色が「中ぐらい」での後部の棘の状態別の体重とサテライト数の関連

図 6 およびスライド 22 から甲羅の色が「中ぐらい」で後体部の棘が「正常」の場合に体重が増えればサテライト数も微増する。「一方破損」および「両方破損」では、体重が増えた場合にサテライト数が急増する。表 5 から、甲羅の色が「やや明るい」場合には、後体部の棘は 12 ツガイ中 9 ツガイが「正常」

で、図 5 から体重が増えてもサテライト数は増えない。甲羅の色が「中ぐらい」に変化すると、体重が増加するとサテライト数も大幅に増える。更に色が「やや暗い、暗い」場合には、更に体重が増えるにつれて、サテライト数が増えるとも言えるが、体重が小さい場合には、サテライト数が減少することが読み取れる。

甲羅の色が「中ぐらい」の場合は、表 5 から後体部の棘の(正常:両方破損)=(24:63)と、ある程度数が揃っているので、サテライト数との関連を甲羅の色が「中ぐらい」に限定して関連を調べる。その結果、後体部の棘が「正常」の場合に体重が増えればサテライト数は微増し、「一方破損」および「両方破損」では、体重が増えた場合にサテライト数が急増することが観察される。

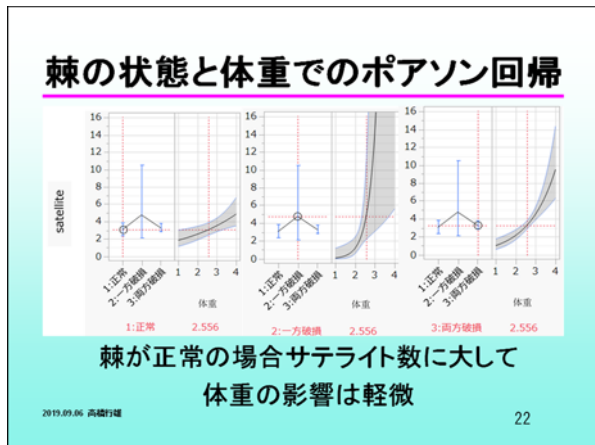
スライド 21

後退部の棘の状態を含めた解析

- ◆ 甲羅の色が「**中ぐらい**」の場合
 - 後体部の棘:(正常:両方破損)=(24:63)
- ◆ サテライト数との関連を甲羅の色が「中ぐらい」に限定して関連を調べる
 - 後体部の棘が「正常」の場合に体重が増えればサテライト数は微増。
 - 「一方破損」および「両方破損」では、体重が増えた場合にサテライト数が急増

2019.09.06 高橋行雄 21

スライド 22



11. 層別散布図行列における回帰の 95%信頼区間

交互作用が疑われるような探索的な解析を行うためには、各種のグラフ表示が欠かせない。これまでも JMP の多彩なグラフ表示を活用し、カブトガニの各種の変数とサテライト数の関連を浮き彫りにしてきたが、満足できるものではなかった。全体を俯瞰できるように結果を 1 枚のグラフで表わすことは、可能なのだろうか。

交互作用が疑われるような探索的な解析を行うためには、各種のグラフ表示が欠かせない。JMP の多彩なグラフ表示を活用し、カブトガニの各種の変数とサテライト数の関連を浮き彫りにしてきたが、今ひとつものたりない。全体を俯瞰するために JMP の「グラフ・ビルダー」が有益である。

スライド 23

交互作用がある場合の解析

- ◆ 交互作用が疑われるような探索的な解析を行うためには、各種のグラフ表示が欠かせない。
- ◆ JMPの多彩なグラフ表示を活用し、カブトガニの各種の変数とサテライト数の関連を浮き彫りにしてきたが、今ひとつである。
- ◆ 全体を俯瞰するためにJMPの「グラフ・ビルダー」が有益である。

2019.09.06 高橋行雄

23

S プラスには、グラフ・ビルダーと同様の機能があり愛用していた。JMP グラフ・ビルダーは、S プラスの機能を大幅に凌駕する探索的な統計解析を支援するツールとして優れている。交互作用がある場合には、図を作成して解釈することが鉄則である。交互作用を把握するための作図は容易ではないが、グラフ・ビルダーで容易に把握できる

JMPのグラフ・ビルダー

- ◆ Sプラスに同様の機能があり愛用していた
- ◆ JMPグラフ・ビルダーは, Sプラスの機能を大幅に凌駕する探索的な統計解析を支援するツールとして優れている.
- ◆ 交互作用がある場合には, 図を作成して解釈することが鉄則である.
 - 交互作用を把握するための作図は容易ではないが, グラフ・ビルダーで容易に把握できる

2019.09.06 高橋行雄

24

JMP の新しい作図機能である「グラフ・ビルダー」を用いた結果を図 7 およびスライド 25 に示す. この図から, これまでの探索的解析の結果がより鮮明に浮彫される. サテライト数は, 甲羅の色が暗くなるにつて後体部の棘の破損が進み, それに伴い, 体重の軽い雌ほど連結する雄のサテライト数が減少することが読み取れる. 甲羅の色が暗くなり, 後体部の棘の状態が悪くなる加齢現象により, 体重の軽い雌ほど連結する雄のサテライト数が減少すると解される. そのため, ゼロ・カウントが多い過分散となったと推測される.

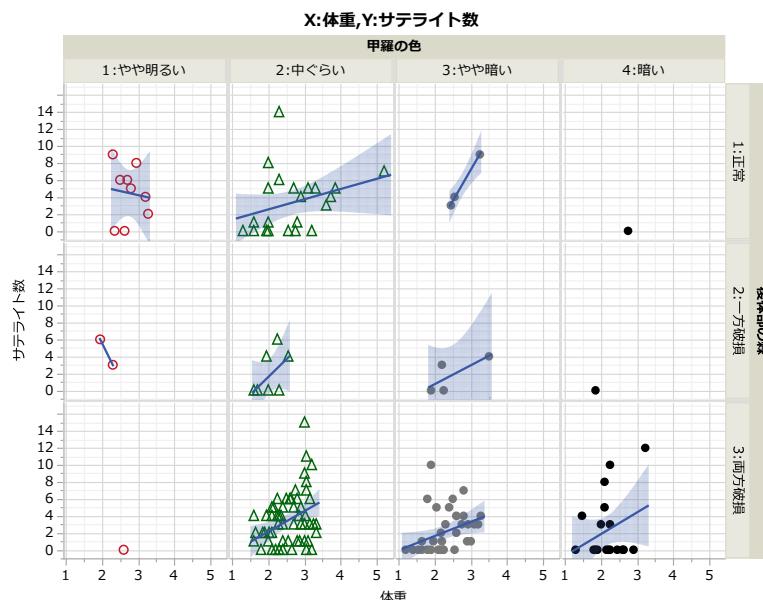
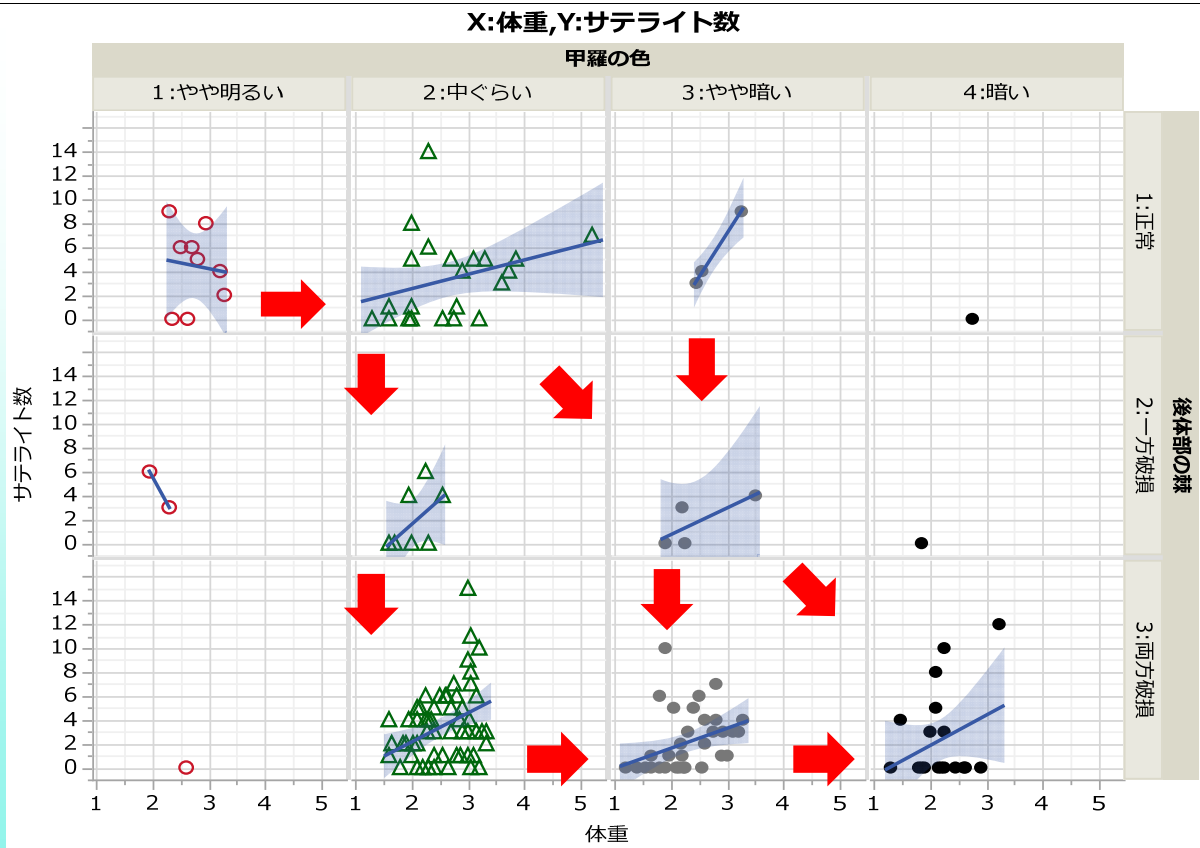


図 7 甲羅の色・棘の状態による層別散布図での回帰の 95%信頼区間の表示

スライド 25 加齢現象による変化



2019.09.06 高橋行雄

25

スライド 26

交互作用の解釈

- ◆ 甲羅の色が暗くなるにつれて後体部の棘の破損が進む
- ◆ それに伴い、体重の軽い雌ほど連結する雄のサテライト数が減少
- ◆ 甲羅の色が暗くなり、後体部の棘の状態が悪くなる“加齢現象”により、体重の軽い雌ほど連結する雄のサテライト数が減少する
- ◆ そのため、ゼロ・カウントが多い過分散となったと推測される。

2019.09.06 高橋行雄

26

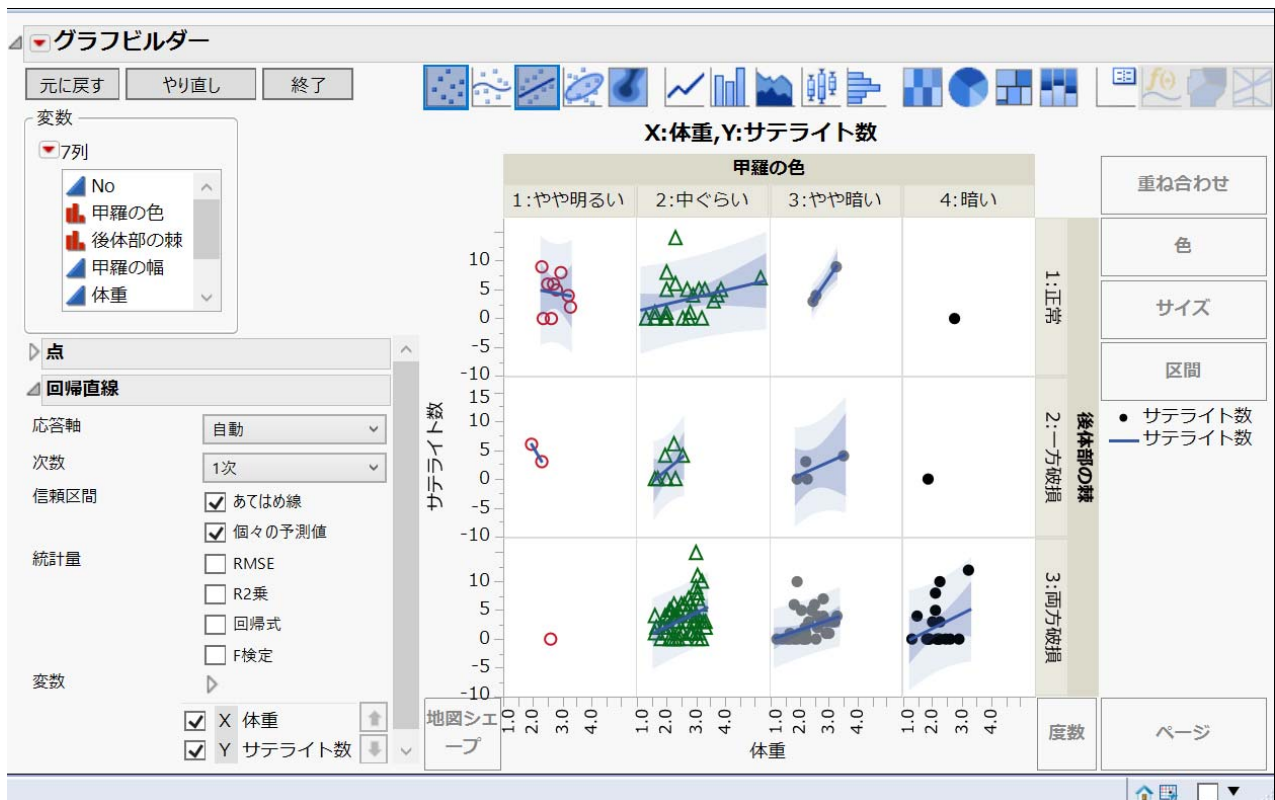
グラフ・ビルダー再考

- ◆ 2つの順序変数を組み合わせて(3×4)のセル別に体重とサテライト数の散布図に回帰直線と95%信頼区間が描けることにより、変数間の関係が一覧でき、内部構造の可視化に有益であった。
- ◆ ただし、ポアソン回帰ではなく、(3×4)のセル別の通常の回帰の結果となっているので、解釈には注意が必要である。

2019.09.06 高橋行雄

27

グラフ・ビルダーで、個別データの予測値を加えたところ



12. Excel による予測プロファイル

単回帰分析を行った場合に、説明変数 X と応答変数 Y の散布図上に回帰直線を描き、回帰直線の 95%信頼区間および個別データの 95%を上書きすることは、多くの統計ソフトで標準的にサポートされていて、回帰分析の結果を解釈する際に大いに役に立つ。説明変数に X^2 の項を入れて 2 次曲線をあてはめた場合の 95%信頼区間についての例示を目にすることはまれであり、その計算式を目にすることもまれである。

JMP では、単回帰のみならず多項式回帰の曲線の信頼区間(平均の信頼区間)および個別データの 95%信頼区間(予測区間)も標準的にサポートされており、それらの計算式を参照することも可能となっている。各種の予測プロファイルの作成には、「推定値の共分散行列」が欠かせない。紛らわしいのは、「多変数の共分散行列」が一般的に知られているが、はっきりと区別しなければならない。なお、Excel による対数リンクのポアソン回帰の共分散行列の算出方法については、第 14 節を参照のこと。

12.1 甲羅の幅と体重の 2 変量ポアソン回帰における予測プロファイル

説明変数が 2 変数のポアソン回帰の場合、95%信頼区間をどのように図示したら良いのであろうか。JMP では、1 変量ごとの「予測プロファイル」、2 変数を組み合わせた「交互作用プロファイル」、「等高線プロファイル」、「曲面プロファイル」など多彩なプロファイルを選択して結果を解釈するための“プロファイル”が描けるようになっている。基本は、1 変量ごとの予測プロファイルであり、(図 4, 図 5, 図 6)に示してきたように、2 変数 (X_1 と X_2) の値について X_1 を固定し X_2 を変化させたときの Y 推定値(平均)と平均の 95%信頼区間が図示される。さらに、 X_2 を固定し X_1 を変化させたときの Y 推定値(平均)と 95%信頼区間が図示される。

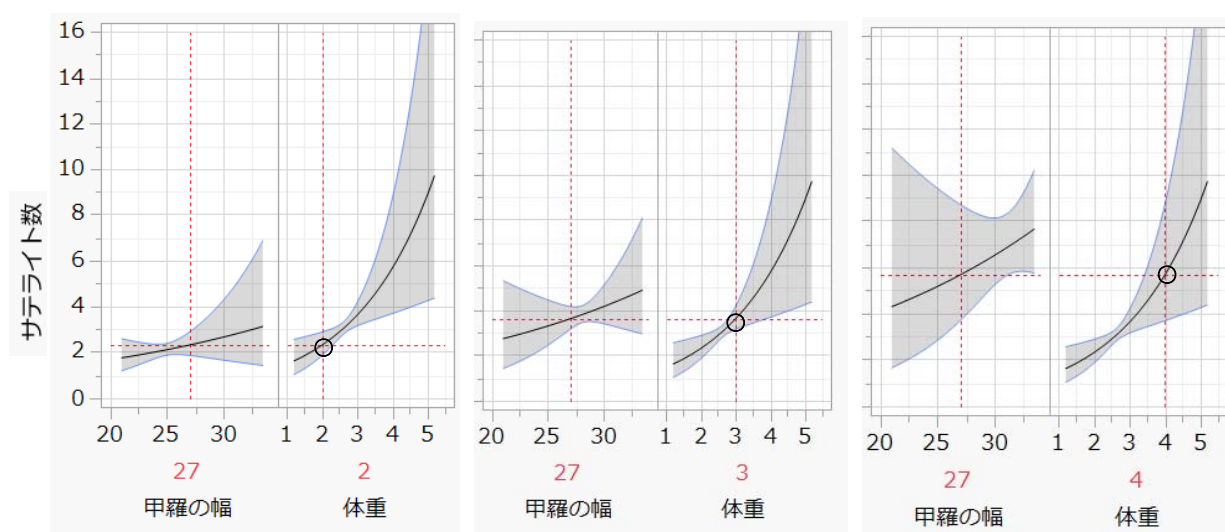


図 4 (再掲) 体重を変化させた場合の甲羅の幅とサテライト数との関連

これまで(図 4, 図 5, 図 6)で示してきた予測プロファイルの作成方法を Excel で示す. 予測プロファイルの作成に必要なのは, 回帰パラメータの推定値および推定値の共分散行列である. 共分散行列は, 標準的には出力されず, 何らかのオプション指定が必要である. JMP の場合は, 「推定値の共分散」を選択する.

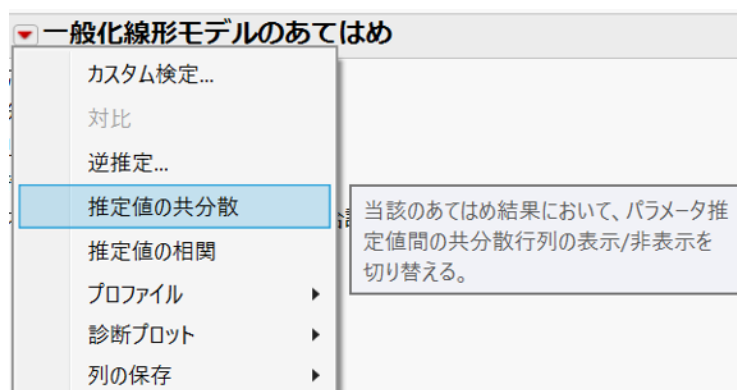


表 8 に示すようにパラメータの推定値は, 表 6 で示した結果に一致する. 「推定値の共分散」には, 「切片」が存在することより, 「多変数の共分散行列」とははっきりと区別することができる. 推定値の標準誤差は, 共分散行列 Σ の対角要素の平方根で得られる. 甲羅の幅の推定値の場合に $SE = 0.8989$ であるが, 共分散行列の対角要素 0.0022 が甲羅の幅の推定値の分散となっていることから, その平方根が SE となる.

表 8 対数リンクのポアソン回帰の推定値および共分散行列

パラメータ推定値			推定値の共分散			
項	推定値	標準誤差	共分散			
切片	-1.2952	0.8989	切片	0.8080	-0.0412	0.1156
甲羅の幅	0.0461	0.0467	甲羅の幅	-0.0412	0.0022	-0.0067
体重	0.4470	0.1586	体重	0.1156	-0.0067	0.0252

パラメータの推定値をベクトル $\hat{\beta}$

$$\hat{\beta} = [-1.2952, 0.0461, 0.4470]^T$$

とし, 切片項 $x_0 = 1$, 甲羅の幅 x_1 , 体重 x_2 をベクトル \mathbf{x}

$$\mathbf{x} = [1, x_1, x_2]$$

としたときに, 推定値 $\ln \hat{y}$ は, ベクトルの積

$$\ln \hat{y} = \mathbf{x} \hat{\beta}$$

となり, 予測値の分散 $Var(\ln \hat{y})$ は, 共分散 Σ を挟んだ \mathbf{x} の 2 次形式

$$Var(\ln \hat{y}) = \mathbf{x} \Sigma \mathbf{x}^T$$

である. 95%信頼区間は,

$$(\ln U95\%, \ln L95\%) = \ln \hat{y} \pm 1.96 \sqrt{\text{Var}(\ln \hat{y})}$$

であり, 元のスケールでは,

$$(U95\%, L95\%) = \exp[\ln \hat{y} \pm 1.96 \sqrt{\text{Var}(\ln \hat{y})}]$$

で計算することができる.

実際の数値で計算過程を示す. Excel には, 行列計算のための関数があり, これらの関数を用いた計算を行う. ベクトルおよび行列は, 矩形で囲んだ表記とする. パラメータの推定値を $\hat{\beta}$, 甲羅の幅 $x_1 = 27$, 体重 $x_2 = 2$ とした場合の \mathbf{x} , 共分散行列を Σ についての Excel シートを示す. なお, 計算結果は, 表 9 にまとめられている.

$\hat{\beta} =$	-1.2952		
	0.0461		
	0.4470		
$\mathbf{x} =$	1	27	2
$\Sigma =$	0.8080	-0.0412	0.1156
	-0.0412	0.0022	-0.0067
	0.1156	-0.0067	0.0252

推定値 $\ln \hat{y} = \mathbf{x} \hat{\beta}$ は, 行列の積の関数 Mmult() を使って 0.8428 と計算される.

$\ln \hat{y} =$	1	27	2	-1.2952	=	0.8428
				0.0461		
				0.4470		
$\ln \hat{y} = \mathbf{x} \hat{\beta} = \text{Mmult}(\mathbf{x}, \hat{\beta})$						

推定値の分散 $\text{Var}(\ln \hat{y}) = \mathbf{x} \Sigma \mathbf{x}^T$ は, 行列の積の関数 Mmult() および 転置の関数 Transpose() を使って 0.0125 と計算される.

$\text{Var}(\ln \hat{y}) =$	1	27	2	0.8080	-0.0412	0.1156	1	=	0.0125
				-0.0412	0.0022	-0.0067	27		
				0.1156	-0.0067	0.0252	2		
$\text{Var}(\ln \hat{y}) = \mathbf{x} \Sigma \mathbf{x}^T = \text{Mmult}(\text{mmult}(\mathbf{x}, \hat{\beta}), \text{Transpose}(\mathbf{x}))$									

推定値の 95%信頼区間は, $(\ln U95\%, \ln L95\%) = \ln \hat{y} \pm 1.96 \sqrt{\text{Var}(\ln \hat{y})}$ なので, 元のスケールでは,

$$\exp(\ln \hat{y}) = \exp(0.8428) = 2.3228$$

$$L95\% = \exp[0.8428 - 1.96 \times \text{Sqrt}(0.0125)] = \exp(0.6234) = 1.8653$$

$$U95\% = \exp[0.8428 + 1.96 \times \text{Sqrt}(0.0125)] = \exp(1.0622) = 2.8927$$

となる. 体重を 1, 2, 3, 4, 5 kg と変化させた場合について, 表 9 に計算結果を示す. 表 10 には, 体重を 2 kg と固定し甲羅の幅を 21~33 cm と変化させた結果を示す.

表 9 甲羅の幅を 27cm と固定し体重を 1~5 kg に変化させた推結果

		推定値 β	共分散 Σ	切片	甲羅の幅	体重	
	切片	-1.2952	β_0	0.8080	-0.0412	0.1156	
	甲羅の幅	0.0461	β_1	-0.0412	0.0022	-0.0067	
	体重	0.4470	β_2	0.1156	-0.0067	0.0252	
x_0	x_1	x_2	$\ln \hat{y}$	$Ver(\ln \hat{y})$	\hat{y}	L95%	U95%
1	28.3	3.05	1.3720	0.0026	3.9433	3.5670	4.3592
1	27	1	0.3958	0.0703	1.4856	0.8836	2.4978
1	27	2	0.8428	0.0125	2.3228	1.8653	2.8927
1	27	3	1.2898	0.0051	3.6319	3.1574	4.1778
1	27	4	1.7367	0.0480	5.6788	3.6963	8.7245
1	27	5	2.1837	0.1412	8.8791	4.2511	18.5457

表 10 体重を 2 kg と固定し甲羅の幅を 21~33 cm に変化させた推定結果

x_0	x_1	x_2	$\ln \hat{y}$	$Ver(\ln \hat{y})$	\hat{y}	L95%	U95%
1	21	2	0.5663	0.0391	1.7618	1.1959	2.5954
1	24	2	0.7046	0.0061	2.0230	1.7352	2.3584
1	27	2	0.8428	0.0125	2.3228	1.8653	2.8927
1	30	2	0.9810	0.0583	2.6672	1.6618	4.2809
1	33	2	1.1193	0.1434	3.0626	1.4581	6.4324

表 8 および表 9 の結果を Excel の「散布図」を用いて描いた結果を図 7 に示す. この結果は, JMP で求めた図 4 の体重を変化させた場合の甲羅の幅とサテライト数との関連の左側に一致する. なお, Excel の「散布図」を使った作図方法については, 第 13 節に示す.

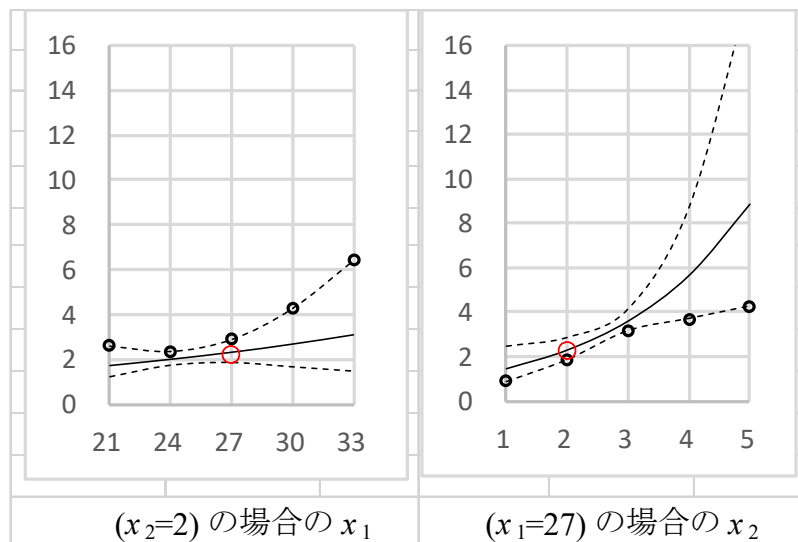


図 8 Excel による予測プロファイル 1

図 4 では、甲羅の幅を 27 cm に固定し、体重を 3 kg および 4 kg に固定した場合についても示してあるので、Excel では、体重を 4 kg に固定した場合についての結果を表 11 に示す。

表 11 体重を 4 kg と固定し甲羅の幅を 21～33 cm に変化させた推定結果

		推定値 β	共分散 Σ	切片	甲羅の幅	体重	
	切片	-1.2952	β_0	0.8080	-0.0412	0.1156	
	甲羅の幅	0.0461	β_1	-0.0412	0.0022	-0.0067	
	体重	0.4470	β_2	0.1156	-0.0067	0.0252	
x_0	x_1	x_2	$\ln y^\wedge$	$Ver(\ln y^\wedge)$	y^\wedge	$L95\%$	$U95\%$
1	28.3	3.05	1.3720	0.0026	3.9433	3.5670	4.3592
1	27	1	0.3958	0.0703	1.4856	0.8836	2.4978
1	27	2	0.8428	0.0125	2.3228	1.8653	2.8927
1	27	3	1.2898	0.0051	3.6319	3.1574	4.1778
1	27	4	1.7367	0.0480	5.6788	3.6963	8.7245
1	27	5	2.1837	0.1412	8.8791	4.2511	18.5457
1	21	4	1.4603	0.2365	4.3071	1.6604	11.1729
1	24	4	1.5985	0.1226	4.9456	2.4899	9.8233
1	27	4	1.7367	0.0480	5.6788	3.6963	8.7245
1	30	4	1.8750	0.0127	6.5206	5.2262	8.1355
1	33	4	2.0132	0.0168	7.4872	5.8061	9.6551

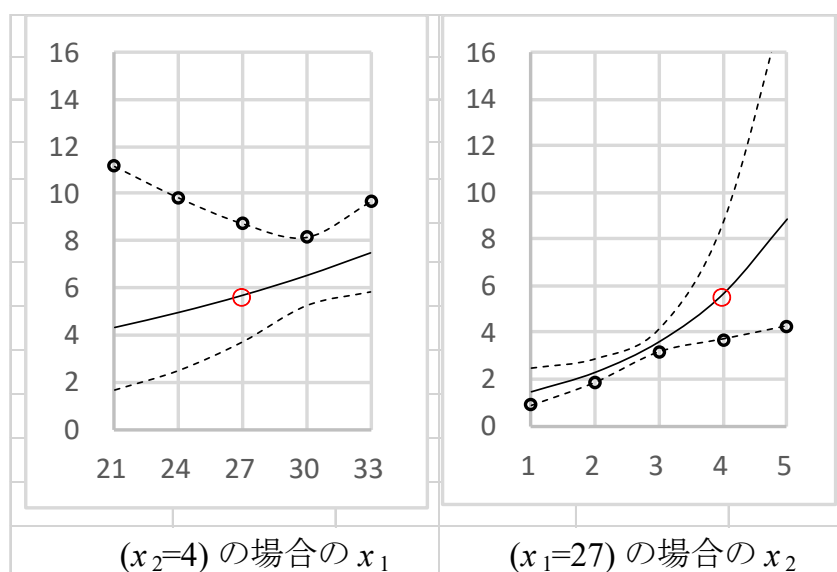


図 9 Excel による予測プロファイル2

図 9 の予測プロファイル右側は、 $x_1 = 27$ に固定したときに x_2 を変化させた場合であり、図 8 の左側と全く同じである。あえて並べてあるのは、右側の $x_2 = 4$ 場合の推定値 5.6788 と左側の $x_1 = 27$ が同じになることを明示するためである。

12.2 順序データ(甲羅の色)と体重の交互作用を含む場合

甲羅の色は(1:やや明るい, 2:中ぐらい, 3:やや暗い, 4:暗い)の4段階の順序データである。JMPでのポアソン回帰の場合は、変数を(名義尺度, 順序尺度, 連続尺度)のどれかに設定するのであるが、順序尺度のダミー変数と名義尺度ダミー変数が異なるので、名義尺度として計算を行った。

JMPの名義尺度のダミー変数は、表12に示すように対比型(ダミー変数 x_i について足し算すると0となる)とする。体重と甲羅の色の交互作用は、体重とそれぞれのダミー変数の積となる。表7に対数リンクによるポアソン回帰の推定結果、図5に予測プロファイルを示した。変数の数は多くなるが、予測プロファイル作成のための計算方法は、2変数のポアソン回帰の場合と考え方は同じである。

表12 対比型のダミー変数

甲羅の色	x_1	x_2	x_3
1:やや明るい	1	0	0
2:中ぐらい	0	1	0
3:やや暗い	0	0	1
4:暗い	-1	-1	-1

実際の数値での計算結果を示す。Excelには、行列計算のための関数があり、これらの関数を用いた計算を行う。ベクトルおよび行列は、矩形で囲んだ表記とする。

表13 パラメータの推定値(表7再掲)

変数名	変数内容			推定値	標準誤差
x_0	切片	切片	β_0	-0.2778	0.3450
x_1	A1	[1:やや明るい]	β_1	2.2221	0.7978
x_2	A2	[2:中ぐらい]	β_2	0.2010	0.3797
x_3	A3	[3:やや暗い]	β_3	-1.1855	0.4865
x_4	W体重	体重	β_4	0.5463	0.1344
x_5	A1×W	[1:やや明るい]*体重	β_5	-0.7518	0.3050
x_6	A2×W	[2:中ぐらい]*体重	β_6	-0.0646	0.1456
x_7	A3×W	[3:やや暗い]*体重	β_7	0.3820	0.1870

パラメータの推定値 $\hat{\beta}$ 、甲羅の色 $x_1=1$, $x_2=0$, $x_3=0$, 体重 $x_4=2.5$, $x_5=1 \times 2.5$, $x_6=0 \times 2.5$, $x_7=0 \times 2.5$ とした場合の \mathbf{x} 、共分散行列を Σ とした場合についてExcelシートを示す。

表 14 共分散行列(計算結果を JMP ファイルに出力し整形)

	行	切片	甲羅の 色[1:や や明るい]	甲羅の 色[2:中 ぐらい]	甲羅の 色[3:や や暗い]	体重	甲羅の 色[1:や や明るい]*体 重	甲羅の 色[2:中 ぐらい]*体 重	甲羅の 色[3:や や暗い]*体 重
1	切片	0.1190	0.1397	-0.1064	-0.0602	-0.0456	-0.0519	0.0413	0.0238
2	甲羅の色[1:やや明るい]	0.1397	0.6365	-0.1523	-0.1986	-0.0519	-0.2407	0.0562	0.0737
3	甲羅の色[2:中ぐらい]	-0.1064	-0.1523	0.1442	0.0476	0.0413	0.0562	-0.0542	-0.0195
4	甲羅の色[3:やや暗い]	-0.0602	-0.1986	0.0476	0.2367	0.0238	0.0737	-0.0195	-0.0893
5	体重	-0.0456	-0.0519	0.0413	0.0238	0.0181	0.0194	-0.0165	-0.0096
6	甲羅の色[1:やや明るい]*体重	-0.0519	-0.2407	0.0562	0.0737	0.0194	0.0931	-0.0210	-0.0279
7	甲羅の色[2:中ぐらい]*体重	0.0413	0.0562	-0.0542	-0.0195	-0.0165	-0.0210	0.0212	0.0080
8	甲羅の色[3:やや暗い]*体重	0.0238	0.0737	-0.0195	-0.0893	-0.0096	-0.0279	0.0080	0.0350

推定値 $\ln \hat{y} = \mathbf{x}\hat{\boldsymbol{\beta}}$ の計算はサイズが大きいが、行列の積の関数 `Mmult()` を使って 1.4306 と計算される。

	切片	甲羅の幅			体重	甲羅の幅×体重					
		A1	A2	A3	W	A1×W	A2×W	A3×W			
	x_0	x_1	x_2	x_3	x_4	x_5	x_6	x_7			
	\mathbf{x}								$\hat{\boldsymbol{\beta}}$		
$\ln y^{\wedge} =$	1	1	0	0	2.5	2.5	0	0	-0.2778	=	1.4306
									2.2221		
									0.2010		
									-1.1855		
									0.5463		
									-0.7518		
									-0.0646		
									0.3820		
									$\ln \hat{y} = \mathbf{x}\hat{\boldsymbol{\beta}} = \text{Mmult}(\mathbf{x}, \hat{\boldsymbol{\beta}})$		

推定値の分散 $\text{Var}(\ln \hat{y}) = \mathbf{x} \boldsymbol{\Sigma} \mathbf{x}^T$ は、行列の積の関数 `Mmult()` および 転置の関数 `Transpose()` を使って 0.0219 と計算される。

	\mathbf{x}								$\boldsymbol{\Sigma}$								\mathbf{x}^T		
$Var(\ln \hat{y})=$	1	1	0	0	2.5	2.5	0	0	0.12	0.14	-0.11	-0.06	-0.05	-0.05	0.04	0.02	1	=	0.0219
									0.14	0.64	-0.15	-0.20	-0.05	-0.24	0.06	0.07	1		
									-0.11	-0.15	0.14	0.05	0.04	0.06	-0.05	-0.02	0		
									-0.06	-0.20	0.05	0.24	0.02	0.07	-0.02	-0.09	0		
									-0.05	-0.05	0.04	0.02	0.02	0.02	-0.02	-0.01	2.5		
									-0.05	-0.24	0.06	0.07	0.02	0.09	-0.02	-0.03	2.5		
									0.04	0.06	-0.05	-0.02	-0.02	-0.02	0.02	0.01	0		
									0.02	0.07	-0.02	-0.09	-0.01	-0.03	0.01	0.03	0		
$Var(\ln \hat{y}) = \mathbf{x} \boldsymbol{\Sigma} \mathbf{x}^T = \text{Mmult}(\text{Mmult}(\mathbf{x}, \boldsymbol{\Sigma}), \text{Transpose}(\mathbf{x}))$																			

推定値の 95%信頼区間は, $(\ln U95\%, \ln L95\%) = \ln \hat{y} \pm 1.96\sqrt{\text{Var}(\ln \hat{y})}$ なので,

$$\exp(\ln \hat{y}) = \exp(1.431) = 4.1810$$

$$L95\% = \exp[0.8428 - 1.96 \times \sqrt{0.0125}] = \exp(1.1402) = 3.1275$$

$$U95\% = \exp[0.8428 + 1.96 \times \sqrt{0.0125}] = \exp(1.7209) = 5.5894$$

となる. . 表 15 には, 体重を 2.5 kg と固定し甲羅の色を(1, 2, 3, 4)と変化させた推定結果を示す.

体重を 1, 2, 3, 4, 5 kg と変化させた場合について, 表 16 に計算結果を示す.

表 15 体重が 2.5 kg の場合の甲羅の色についての予測値の計算

			推定値	共分散		甲羅の幅			体重	甲羅の幅×体重			
		項	β	Σ	切片	A1	A2	A3	W	A1×W	A2×W	A3×W	
		x_0 切片	-0.278	β_0	0.119	0.140	-0.106	-0.060	-0.046	-0.052	0.041	0.024	
		x_1 A1	2.222	β_1	0.140	0.637	-0.152	-0.199	-0.052	-0.241	0.056	0.074	
		x_2 A2	0.201	β_2	-0.106	-0.152	0.144	0.048	0.041	0.056	-0.054	-0.020	
		x_3 A3	-1.185	β_3	-0.060	-0.199	0.048	0.237	0.024	0.074	-0.020	-0.089	
		x_4 W体重	0.546	β_4	-0.046	-0.052	0.041	0.024	0.018	0.019	-0.017	-0.010	
		x_5 A1×W	-0.752	β_5	-0.052	-0.241	0.056	0.074	0.019	0.093	-0.021	-0.028	
		x_6 A2×W	-0.065	β_6	0.041	0.056	-0.054	-0.020	-0.017	-0.021	0.021	0.008	
		x_7 A3×W	0.382	β_7	0.024	0.074	-0.020	-0.089	-0.010	-0.028	0.008	0.035	
色	x_0	x_1	x_2	x_3	x_4	x_5	x_6	x_7	$\ln y^\wedge$	$V(\ln y^\wedge)$	y^\wedge	$L\ 95\%$	$U\ 95\%$
1	1	<i>I</i>	<i>0</i>	<i>0</i>	<i>2.5</i>	<i>2.5</i>	<i>0</i>	<i>0</i>	<i>1.431</i>	<i>0.022</i>	<i>4.181</i>	<i>3.128</i>	<i>5.589</i>
2	1	0	1	0	2.5	0	2.5	0	1.127	0.004	<i>3.087</i>	2.746	3.471
3	1	0	0	1	2.5	0	0	2.5	0.857	0.010	<i>2.357</i>	1.930	2.879
4	1	-1	-1	-1	2.5	-2.5	-2.5	-2.5	0.936	0.024	<i>2.550</i>	1.887	3.447

甲羅の色を名義尺度として扱ったので, それぞれの推定値に 95%信頼区間を付け加えた. JMP で作成した予測プロファイルと同じ結果となっていることが確認される.

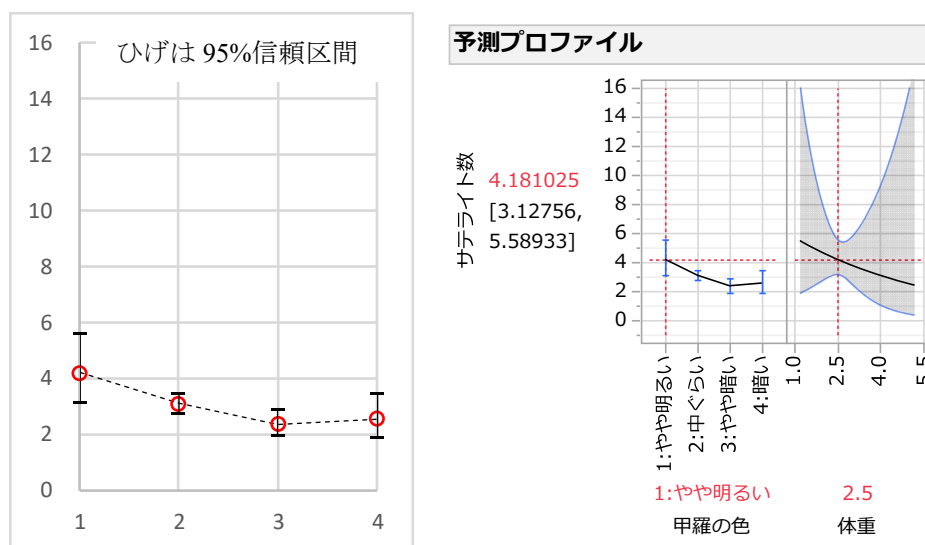


図 10 体重が 2.5 kg の場合の甲羅の色についての予測プロファイル

甲羅の色を1:やや明るいに固定し、体重を 1, 2, 3, 4, 5 kg と変化させた場合について、表 16 に計算結果を示す。甲羅の色を順次変えた場合について、表 17 および表 18 に示す。

表 16 甲羅の色が 1:やや明るい場合の場合の体重についての予測値

色	x_0	x_1	x_2	x_3	x_4	x_5	x_6	x_7	$\ln y^\wedge$	$V(\ln y^\wedge)$	y^\wedge	$L95\%$	$U95\%$
1	1	1	0	0	1	1	0	0	1.739	0.405	5.690	1.635	19.803
	1	1	0	0	2	2	0	0	1.533	0.075	4.633	2.713	7.913
	1	1	0	0	3	3	0	0	1.328	0.044	3.773	2.498	5.699
	1	1	0	0	4	4	0	0	1.122	0.314	3.072	1.024	9.212
	1	1	0	0	5	5	0	0	0.917	0.883	2.501	0.396	15.786

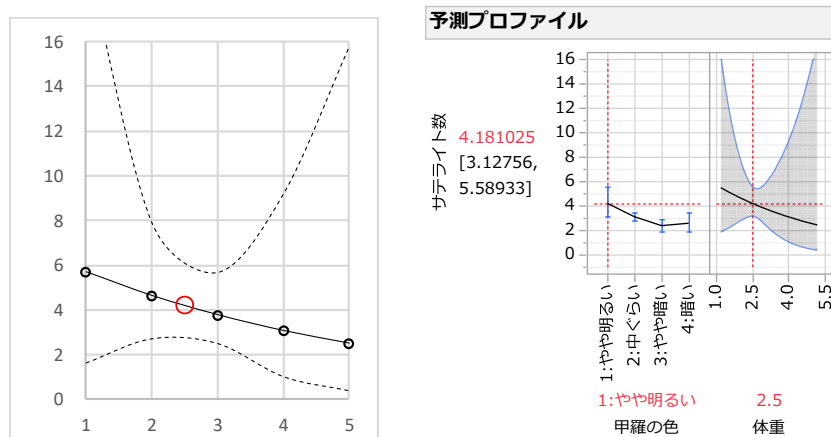


図 11 甲羅の色が 1:やや明るい場合の場合の体重についての予測プロファイル

表 17 甲羅の色が 2:中ぐらい場合の場合の体重についての予測値

色	x_0	x_1	x_2	x_3	x_4	x_5	x_6	x_7	$\ln y^\wedge$	$V(\ln y^\wedge)$	y^\wedge	$L95\%$	$U95\%$
2	1	0	1	0	1	0	1	0	0.405	0.022	1.499	1.119	2.008
	1	0	1	0	2	0	2	0	0.887	0.007	2.427	2.068	2.848
	1	0	1	0	3	0	3	0	1.368	0.004	3.928	3.492	4.419
	1	0	1	0	4	0	4	0	1.850	0.013	6.359	5.082	7.957
	1	0	1	0	5	0	5	0	2.332	0.035	10.293	7.130	14.861

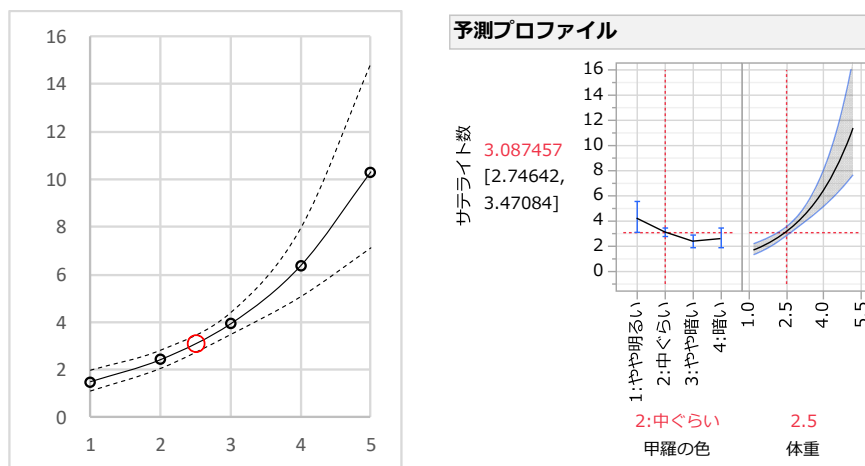


図 12 甲羅の色が 2:中ぐらい場合の場合の体重についての予測プロファイル

表 18 甲羅の色が 3:やや暗い, 4:暗い場合の体重についての予測値

色	x_0	x_1	x_2	x_3	x_4	x_5	x_6	x_7	$\ln y^\wedge$	$V(\ln y^\wedge)$	y^\wedge	$L95\%$	$U95\%$
3	1	0	0	1	1	0	0	1	-0.535	0.095	0.586	0.320	1.070
	1	0	0	1	2	0	0	2	0.393	0.022	1.482	1.111	1.976
	1	0	0	1	3	0	0	3	1.322	0.016	3.750	2.922	4.811
	1	0	0	1	4	0	0	4	2.250	0.078	9.487	5.480	16.424
	1	0	0	1	5	0	0	5	3.178	0.208	24.005	9.813	58.719
4	1	-1	-1	-1	1	-1	-1	-1	-0.535	0.211	0.586	0.238	1.442
	1	-1	-1	-1	2	-2	-2	-2	0.446	0.037	1.562	1.073	2.272
	1	-1	-1	-1	3	-3	-3	-3	1.426	0.060	4.164	2.575	6.734
	1	-1	-1	-1	4	-4	-4	-4	2.407	0.282	11.101	3.922	31.420
	1	-1	-1	-1	5	-5	-5	-5	3.388	0.701	29.598	5.733	152.81

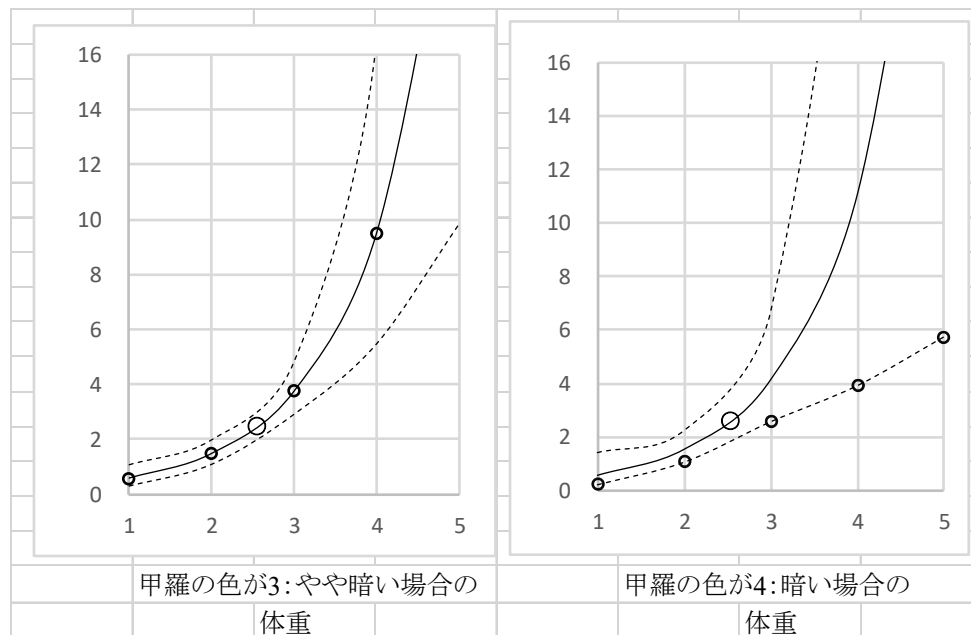


図 13 甲羅の色が 3:やや暗い, 4:暗い場合の体重についての予測プロフィール

13. 甲羅の幅に対する信頼区間と予測区間

対数リンクのポアソン回帰に対する信頼区間は、予測値(平均)に対する 95%信頼区間を意味し、予測区間は、個別データの 95%信頼区間を意味する。第 5 節では、信頼区間と予測区間は、JMP の「列の保存」の機能を使い、JMP ファイルに出力された結果を「重ね合わせプロット」で整形する方法を示した。対数リンクのポアソン回帰における信頼区間と予測区間の計算方法については、説明をしておかなかった。

そこで、JMP で求められた推定値と共分散行列を用いて、Excel による信頼区間と予測区間を計算し、グラフ化する方法を示す。これまで、対数リンクのポアソン回帰については JMP に丸投げして推定値および共分散行列を得ていた。

第 4 節で、カブトガニの甲羅の幅を説明変数、サテライト数を応答変数とする対数リンクのポアソン回帰を行い、スライド 7 に散布図上に信頼区間と予測区間を上書きした結果を示した。表 19 に解析に用いたデータリストの一部を示す。

表 19 甲羅の幅とサテライト数のデータリスト(付表 A およびスライド 4 抜粋)

	甲羅の幅	サテライト
No	x	y
1	28.3	8
2	22.5	0
3	26.0	9
4	24.8	0
5	26.0	4
:		
170	29.0	4
171	28.0	0
172	27.0	0
173	24.5	0

表 20 に対数リンクによるポアソン回帰のパラメータ推定値および推定値の共分散行列を示す。これらの推定値があれば、信頼区間と予測区間を求めることができる。さて、推定値および共分散行列は、どのようにして計算されるのだろうか。表 20 の左のパラメータの推定値は表 2 の再掲であるが、表 20 の右の推定値の共分散は、JMP で改めて求めた結果である。

表 20 甲羅の幅に対する対数リンクのポアソン回帰の結果

パラメータ推定値			推定値の共分散		
項	推定値	標準誤差	共分散		
切片	-3.304757	0.542242	切片	0.294026	甲羅の幅
甲羅の幅	0.164045	0.019965	甲羅の幅	-0.010790	0.000399

表 21 に示すのは, JMP で得られた推定値を用いて信頼区間と予測区間を計算した結果である. 表の左側は, 表 19 で示した 173 個のデータリストである. 表の右に, 切片としての x_0 , 甲羅の幅 x_1 を与え, $\ln(\hat{y})$, $Var[\ln(\hat{y})]$ の計算を行っている.

パラメータの推定値 $\hat{\beta}$ は,

$$\hat{\beta} = [-3.3048, 0.1640]^T$$

であり, 切片項 $x_0 = 1$, 甲羅の幅 x_1 とする横ベクトル \mathbf{x}

$$\mathbf{x} = [1, x_1]$$

としたときに, 推定値 $\ln(\hat{y})$ は,

$$\ln(\hat{y}) = \mathbf{x}\hat{\beta}$$

として推定される. 予測値の分散 $Var(\ln \hat{y})$ は, 共分散行列 Σ を挟んだ \mathbf{x} の 2 次形式

$$Var(\ln \hat{y}) = \mathbf{x} \Sigma \mathbf{x}^T$$

で求められる. 95%信頼区間は, $\ln \hat{y}$ について正規分布の両側 5%点である 1.96 を用いて計算し, 指数を取り

$$(U95\%, L95\%) = \exp\left[\ln \hat{y} \pm 1.96\sqrt{Var(\ln \hat{y})}\right]$$

で求めている.

JMP では, 予測区間(個別データの 95%信頼区間)は, 推定値 \hat{y} に対するポアソン分布の 2.5%点および 97.5%点より推定している. Excel には, ポアソン回帰の確率関数および累積分布関数があるが, パーセント点を与えて分位点を求める関数がないので, 推定値 \hat{y} をポアソン回帰の母数 μ として与え JMP の Poisson Quantile() で分位点を求め Excel シートに戻している.

ポアソン回帰における JMP の予測区間の計算方法は, 良く知られている通常の回帰分析とは異なる. 通常の回帰分析では, 甲羅の幅に対して, 回帰直線の推定値に対する分散に個別データの分散を加えた合成分散を使っている. 対数リンクのポアソン回帰では, 個別データの分散は, $\exp(\hat{y}) = \hat{y}$ と甲羅の幅のサイズによって変化し, 推定値 \hat{y} の信頼区間の推定で用いた分散は対数で設定されていて, それらを単純に加えて合成分散することができないためである.

表 21 甲羅の幅に対する対数リンクのポアソン回帰の信頼区間と予測区間

					$\beta^0_0 = -3.3048$	$\Sigma =$	0.2940	-0.0108			
					$\beta^1_1 = 0.1640$		-0.0108	0.0004			
	甲羅の幅	サテライト						推定値		個別データ	
No	x	y	x_0	x_1	$\ln y^{\wedge}$	$Var(\ln y^{\wedge})$	y^{\wedge}	L 95%	U 95%	L 95%	U 95%
1	28.3	8	1	20.5	0.058	0.019	1.060	0.808	1.390	0	3
2	22.5	0	1	21.0	0.140	0.017	1.150	0.893	1.482	0	4
3	26.0	9	1	22.0	0.304	0.012	1.356	1.092	1.684	0	4
4	24.8	0	1	23.0	0.468	0.009	1.597	1.332	1.915	0	4
5	26.0	4	1	24.0	0.632	0.006	1.882	1.622	2.183	0	5
6	23.8	0	1	25.0	0.796	0.004	2.217	1.969	2.498	0	6
7	26.5	0	1	26.0	0.960	0.002	2.613	2.372	2.878	0	6
8	24.7	0	1	27.0	1.124	0.002	3.079	2.821	3.360	0	7
9	23.7	0	1	28.0	1.289	0.002	3.627	3.300	3.987	0	8
10	25.6	0	1	29.0	1.453	0.003	4.274	3.808	4.797	1	9
11	24.3	0	1	30.0	1.617	0.005	5.036	4.360	5.817	1	10
12	25.8	0	1	31.0	1.781	0.008	5.934	4.971	7.082	2	11
13	28.2	11	1	32.0	1.945	0.012	6.991	5.657	8.641	2	13
14	21.0	0	1	33.0	2.109	0.016	8.238	6.428	10.557	3	14
15	26.0	14	1	33.5	2.191	0.018	8.942	6.850	11.672	4	15
:											
170	29.0	4									
171	28.0	0									
172	27.0	0									
173	24.5	0									

これらの計算結果を Excel の「散布図(x, y)」を使って描いたのが図 14 である。

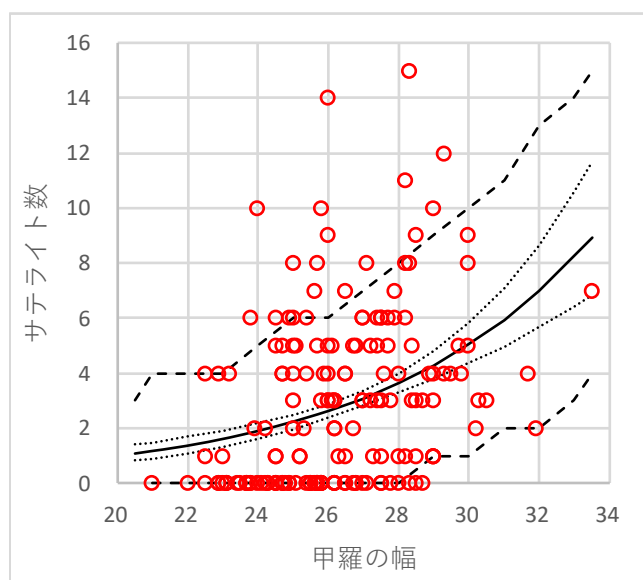
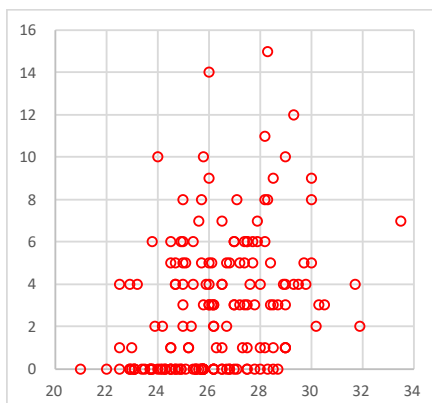
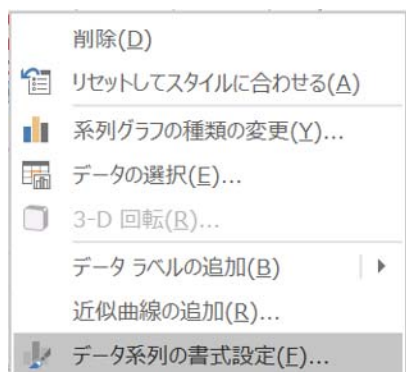


図 14 甲羅の幅に対する信頼区間と予測区間

図 14 は、「散布図(x, y)」の機能を使って整形したものである。以下に整形方法を示す。まず、甲羅の幅とサテライト数のデータを選択して点をプロットする。デフォルトは●印なので○印への変更する。

マーカーの変更は、●印を選択して「データ系列の書式設定」で行う。 軸の目盛りは、変更したい軸を選択し「軸の書式設定」で行う。



データ系列の書式設定

系列のオプション ▼

線 マーカー

▲ マーカーのオプション

☐ 自動(U)

☐ なし(Q)

☒ 組み込み

種類 ○

サイズ 5

▲ 塗りつぶし

☒ 塗りつぶしなし(N)

☐ 塗りつぶし (単色)(S)

☐ 塗りつぶし (グラデーション)(G)

☐ 塗りつぶし (図またはテキストチャ)(P)

☐ 塗りつぶし (パターン)(A)

☐ 自動(U)

☐ 要素を塗り分ける(V)

▲ 枠線

☐ 線なし(N)

☒ 線 (単色)(S)

☐ 線 (グラデーション)(G)

☐ 自動(U)

色(C) 色

透明度(I) 0%

幅(W) 0.75 pt

軸の書式設定

軸のオプション ▼ 文字のオプション

軸のオプション

境界値

最小値(N) 20.0

最大値(X) 34.0

単位

主(I) 2.0

補助(I) 1.0

▲ 表示形式

カテゴリ(C)

数値

小数点以下の桁数(D): 0

☐ 桁区切り (,) を使用する(U)

負の数の表示形式(N):

(1234)

(1234)

1234

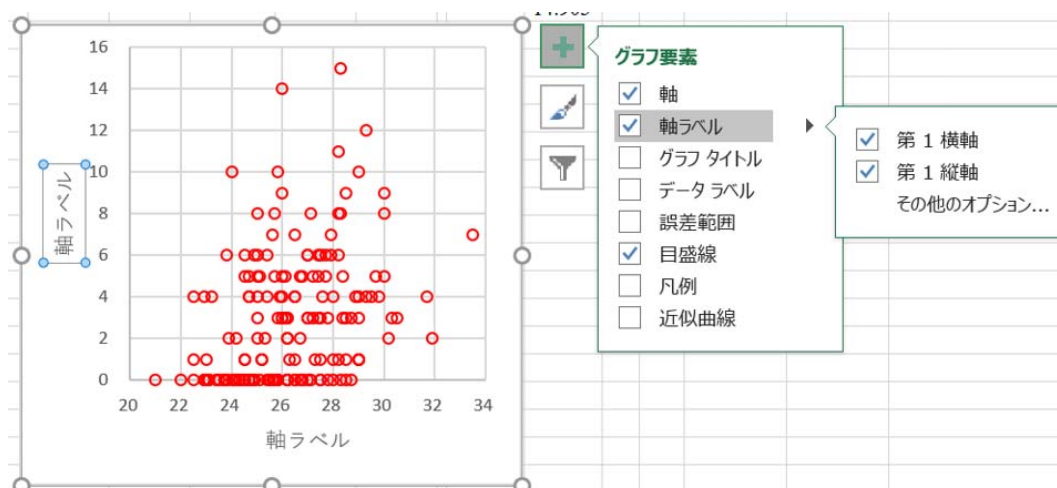
-1234

-1234

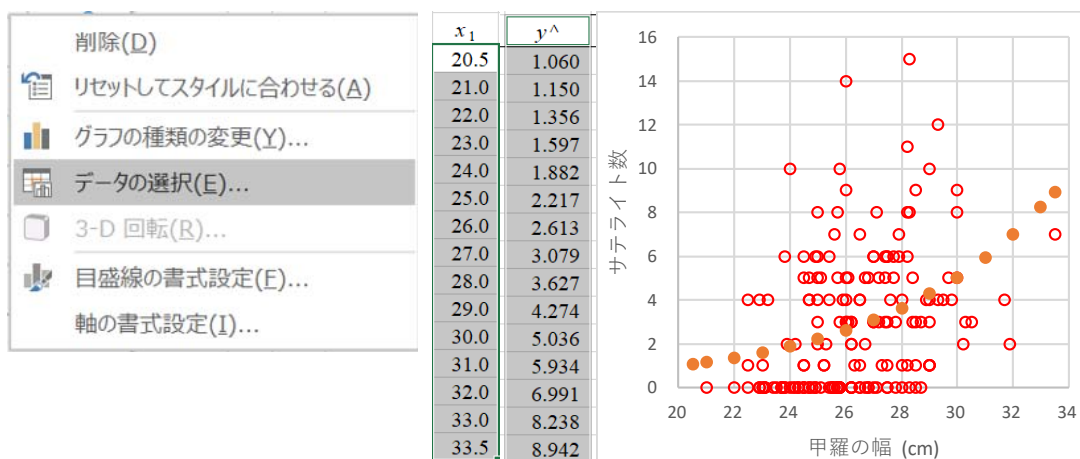
△ 1234

▲ 1234

タイトルおよび軸ラベルなどは、「グラフ要素」を使って適宜与える。



推定値, 信頼区間, 予測値の上書きは, 「データの選択」を使う。「データソースの選択」, 「追加」で「系列名」に「推定値」を入力, 「系列 X の値」で X の範囲を選択, 「系列 Y の値」で Y の範囲を選択する。



データソースの選択

グラフデータの範囲(D):

データ範囲が複雑すぎるため、表示できません。データ範囲を選択し直すと、[系列] タブのすべての系列が置き換えられます。

行/列の切り替え(W)

凡例項目 (系列)(S)

追加(A) 編集(E) 削除(R)

☒ y

☒ 推定値

非表示および空白のセル(H)

横 (項目)

21

21

22

23

24

系列の編集

系列名(N):

推定値 = 推定値

系列 X の値(X):

= 'カブト グラフ'!\$G\$11:\$G\$25 = 20.5 , 21.0 , ...

系列 Y の値(Y):

= 'カブト グラフ'!\$J\$11:\$J\$25 = 1.059892033, 1

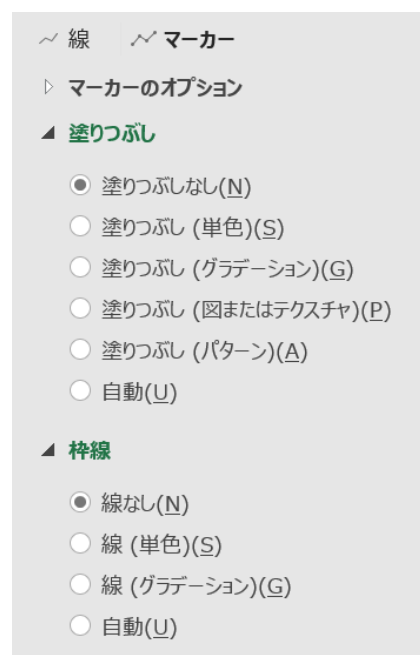
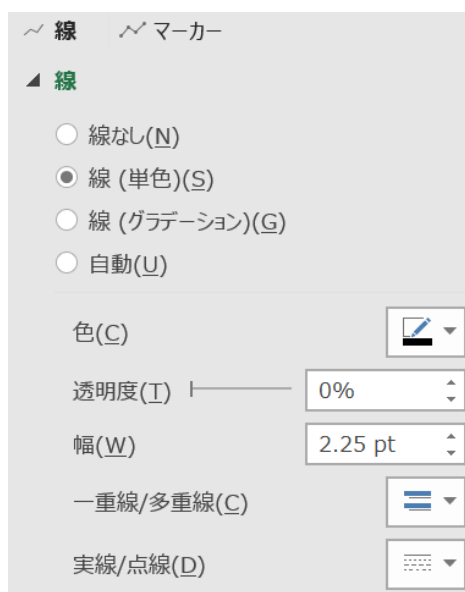
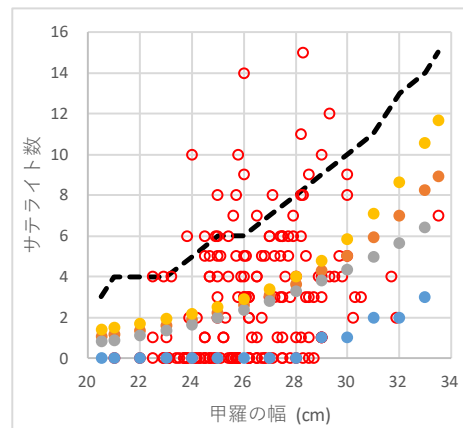
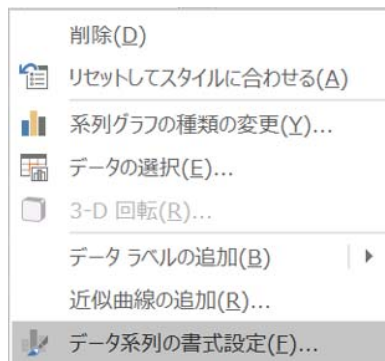
OK キャンセル

さらに、

推定値		個別データ	
L95%	U95%	L95%	U95%

についても順次「追加」し、「系列名」、「系列 X の値」、

「系列 Y の値」でデータの選択を行う。デフォルトは、●印などであるので、線で結び、●印を消して、線種、線の幅などを適宜変更する。



統計解析の結果は、何等かの図を作成し、その図を用いて解釈をすることが相互理解の助けになり、新しい知見を得ることになる。統計ソフトでカバーできる範囲は、できる限り活用することが望ましいが、努力の割には見栄えが悪いばあいには、Excel の図などを活用すること薦める。特にここに示した「散布図(x, y)」は、多機能であるが故に使いづらいことも確かであるが、データが代わっても、それに対応して自動的に更新されることは、探索的な統計解析の結果の表示に優れている。

14. 最尤法による対数リンクのポアソン回帰

JMP では、最尤法により対数リンクのポアソン回帰を行っているが、どのような計算を行っているのだろうか。紋きり型な言い方、「対数リンクのポアソン回帰の対数尤度を最大にするように回帰パラメータを求める方法であり、統計ソフトで簡単に求められる」が蔓延している。Excel は、現代の算盤であり「最尤法による対数リンクのポアソン回帰」を Excel の基本機能のみで実現できる。説明変数を甲羅の幅を x_i 、反応変数をサテライト数を y_i 、データ数は 173 ツガイとする。

対数リンク

変数 x_i が増大するにつれ、観測データ y_i が指数関数的に増加することもしばしば経験する。位置パラメータ μ_i （推定値 \hat{y}_i ）も指数関数的に増加すると仮定すると、回帰式 $\beta_0 + \beta_1 x_i$ は、

$$\mu_i = \exp(\beta_0 + \beta_1 x_i) + \varepsilon_i, \quad i=1, 2, \dots, n \quad (13.1)$$

と定義され、誤差 ε_i の分布がポアソン分布に従うとしたときに分布の確率 p_i は、

$$\begin{aligned} p_i &= \frac{\mu_i^{y_i} e^{-\mu_i}}{y_i!} \\ &= \frac{\exp(\beta_0 + \beta_1 x_i)^{y_i} e^{-\exp(\beta_0 + \beta_1 x_i)}}{y_i!} \end{aligned} \quad (13.2)$$

である。両辺に対数をとって、

$$\ln(\mu_i) = \beta_0 + \beta_1 x_i, \quad i=1, 2, \dots, n \quad (13.3)$$

のように線形化する。一般化線形モデルでは、この線形化する変換をリンク関数と言う。一般化線形モデルに対する統計ソフトでは、元の指数関数ではなく、線形化する関数名を使っている。この場合のリンク関数は「対数」である。

対数リンクの場合の偏微分式

ニュートン・ラフソン法による最尤法は、指数関数に対して直接計算するので、対数変換した式を使う必要はない。

式 (13.2) の p_i に対する対数尤度 $\ln L_i$ は、

$$\begin{aligned} \ln L_i &= \ln \left[\frac{\exp(\beta_0 + \beta_1 x_i)^{y_i} e^{-\exp(\beta_0 + \beta_1 x_i)}}{y_i!} \right] \\ &= y_i(\beta_0 + \beta_1 x_i) - \exp(\beta_0 + \beta_1 x_i) - \ln(y_i!) \end{aligned} \quad (13.4)$$

であり、パラメータ β_0 および β_1 で偏微分すると

$$U_{1i} = \frac{\partial \ln L_i}{\partial \beta_0} = y_i - \exp(\beta_0 + \beta_1 x_i) = y_i - \mu_i$$

$$U_{2i} = \frac{\partial \ln L_i}{\partial \beta_1} = y_i x_i - \exp(\beta_0 + \beta_1 x_i) x_i = (y_i - \mu_i) x_i \quad (13.5)$$

となり，さらに β_0 および β_1 で偏微分すると

$$\begin{aligned} H_{1,1,i} &= \frac{\partial^2 \ln L_i}{\partial \beta_0 \partial \beta_0} = -\exp(\beta_0 + \beta_1 x_i) = -\mu_i \\ H_{1,2,i} &= \frac{\partial^2 \ln L_i}{\partial \beta_0 \partial \beta_1} = \frac{\partial \ln L_i}{\partial \beta_0} x_i = -\mu_i x_i \\ H_{2,1,i} &= \frac{\partial^2 \ln L_i}{\partial \beta_1 \partial \beta_0} = \frac{\partial \ln L_i}{\partial \beta_0} x_i = -\mu_i x_i \\ H_{2,2,i} &= \frac{\partial^2 \ln L_i}{\partial \beta_1 \partial \beta_1} = \frac{\partial \ln L_i}{\partial \beta_0} x_i^2 = -\mu_i x_i^2 \end{aligned} \quad (13.6)$$

となる．これらを i について加えスコアベクトル U とヘッセ行列 H にまとめる．

$$U = \begin{bmatrix} \sum_i U_{1i} \\ \sum_i U_{2i} \end{bmatrix} = \begin{bmatrix} \sum_i (y_i - \mu_i) \\ \sum_i (y_i - \mu_i) x_i \end{bmatrix} \quad (13.7)$$

$$H = \begin{bmatrix} \sum_i H_{1,1,i} & \sum_i H_{1,2,i} \\ \sum_i H_{2,1,i} & \sum_i H_{2,2,i} \end{bmatrix} = \begin{bmatrix} \sum_i -\mu_i & \sum_i -\mu_i x_i \\ \sum_i -\mu_i x_i & \sum_i -\mu_i x_i^2 \end{bmatrix} \quad (13.8)$$

反復計算

ニュートン・ラフソン法は，パラメータ $\beta = [\beta_0, \beta_1]^T$ の最初の初期値を $\hat{\beta}^{(0)}$ としたときに， $U^{(0)}$ と $H^{(0)}$ を計算し

$$\hat{\beta}^{(1)} = \hat{\beta}^{(0)} + (-H^{(0)})^{-1} U^{(0)} \quad (13.9)$$

$\hat{\beta}^{(1)}$ を $\hat{\beta}^{(1)}$ に代入し， $U^{(1)}$ と $H^{(1)}$ を計算し直し，

$$\hat{\beta}^{(2)} = \hat{\beta}^{(1)} + (-H^{(1)})^{-1} U^{(1)} \quad (13.10)$$

$\hat{\beta}^{(2)}$ を改めて計算する．これを，繰り返し計算して，対数尤度 $\ln L$ を最大化する方法である．最大化したとの判断は， $\hat{\beta}^{(m)}$ と $\hat{\beta}^{(m+1)}$ の差が 10^{-6} 以下になった場合など適宜設定する．

気が遠くなるような計算であるが，現在の Excel は，このくらいの計算は瞬時である．私も昔は，行列計算言語（SAS の IML）などを使っていたが，計算結果の可視化似に優れた Excel を使って解説書を作成すようになった．

Excel による計算

表 22 に Excel によるニュートン・ラフソン法で，対数尤度を最大化する方法を示す．初期値として， $\hat{\beta}_0 = -3.00$ ， $\hat{\beta}_1 = 0.20$ を与えた反復で，第 6 反復で収束した結果である．実際の反復計算は，表 22 の「元のパラメータ」から得られた「新たなパラメータ」をコピーし，値のみを元のパラメータにペーストする．変化量が共に 0 になったら収束したと判断する．

	元の(m-1)	変化量	新たな(m)
	パラメータ	$(-H)^{-1}U$	パラメータ
$\beta^0=$	-3.0000	-0.4539	-3.4539
$\beta^1=$	0.2000	-0.0099	0.1901

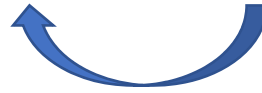


表 22 ニュートン・ラフソン法によるポアソン回帰のあてはめ

			元の(m-1)	変化量	新たな(m)	1階の	2 階の偏微分		負の逆行列	
			パラメータ	$(-H)^{-1}U$	パラメータ	偏微分U	H		$(-H)^{-1}$	
	$\beta^0=$		-3.0000	-0.4539	-3.4539	-1313.50	-1818.5	-49551.4	0.080561	-0.002936
	$\beta^1=$		0.2000	-0.0099	0.1901	-35882.33	-49551.4	-1359484	-0.002936	0.000108
				計	-1129.71	-1313.50	-35882.33	-1818.50	-49551.43	-1359484
i	x	y	μ^{\wedge}	p	$\ln Li$	$\partial \beta_0$	$\partial \beta_1$	$\partial \beta_0 \partial \beta_0$	$\partial \beta_0 \partial \beta_1$	$\partial \beta_1 \partial \beta_1$
1	28.3	8	14.2963	0.0268	-3.6209	-6.2963	-178.1850	-14.2963	-404.5850	-11449.75
2	22.5	0	4.4817	0.0113	-4.4817	-4.4817	-100.8380	-4.4817	-100.8380	-2268.86
3	26.0	9	9.0250	0.1318	-2.0268	-0.0250	-0.6504	-9.0250	-234.6504	-6100.91
4	24.8	0	7.0993	0.0008	-7.0993	-7.0993	-176.0633	-7.0993	-176.0633	-4366.37
5	26.0	4	9.0250	0.0333	-3.4031	-5.0250	-130.6504	-9.0250	-234.6504	-6100.91
6	23.8	0	5.8124	0.0030	-5.8124	-5.8124	-138.3360	-5.8124	-138.3360	-3292.40
7	26.5	0	9.9742	0.0000	-9.9742	-9.9742	-264.3158	-9.9742	-264.3158	-7004.37
8	24.7	0	6.9588	0.0010	-6.9588	-6.9588	-171.8811	-6.9588	-171.8811	-4245.46
9	23.7	0	5.6973	0.0034	-5.6973	-5.6973	-135.0270	-5.6973	-135.0270	-3200.14
10	25.6	0	8.3311	0.0002	-8.3311	-8.3311	-213.2771	-8.3311	-213.2771	-5459.89
:										
170	29.0	4	16.4446	0.0002	-8.4227	-12.4446	-360.8948	-16.4446	-476.8948	-13829.95
171	28.0	0	13.4637	0.0000	-13.4637	-13.4637	-376.9847	-13.4637	-376.9847	-10555.57
172	27.0	0	11.0232	0.0000	-11.0232	-11.0232	-297.6258	-11.0232	-297.6258	-8035.90
173	24.5	0	6.6859	0.0012	-6.6859	-6.6859	-163.8044	-6.6859	-163.8044	-4013.21

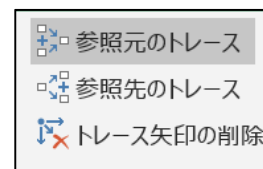
表 23 ニュートン・ラフソン法の反復

		元の(m)	変化量	新たな(m+1)	共分散行列	
反復		パラメータ	$(-H)^{-1}U$	パラメータ	負の逆行列 $(-H)^{-1}$	
1	$\beta^0=$	-3.0000	-0.4539	-3.4539	0.080561	-0.002936
	$\beta^1=$	0.2000	-0.0099	0.1901	-0.002936	0.000108
2	$\beta^0=$	-3.4539	-0.0264	-3.4802		
	$\beta^1=$	0.1901	-0.0148	0.1754	:	
3	$\beta^0=$	-3.4802	0.1451	-3.3351	:	
	$\beta^1=$	0.1754	-0.0099	0.1655	:	
4	$\beta^0=$	-3.3351	0.0301	-3.3051	:	
	$\beta^1=$	0.1655	-0.0014	0.1641	:	
5	$\beta^0=$	-3.3051	0.0003	-3.3048	:	
	$\beta^1=$	0.1641	0.0000	0.1640		
6	$\beta^0=$	-3.3048	0.0000	-3.3048	0.294026	-0.010790
	$\beta^1=$	0.1640	0.0000	0.1640	-0.010790	0.000399

反復が 6 の $(-H)^{-1}$ が表 20 で示した共分散行列に一致する。

Excel での計算は、計算式がセル裏に隠れていて見えないことが利点でもあり、最大の欠点でもある。この Excel のシートの結果が正しいのかは、統計ソフトの結果と常に照合する必要がある。新たなデータで計算した場合に、何らかのミス操作で計算式が壊れ正しい結果が得られないことは、しばしば経験する。

利点は、実際の計算シートを用いれば、どのような計算をしているのか、確認することが容易な点である。Excel の計算式は、セルの位置で示されているので難解であるので、「参照元のトレース」、「参照先のトレース」を使うことにより可視化するとよい。



1 階の偏微分 U の参照元			共分散行列 (負の逆行列) 参照元			
1階の 偏微分 U	2 階の偏微分 H		2 階の偏微分 H	負の逆行列 $(-H)^{-1}$		
0.00	-505.0	-13669.1	-505.0	-13669.1	0.294026	0.010790
0.00	-13669.1	-372497	-13669.1	-372497	-0.010790	0.000399
0.00	0.00	-505.00	0.00	-505.00	-3669.10	-372497
$\partial \beta_0$	$\partial \beta_1$	$\partial \beta_0 \partial \beta_0$	$\partial \beta_1$	$\partial \beta_0 \partial \beta_0$	$\partial \beta_0 \partial \beta_1$	$\partial \beta_1 \partial \beta_1$
4.1897	118.5673	-3.8103	118.5673	-3.8103	-107.8327	3051.66
-1.4715	-33.1079	-1.4715	-33.1079	-1.4715	-33.1079	-744.93
6.3872	166.0677	-2.6128	166.0677	-2.6128	-67.9323	-1766.24

変化量 の参照元 $(-H)^{-1}U$									
	元の(m-1) パラメータ	変化量 $(-H)^{-1}U$	新たな(m) パラメータ	1階の 偏微分 U	2 階の偏微分 H	負の逆行列 $(-H)^{-1}$			
$\beta^0_0 =$	-3.0000	-0.4539	-3.4539	1313.50	-1818.5	-49551.4	0.080561	0.002936	
$\beta^0_1 =$	0.2000	-0.0099	0.1901	-35882.33	-49551.4	-1359484	-0.002936	0.000108	
		計	-1129.71	1313.50	-35882.33	-1818.50	-49551.43	-1359484	
x	y	μ^{\wedge}	p	$\ln Li$	$\partial \beta_0$	$\partial \beta_1$	$\partial \beta_0 \partial \beta_0$	$\partial \beta_0 \partial \beta_1$	$\partial \beta_1 \partial \beta_1$

対数尤度 $\sum_i \ln L_i$ の参照元				
	元の(m-1) パラメータ	変化量 $(-H)^{-1}U$	新たな(m) パラメータ	
$\beta^0_0 =$	3.0000	-0.4539	-3.4539	
$\beta^0_1 =$	0.2000	-0.0099	0.1901	
		計	-1129.71	
x	y	μ^{\wedge}	p	$\ln Li$
28.3	8	4.2963	0.0268	3.6209
22.5	0	4.4817	0.0113	-4.4817
26.0	9	9.0250	0.1318	-2.0268

$\sum_i \ln L_i$ の参照元は、 $\ln L_i \leftarrow$ ポアソン確率 $p \leftarrow (y, \hat{\mu}), \hat{\mu} \leftarrow (\beta_0, \beta_1, x)$ とたどることができる。

対数尤度 $\sum_i \ln L_i$ は、計算されているだけで、最尤法の反復計算では使われていない。結果として最大になっていると期待されているだけである。もちろん対数尤度 $\sum_i \ln L_i$ の偏微分式が使われているのであるが、対数尤度 $\sum_i \ln L_i$ そのものを使い、Excel のソルバーにより、

「 $\sum_i \ln L_i$ を最大にするように (β_0, β_1) を変化させよ」

とすることにより、手作業での反復計算なしに、直接「反復 6」の結果を得ることができる。このことから、間接的に「最大」であることを実感できる。この場合は、偏微分式 U 、 H および $(-H)^{-1}U$ などとも必要としない。パラメータが直接推定することができる。ただし、共分散行列の推定ができないことになる。

表 24 ソルバーで $\sum_i \ln L_i$ を最大化するように (β_0, β_1) を変化させた結果

			元の(m-1)	変化量	新たな(m)	1階の	2 階の偏微分		負の逆行列	
			パラメータ	$(-H)^{-1}U$	パラメータ	偏微分U	H		$(-H)^{-1}$	
		$\beta^0=$	-3.3048	0.0000	-3.3048	0.00	-505.0	-13669.1	0.294027	-0.010790
		$\beta^1=$	0.1640	0.0000	0.1640	0.04	-13669.1	-372496	-0.010790	0.000399
			計		-461.59	0.00	0.04	-505.00	-13669.06	-372496
i	x	y	μ^{\wedge}	p	$\ln Li$	$\partial \beta_0$	$\partial \beta_1$	$\partial \beta_0 \partial \beta_0$	$\partial \beta_0 \partial \beta_1$	$\partial \beta_1 \partial \beta_1$
1	28.3	8	3.8103	0.0244	-3.7132	4.1897	118.5677	-3.8103	-107.8323	-3051.65
2	22.5	0	1.4715	0.2296	-1.4715	-1.4715	-33.1078	-1.4715	-33.1078	-744.93
3	26.0	9	2.6128	0.0011	-6.7709	6.3872	166.0679	-2.6128	-67.9321	-1766.23
:										
170	29.0	4	4.2740	0.1936	-1.6419	-0.2740	-7.9456	-4.2740	-123.9456	-3594.42
171	28.0	0	3.6273	0.0266	-3.6273	-3.6273	-101.5658	-3.6273	-101.5658	-2843.84
172	27.0	0	3.0785	0.0460	-3.0785	-3.0785	-83.1207	-3.0785	-83.1207	-2244.26
173	24.5	0	2.0428	0.1297	-2.0428	-2.0428	-50.0497	-2.0428	-50.0497	-1226.22

$$p_i = \frac{\mu_i^{y_i} e^{-\mu_i}}{y_i!}, \quad U = \begin{bmatrix} \sum_i (y_i - \mu_i) \\ \sum_i (y_i - \mu_i) x_i \end{bmatrix}, \quad H = \begin{bmatrix} \sum_i -\mu_i & \sum_i -\mu_i x_i \\ \sum_i -\mu_i x_i & \sum_i -\mu_i x_i^2 \end{bmatrix}, \quad \hat{\beta}^{(1)} = \hat{\beta}^{(0)} + (-H^{(0)})^{-1} U^{(0)}$$

15. まとめ

一般化線形モデルで定式化されているポアソン回帰については、ドブソン著、田中・森川・山中ら訳(2008)に丁寧な解説があり必読の書である。反復重み付き回帰による最尤法によるポアソン回帰について、人工データを用いた丁寧な導入があり、これに基づき高橋(2019)で Excel を用いたポアソン回帰による勾配比検定を行うことができた。

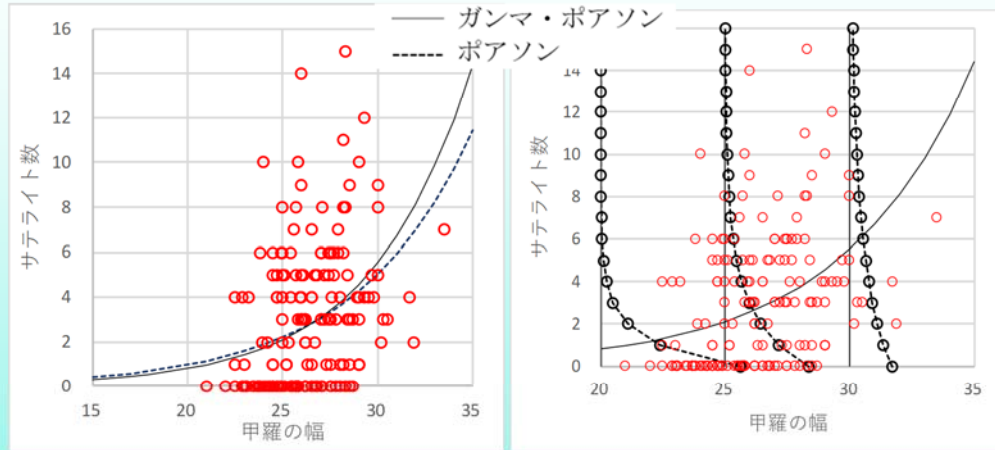
通常の回帰分析に対応したポアソン回帰は、恒等リンクとした場合であるが、各種の応用例で散見するのは、対数リンクでオフセットがあり、あるいは過分散を考慮する場合など多彩である。JMP の一般化線形モデル、SAS の GENMOD プロシジャでポアソン回帰ができるようになっているが、適当な参考文献は見当たらない。ドブソンの訳本でも、ポアソン回帰の例示は、対数リンクでオフセットがあり、2乗項もあり、さらに交互作用も含めた事例が示されている。もちろん統計ソフトを使う前提で、追試も容易にできるのだが、その結果の解釈は難解である。

アグレスティの訳本で例示されているカブトガニのサテライト数の事例は、通常の回帰分析と対比しやすいので探索的なポアソン回帰の例示として取り上げた。全データに対するサテライト数の分布について、ゼロ過剰ポアソン分布よりも、さらにゼロ過剰ガンマ・ポアソン分布のあてはめが良好であったが、それらの分布を用いた回帰分析には難点がある。これは、図 2 にも示したように、甲羅の幅が大きい時にはサテライト数のゼロが存在しなくなるので、ゼロ過剰ガンマ・ポアソン分布を仮定して回帰分析を行うと、体重が重い場合にも過剰なゼロが存在を仮定することになり、現実のデータとの乖離を無視できなくなるためである。

スライド 28 左は、Excel で最尤法による対数リンクのポアソン回帰およびガンマ・ポアソン回帰の指数曲線を描いたもので、スライド 28 右は、甲羅の幅が(20, 25, 30)の場合についてガンマ・ポアソン分布を上書きしたもので、甲羅の幅が 30 cm の場合に実際の観察データでサテライト数が 0 となるカブトガニがないにもかかわらず、分布関数では大きな確率となっている。

スライド 29 は、ガンマ・ポアソン回帰とゼロ過剰ガンマ・ポアソン回帰を比較したものであり、甲羅の幅が 30 cm の場合に実際の観察データでサテライト数が 0 となるカブトガニがないにもかかわらず、分布関数では大きな確率となっていて、ゼロ過剰ガンマ・ポアソン回帰の適用には、無理がある。

ガンマ・ポアソン回帰

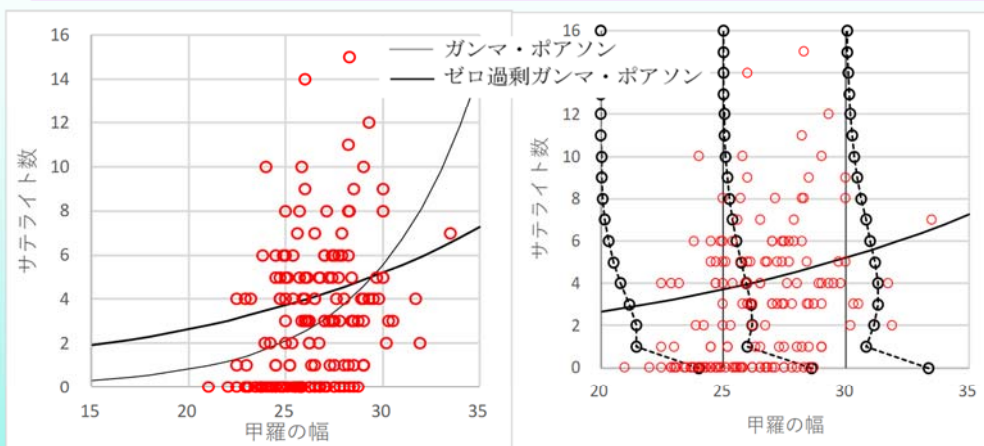


甲羅の幅が大きい時にガンマ・ポアソン分布の仮定は否定的

2019.09.06 高橋行雄

28

ゼロ過剰ガンマ・ポアソン回帰



甲羅の幅が大きい時にゼロ過剰ガンマ・ポアソン分布の仮定は否定的

2019.09.06 高橋行雄

29

対数リンクによるポアソン回帰は、元データには指数曲線のあてはめ、両辺に対数を取るモデルであり、ゼロ・データに対しては対数変換が行われないうように調整する仕組みになっている。この仕組みは、一般化線形モデルで分布を正規とし、対数リンクとした場合でも適用され、ゼロを含むようなデータに対し指数曲線をあてはめることが可能となる。なお、ポアソン回帰を行っても過分散が解消されないような場合に、正規分布を仮定し、対数リンクによる指数曲線をあてはめる場合にも、ゼロ・データに対する調整が行われる。

探索的ポアソン回帰は、表 1 でも示したように過剰なゼロが、どのような状況で発生するかを念頭にし、「甲羅の色」、「後体部の棘」とサテライト数の関係から、甲羅の色が暗くなるにつれゼロ・カウントが増加するが、後体部の棘については、関連が見いだされなかった。さらに、甲羅の色と後体部の棘を組み合わせても過分散は解消しなかった。

甲羅の幅と体重の 2 変数間には 0.89 と高い相関があり、2 変数のポアソン回帰に引き続き、図 4 に示したように体重を段階的に変化させた場合の甲羅の幅の推定曲線と 95%信頼区間のプロファイルから、甲羅の幅をポアソン回帰の説明変数に加える必要がないことが、視覚的に見いだされた。もちろん、2 変数のポアソン回帰の尤度比検定で、甲羅の幅の p 値は 0.3257 と有意ではないことから推測されることではあるが、JMP のプロファイル機能は、視覚的に変数相互の関連を見出し、より具体的な相互関係の理解するために有益である。

このプロファイル機能により、図 5 に示したように 4 水準の甲羅の色と体重の 2 変数に交互作用を加えたポアソン回帰で、甲羅の色が「やや明るい」場合に、体重が増えてもサテライト数が増えないことが図示され、甲羅の色が「中ぐらい、やや暗い、暗い」場合とは、全く異なるプロファイルであることが明示された。他方、図 6 に示すように後体部の棘と体重の関連には、交互作用を示唆するような兆候は見いだせなかった。

甲羅の色と後体部の棘に体重、さらにそれらの交互作用を含めたポアソン回帰は、観察データなので、データが不均一であり、解を得ることができなかった。これらの変数とサテライト数の関連を見出すためには、図 7 に示すように JMP のグラフ・ビルダーが役に立つ。最初に体重とサテライト数の散布図を描き、回帰直線と 95%信頼区間を上書きする。ここまでならば、JMP の伝統的な二変数の関係での対応と同じであるが、これに 4 水準の甲羅の色、3 水準の後体部を組み合わせた 4×3 の場合についてタイル状に体重とサテライト数の回帰直線と 95%信頼区間を並べて表示できた。グラフ・ビルダーで対数リンクのポアソン回帰が実施できれば申し分ないのであるが、残念ながら現在のバージョン 14 では対応していない。

伝統的な回帰分析であっても、名義尺度の水準ごとの散布図行列上に回帰直線の 95%信頼区間が表示されるだけでも、結果を総合的に俯瞰するために有益である。これに類似する機能が S プラスにあり、以前は愛用していたのであるが、JMP グラフ・ビルダーは、S プラスの機能を大幅に凌駕する探索的な統計解析を支援するツールとして優れている。

スライド 30

まとめ

- ◆ 探索的ポアソン回帰にJMPのプロファイル機能が、データの内部構造を可視化するために有効であった。
- ◆ 過分散の場合にポアソン回帰の結果は誇張されすぎるので、過分散の調整は必須である。
- ◆ ガンマ・ポアソン回帰、ゼロ過剰ガンマ・ポアソン回帰などの適用には注意が必要である。

2019.09.06 高橋行雄

30

参考文献

- 1) 高橋行雄(2002), GENMOD プロシジャによる計数データの解析, SAS ユーザ総会論文集:193-202.
- 2) 高橋行雄(2004), ポアソン回帰分析入門ー細胞数をカウントしたデータの解析ー, <https://www.yukms.com/biostat/takahasi/rec/017.htm>, 2019 年 7 月 19 日アクセス.
- 3) 高橋行雄(2019), ポアソン回帰を用いた勾配比検定, 2019 年度日本計量生物学会講予稿集; 65-70.
- 4) 久保拓弥(2012), データ解析のための統計モデリング入門 一般化線形モデル・階層ベイズモデル・MCMC, 岩波書店:39-65.
- 5) アグレスティ著, 渡邊裕之・菅波秀樹・吉田光弘ら訳(2003), カategoricalデータ解析入門, サイエンス社:110-127, 168-179.
- 6) Agresti,A. (2019), An Introduction to Categorical Data Analysis 3ed., Wiley.
- 7) Agresti,A. (2013), Categorical Data Analysis 3ed., Wiley.
- 8) ドブソン著, 田中豊・森川義彦・山中竹春・富田誠 訳(2008), 一般化線形モデル入門, 原著 第 2 版, 共立出版:67-80, 186-189.
- 9) Dobson,A.J., Barentt,A.G.,(2008), An Introduction to Ggeneralized Linear Models. CRP Press.

Excel, JMP ファイル一覧

サイズ	名前	種類
 23 KB	E探索的ポアソン8_02節_データリスト	Microsoft Excel ワークシート
 11 KB	E探索的ポアソン8_07節_表5	Microsoft Excel ワークシート
 67 KB	E探索的ポアソン8_12節_予測プロフィール	Microsoft Excel ワークシート
 67 KB	E探索的ポアソン8_13_14節_Excelでポアソン回帰	Microsoft Excel ワークシート
 16 KB	J探索的ポアソン8_03_04_06節_表1~4	JMP ファイル
 19 KB	J探索的ポアソン8_05節_図2	JMP ファイル
 15 KB	J探索的ポアソン8_07節_図3_分割表検定	JMP ファイル
 22 KB	J探索的ポアソン8_08節_体重 幅 プロファイル	JMP ファイル
 12 KB	J探索的ポアソン8_08節_等高線	JMP ファイル
 21 KB	J探索的ポアソン8_09節_色と体重	JMP ファイル
 8 KB	J探索的ポアソン8_10節_棘と体重	JMP ファイル
 12 KB	J探索的ポアソン8_11節_グラフビルダー	JMP ファイル
 11 KB	J探索的ポアソン8_12節_共分散	JMP ファイル
 10 KB	J探索的ポアソン8_13節_共分散	JMP ファイル
 3 KB	J探索的ポアソン8_13節_予測区間	JMP ファイル

非売品, 無断複製を禁ずる

第 7 回 続高橋セミナー

最尤法による探索的ポアソン回帰

BioStat 研究所(株)

〒105-0014 東京都 港区 芝 1-12-3 の 1005

2019 年 10 月 30 日 高橋 行雄

takahashi.stat@nifty.com , FAX : 03-342-8035