

```
1
2 # 2日目
3 # 一般化線形モデル完全攻略
4
5 # 2013年8月28日最終改訂
6 # 執筆者 馬場真哉 (北大水産 M2)
7 # ウェブサイト http://logics-of-blue.com/
8 # ミスなどのご連絡は logics.of.blue★gmail.com までお願いします。
9 # (星を@に変更)
10
11 #=====
12 # 予測と確率分布
13 # 実装する部分はありません。
14 # スライドに載せたグラフを描くプログラムだけです
15 #=====
16
17 # 正規分布じゃダメな理由
18
19 lm.dist <- lm(dist ~ speed, data=cars)
20
21 new <- data.frame(
22   speed = seq(min(cars$speed), max(cars$speed), by = 0.1))
23 new
24 pred.lm <- predict(lm.dist, new, interval="prediction")
25 head(pred.lm)
26
27 plot(cars, ylim = c(-50,120), main="正規分布", cex.main=2, pch=16)
28 lines(new$speed, pred.lm[,1], lwd=3)
29 lines(new$speed, pred.lm[,2], lty = 2, lwd=2, col=4)
30 lines(new$speed, pred.lm[,3], lty = 2, lwd=2, col=4)
31 abline(h=0, lwd=4)
32
33
```

```
34 #=====
35 # 一般化線形モデル基礎
36 #=====
37
38 #-----
39 # スライドに載せた数値の計算方法
40 #-----
41
42 # サンプルデータ
43 y <- c(7, 9, 8, 11)
44 plot(y, pch=16, ylim=c(0,15), cex=2)
45
46
47 # ポアソン分布の期待値が 5 だった時、各データが生じる確率
48 lambda <- 5
49 dpois(7, lambda=lambda)
50 dpois(9, lambda=lambda)
51 dpois(8, lambda=lambda)
52 dpois(11, lambda=lambda)
53
54 # ポアソン分布の定義通り計算
55 kaijou <- function(x) gamma(x+1)
56 kaijou(3)
57
58 exp(-lambda) * lambda^7 / kaijou(7)
59 dpois(7, lambda=lambda)
60
61 # ポアソン分布の期待値が 9 だった時、各データが生じる確率
62 lambda <- 9
63 dpois(7, lambda=lambda)
64 dpois(9, lambda=lambda)
65 dpois(8, lambda=lambda)
66 dpois(11, lambda=lambda)
67
68
69 # 期待値別、ポアソン分布のグラフ
```

```

70 x <- 0:20
71 par(mfrow=c(1,2))
72 #  $\lambda=5$  のポアソン分布
73 lambda <- 5
74 plot(dpois(x, lambda=lambda)~x, type="b", ylab="確率", main="ダメそう
75  $\lambda=5$ ")
76 points(7,dpois(7, lambda=lambda), pch=16, cex=2, col=2)
77 points(9,dpois(9, lambda=lambda), pch=16, cex=2, col=2)
78 points(8,dpois(8, lambda=lambda), pch=16, cex=2, col=2)
79 points(11,dpois(11, lambda=lambda), pch=16, cex=2, col=2)
80
81 #  $\lambda=9$  のポアソン分布
82 lambda <- 9
83 plot(dpois(x, lambda=lambda)~x, type="b", ylab="確率", main="よさそう
84  $\lambda=9$ ")
85 points(7,dpois(7, lambda=lambda), pch=16, cex=2, col=2)
86 points(9,dpois(9, lambda=lambda), pch=16, cex=2, col=2)
87 points(8,dpois(8, lambda=lambda), pch=16, cex=2, col=2)
88 points(11,dpois(11, lambda=lambda), pch=16, cex=2, col=2)
89 par(mfrow=c(1,1))
90
91
92 #  $\lambda=5$  の時の尤度
93 dpois(lambda=5, y)
94 prod(dpois(lambda=5, y))
95
96 # 対数尤度。総積に対してログをとる
97 log(prod(dpois(lambda=5, y)))
98
99 # 先に対数をとると、総和で計算できる
100 sum(log(dpois(lambda=5, y)))
101
102
103 #  $\lambda=9$  の時の尤度
104 dpois(lambda=9, y)
105 prod(dpois(lambda=9, y))

```

```

106
107 # 対数尤度。総積に対してログをとる
108 log(prod(dpois(lambda=9, y)))
109
110 # 先に対数をとると、総和で計算できる
111 sum(log(dpois(lambda=9, y)))
112
113
114
115 #-----
116 # R の optim 関数による最尤法の実装
117 #-----
118
119 # 最尤推定値をグラフで確認
120 x <- seq(0,15, by=0.01)
121 y <- c(7,9,8,11)
122 kekka <- numeric()
123
124 # λ を変えることによる尤度の変化
125 for(i in 1:length(x)){
126   kekka[i] <- prod(dpois(lambda=x[i], y))
127 }
128
129 kekka
130 plot(kekka~x, type="l", xlab="λ", ylab="確率", lwd=4)
131
132 # 力技で、最も尤度が高くなる λ を見つけ出す
133 max(kekka)
134 kekka.list <- data.frame(x=x, kekka=kekka)
135 kekka.list
136 subset(kekka.list, kekka.list$kekka==max(kekka))
137
138
139 # optim を使って最尤法してみる
140 # 尤度にマイナスをかけていることに注意
141 # optim は「最小化」するための関数。最大化したければマイナス 1 をかける

```

```

142 f.Lik <- function(para){
143   f.Lik <- -prod(dpois(lambda=para, y))
144   return(f.Lik)
145 }
146
147 # 小さいほうが尤度が高いことを意味している
148 f.Lik(9)
149 f.Lik(10)
150 f.Lik(15)
151
152 # 実はこの最適化はやや困難。methodを変えると結果が変わる
153 optim(fn=f.Lik, par=c(8), method="Brent", lower=0, upper=10)
154 optim(fn=f.Lik, par=c(8), method="BFGS")
155 optim(fn=f.Lik, par=c(8))
156
157
158
159 # 対数をとってから最尤法
160 # 対数の方が小さい値 (10 の-10 乗とか) に強いため、
161 # 対数をとってから計算するのがセオリー
162 f.logLik <- function(para){
163   f.logLik <- -sum(log(dpois(lambda=para, y)))
164   return(f.logLik)
165 }
166
167 f.logLik(9)
168 f.logLik(10)
169
170
171 optim(fn=f.logLik, par=c(8), method="Brent", lower=0, upper=10)
172 optim(fn=f.logLik, par=c(8), method="BFGS")
173 optim(fn=f.logLik, par=c(8))
174
175 # 実は最尤推定値は単なる期待値に等しい
176 mean(y)
177

```

```

178
179 # ★★★★★★★★★★★★★★★★★★
180 # ★          写経はじめ          ★
181 # ★★★★★★★★★★★★★★★★★★
182
183
184 # サンプルデータ
185 y <- c(7,9,8,11)
186 plot(y, pch=16, ylim=c(0,15), cex=2)
187
188
189 # ポアソン分布の期待値が 5 だった時、各データが生じる確率
190 lambda <- 5
191 dpois(7, lambda=lambda)
192 dpois(9, lambda=lambda)
193 dpois(8, lambda=lambda)
194 dpois(11, lambda=lambda)
195
196
197 # 対数をとってから最尤法
198 # 対数の方が小さい値 (10 の-10 乗とか) に強いため、
199 # 対数をとってから計算するのがセオリー
200 f.logLik <- function(para){
201   f.logLik <- -sum(log(dpois(lambda=para, y)))
202   return(f.logLik)
203 }
204
205 f.logLik(9)
206 f.logLik(10)
207
208 # 最適化
209 optim(fn=f.logLik, par=c(8), method="BFGS")
210
211 # 実は最尤推定値は単なる期待値に等しい
212 mean(y)
213

```

```

214
215 # ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪
216 # ♪          写経終わり           ♪
217 # ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪
218
219
220 #=====
221 # xによってyが変動することを加味したポアソン回帰
222 # スライド参照
223 #=====
224
225 # ★★★★★★★★★★★★★★★★★★★★★★
226 # ★          写経はじめ           ★
227 # ★★★★★★★★★★★★★★★★★★★★★★
228
229 x <- c(1, 2, 3, 4)
230 y <- c(5, 7, 10, 15)
231
232 plot(y ~ x, pch=16, cex=2)
233
234 # 尤度や対数尤度を計算する
235
236 # 切片1、傾きは0.2だと仮定したときの尤度は？
237 para <- c(1,0.2)
238
239 a <- para[1]
240 b <- para[2]
241
242 # expの中身が線形予測子
243 # link関数がlogなので、その逆関数であるexpをとった
244 lambda <- exp(a + b*x)
245 lambda
246
247 # 確率の計算はおなじみの d○○を使う
248 # 今回はポアソン分布なので dpois
249 Lik <- dpois(y, lambda=lambda)

```

```
250 Lik
251
252 # 対数尤度
253 logLik <- log(dpois(y, lambda=lambda))
254 logLik
255
256 # 尤度
257 prod(Lik)
258
259 # 対数尤度
260 # 普通の尤度よりも扱いやすそうな数値になっていることに注目
261 sum(logLik)
262 log(prod(Lik))
263
264 # 対数尤度で最適化するための関数
265 ML <- function(para) {
266   a <- para[1]
267   b <- para[2]
268   lambda <- exp(a + b*x)
269   logLik <- -sum(log(dpois(y, lambda=lambda)))
270   return(logLik)
271 }
272
273
274 # 対数尤度を最大化
275 result <- optim(fn=ML, par=c(1, 0.2))
276 result
277
278 # 推定されたλ
279 result.par <- result$par
280 result.par
281
282 best.lambda <- exp(result.par[1] + result.par[2]*x)
283 best.lambda
284
285 plot(y ~ x, pch=16, cex=3)
```



```

286 lines(best.lambda ~ x, lwd=4, col=2)
287
288 # ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪
289 # ♪          写経終わり          ♪
290 # ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪ ♪
291
292
293 # 以下参考
294
295 # x=1 の時の Y の確率分布
296 # x=1 の時の推定された λ
297 best.lambda[1]
298
299 y <- 0:20
300 dpois(y, best.lambda[1])
301
302 plot(dpois(y, best.lambda[1]) ~ y, type="b", lwd=3, pch=16, cex=2)
303
304 dpois(0, best.lambda[1])
305 dpois(1, best.lambda[1])
306 dpois(5, best.lambda[1])
307 dpois(10, best.lambda[1])
308
309 # x が 1 ~ 4 の時の、各々の Y の確率分布
310 par(mfrow=c(2,2))
311 plot(dpois(y, best.lambda[1]) ~ y, type="b", lwd=3,
312      pch=16, cex=1.5, main="x=1 の時の確率分布",
313      xlab="", ylab="", cex.main=2)
314 plot(dpois(y, best.lambda[2]) ~ y, type="b", lwd=3,
315      pch=16, cex=1.5, main="x=2 の時の確率分布",
316      xlab="", ylab="", cex.main=2)
317 plot(dpois(y, best.lambda[3]) ~ y, type="b", lwd=3,
318      pch=16, cex=1.5, main="x=3 の時の確率分布",
319      xlab="", ylab="", cex.main=2)
320 plot(dpois(y, best.lambda[4]) ~ y, type="b", lwd=3,
321      pch=16, cex=1.5, main="x=4 の時の確率分布",

```

```
322     xlab="", ylab="", cex.main=2)
323     par(mfrow=c(1,1))
324
325
```

```

326 #=====
327 # 残差の説明と deviance
328 #=====
329 x <- c(1, 2, 3, 4)
330 y <- c(5, 7, 10, 15)
331 plot(y ~ x, pch=16, ylim=c(3,17), cex=2)
332
333 # もし、データと完全に一致する予測値を出せていたとしたら？
334 # データが 5 だった時は期待値 5 のポアソン分布…
335 dpois(5, lambda=5)
336 dpois(7, lambda=7)
337 dpois(10, lambda=10)
338 dpois(15, lambda=15)
339
340 # 最尤法で推定された予測値から、今回のデータが得られたとしたら？
341 dpois(5, best.lambda[1])
342 dpois(7, best.lambda[2])
343 dpois(10, best.lambda[3])
344 dpois(15, best.lambda[4])
345
346
347 # データを完全に予測できていた時の対数尤度
348 logLik.full <-
349   log(dpois(5, lambda=5)) +
350   log(dpois(7, lambda=7)) +
351   log(dpois(10, lambda=10)) +
352   log(dpois(15, lambda=15))
353
354 # 最尤法の結果得られた最大化対数尤度
355 logLik.ML <-
356   log(dpois(5, lambda=best.lambda[1])) +
357   log(dpois(7, lambda=best.lambda[2])) +
358   log(dpois(10, lambda=best.lambda[3])) +
359   log(dpois(15, lambda=best.lambda[4]))
360
361

```

```
362 # 完璧な予測の対数尤度と、制約 ( $\lambda$  一定) の時の最大化対数尤度の比較
363 logLik.full
364 logLik.ML
365 logLik.full - logLik.ML
366
367 # residual deviance とは「完璧予測の対数尤度」と
368 # 「最尤法で計算された最大化対数尤度」の差
369 # 雰囲気は残差平方和に似ている
370 # 2倍してあるのは尤度比検定の都合。
371 logLik.full - logLik.ML
372 (logLik.full - logLik.ML) * 2
373
374 # R の関数で計算してみた結果
375 model.poisson <- glm(y ~ x, family="poisson")
376 summary(model.poisson)
377
378 # glm 関数の結果と optim で推定した結果がだいたい等しいことを確認
379 coef(model.poisson)
380 result.par
381
382 # deviance
383 # カンペキな予測の対数尤度と、計算されたモデルの対数尤度の差の 2 倍
384 deviance(model.poisson)
385 (logLik.full - logLik(model.poisson)) * 2
386
387 # 残差
388 # Type = respons は、Type = respons の予測値と実測値との差
389 resid(model.poisson, type="respons")
390
391 # Type = deviance が本命
392 resid(model.poisson, type="deviance")
393
394
395 # deviance 残差の計算方法
396 # 平方根をとるという行為を除き、基本的には前回計算した deviance と同じ雰囲気であ
397 ることがわかるはず
```

```
398 best.lambda <- predict(model.poisson, type="response")
399 best.lambda
400
401 sqrt( 2 * (log(dpois(5, lambda=5)) -
402   log(dpois(5, lambda=best.lambda[1]))) )
403 sqrt( 2 * (log(dpois(7, lambda=7)) -
404   log(dpois(7, lambda=best.lambda[2]))) )
405 sqrt( 2 * (log(dpois(10, lambda=10)) -
406   log(dpois(10, lambda=best.lambda[3]))) )
407 sqrt( 2 * (log(dpois(15, lambda=15)) -
408   log(dpois(15, lambda=best.lambda[4]))) )
409
410 # 正負の符号の確認
411 sign(y - best.lambda)
412
413 # devance 残差は type の指定をしなくても、デフォルトに指定されている
414 resid(model.poisson)
415
416 # deviance と残差(Type=deviance) の比較
417 # かなり普通の「残差」っぽく扱える
418 sum(resid(model.poisson)^2)
419 deviance(model.poisson)
420
421 # summary の結果
422 # Wald 検定
423 # t 検定の GLM バージョンだと思いとわかりよい
424 # GLM の場合は、各係数は正規分布に従うため、
425 # t value ではなく、z value になっている
426 summary(model.poisson)
427
428
429
430
431
432
433
```

```

434 #=====
435 # 尤度比検定
436 # 本当はサンプルサイズが大きい時しか使ってはいけない。
437 # サンプルサイズ 4 で実行したのはあくまで勉強のため
438 #=====
439
440 # ★★★★★★★★★★★★★★★★★★
441 # ★          写経はじめ          ★
442 # ★★★★★★★★★★★★★★★★★★
443
444
445 x <- c(1, 2, 3, 4)
446 y <- c(5, 7, 10, 15)
447 plot(y ~ x, pch=16, ylim=c(3,17), cex=2)
448
449 # x によって y が変わると仮定したモデル
450 # 再掲
451 model.poisson <- glm(y ~ x, family="poisson")
452 model.poisson
453 summary(model.poisson)
454
455 # 予測値 ( $\lambda \rightarrow$  平均値) がずっと一定だと仮定した NULL モデル
456 # ナイーブな予測を返す
457 # NULL モデルの deviance が Null deviance
458 model.null <- glm(y ~ 1, family="poisson")
459 summary(model.null)
460
461 # deviance の違い
462 deviance(model.null)
463 deviance(model.poisson)
464
465 # 最大化対数尤度の違い
466 logLik(model.poisson)
467 logLik(model.null)
468
469

```

```
470 # 尤度比検定！！
471 2 * (logLik(model.poisson) - logLik(model.null))
472 deviance(model.null) - deviance(model.poisson)
473
474 # deviance の差→対数尤度の差の 2 倍→尤度の比 が統計量
475 # この値が大きければ x を入れたおかげで deviance が大きく減ったことを意味する
476 # ということは、この値が大きければ有意になりやすい
477 # 初日にやった F 検定と「発想」はとてもよく似ている。
478 dev <- deviance(model.null) - deviance(model.poisson)
479 dev
480
481 # p 値
482 1 - pchisq(dev, df=1)
483
484 # R の関数で確認
485 anova(model.poisson, test="Chisq")
486
487
488
489 # AIC
490 # AIC = (最大化対数尤度 - パラメタ数) の 2 倍
491 AIC(model.poisson)
492 logLik(model.poisson)
493 (logLik(model.poisson) - 2) * -2
494 summary(model.poisson)
495
496
497
498 # ♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪
499 # ♪          写経終わり          ♪
500 # ♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪
501
502
503
504
505
```

```
506 #=====
507 # もうちょっと本物っぽいデータを使って確認
508 #=====
509
510 data4 <- read.csv("data4.csv")
511
512 plot(data4$y~data4$x)
513
514 # 尤度や対数尤度を計算する
515
516 # 切片 1、傾きは 0.2 だと仮定したときの尤度は？
517 para=c(1,0.2)
518
519 a <- para[1]
520 b <- para[2]
521
522 # exp の中身が線形予測子
523 # link 関数が log なので、その逆関数である exp をとった
524 lambda <- exp(a + b*data4$x)
525
526 # 確率の計算はおなじみの d○○を使う
527 # 今回はポアソン分布なので dpois
528 Lik <- dpois(data4$y,lambda=lambda)
529 Lik
530
531 # 対数尤度
532 logLik <- log(dpois(data4$y,lambda=lambda))
533 logLik
534
535 # 尤度
536 prod(Lik)
537 # 対数尤度
538 # 普通の尤度よりも圧倒的に扱いやすそうな数値になっていることに注目
539 sum(logLik)
540
541
```



```

542 # 尤度で最適化するための関数
543 prod.ML <- function(para) {
544   a <- para[1]
545   b <- para[2]
546   lambda <- exp(a + b*data4$x)
547   Lik <- -prod(dpois(data4$y, lambda=lambda))
548   return(Lik)
549 }
550
551 # 対数尤度で最適化するための関数
552 ML <- function(para) {
553   a <- para[1]
554   b <- para[2]
555   lambda <- exp(a + b*data4$x)
556   logLik <- -sum(log(dpois(data4$y, lambda=lambda)))
557   return(logLik)
558 }
559
560
561 # 対数尤度を最大化
562 optim(fn=ML, par=c(1, 0.2))
563 # 普通の尤度を最大化。うまくいってない。対数とるのは大事。
564 optim(fn=prod.ML, par=c(1, 0.2))
565
566
567 # R の glm 関数で計算
568 model.poisson <- glm(y~x, data=data4, family="poisson")
569 model.poisson
570 summary(model.poisson)
571
572 # 係数が optim の結果と一致していることを確認
573 coef(model.poisson)
574 optim(fn=ML, par=c(1, 0.2))$par
575
576 # 最大化対数尤度の値も一致している
577 # 正負が違うのは、optim において、最小化するためマイナス 1 をかけているから

```

```
578 logLik(model.poisson)
579 optim(fn=ML,par=c(1,0.2))$value
580
581 # 予測してみる
582 new <- data.frame(x=seq(0,10,by=0.01))
583 new
584 pred <- predict(model.poisson, new, type="respons")
585
586 plot(data4$y~data4$x)
587 lines(pred~new$x, lwd=2)
588
589
590 # predict いろいろ
591 # 基本的には respons を使うことが多いと思う
592 # しかしデフォルトはなぜか link。
593 predict(model.poisson, type="respons")
594 predict(model.poisson, type="link")
595 predict(model.poisson, type="term")
596 predict(model.poisson, type="term") + 2.7184 -
597   predict(model.poisson, type="link")
598 exp(predict(model.poisson, type="link")) -
599   predict(model.poisson, type="respons")
600
601
602
603 # 予測区間
604 # 95%の確率で「データが」収まる区間
605 # 信頼区間（「引っ張られた線が」収まる区間）を出すのは省略
606
607 p.lower <- qpois(0.025, pred)
608 p.upper <- qpois(0.975, pred)
609
610 plot(data4$y~data4$x, xlim=c(0,10),ylim=c(-1,55),
611       xlab="x", ylab="y", main="予測区間", cex.main=2)
612 lines(pred~new$x, lwd=2)
613 lines(p.lower~new$x, lwd=2, col=4)
```



```
649 #=====
650 # ポアソン回帰
651 # CPUE の標準化
652 #=====
653 # データの読み込み
654 data5 <- read.csv("data5_CPUE.csv")
655 summary(data5)
656 pairs(data5)
657
658 # CPUE の計算
659 CPUE <- data5$Catch / data5$Effort
660 CPUE
661 data5 <- cbind(data5, CPUE)
662 head(data5)
663
664 # CPUE との関係性
665 par(mfrow=c(2,2))
666 plot(data5$Catch ~ data5$SST, main="漁獲量と海面水温の関係")
667 plot(data5$CPUE ~ data5$SST, main="CPUE と海面水温の関係")
668 plot(log(data5$CPUE)~data5$SST, main="CPUE にログをとった")
669 boxplot(data5$CPUE~data5$Year, main="CPUE と年の関係")
670 par(mfrow=c(1,1))
671
672
673 # CPUE の年変動
674 Mean.CPUE <- as.numeric( tapply(data5$CPUE, data5$Year, mean) )
675 Mean.CPUE
676
677 Y <- 2001:2010
678 Y
679 plot(Mean.CPUE ~ Y, type ="l", xlim = c(2001,2010), lwd=2)
680
681
682
683
684
```

```

685 # ★★★★★★★★★★★★★★★★★★★★★★
686 # ★          写経はじめ          ★
687 # ★★★★★★★★★★★★★★★★★★★★★★
688
689 # データの読み込み
690 data5 <- read.csv("data5_CPUE.csv")
691 summary(data5)
692 pairs(data5)
693
694 # CPUE の計算
695 CPUE <- data5$Catch / data5$Effort
696 CPUE
697 data5 <- cbind(data5, CPUE)
698
699 # CPUE のカウントモデル
700 # 努力量が倍になれば、漁獲量は倍になるはずなので、オフセット項を入れる。
701 model.CPUE <- glm(
702   Catch ~ as.factor(Year) + offset(log(Effort)) + SST,
703   family="poisson",
704   data=data5)
705 summary(model.CPUE)
706 library(car)
707 Anova(model.CPUE)
708
709 # モデルの診断
710 # Cook 距離は 0.5 以上ならやや影響強い。
711 # 1 を超えているデータは一度削除してからもう一度分析するのがよい
712 # 今回はシミュレーションデータを使ったので全く問題なし
713 par(mfrow=c(2,2))
714 plot(model.CPUE, which=1:4)
715 par(mfrow=c(1,1))
716
717 # CPUE の標準化
718 # CPUE = Catch / Effort より、Effort=1 の時の Catch が求めればよい
719 # SST の影響を取り除くため、毎年 SST は平年値で一定と仮定する
720

```

```

721 new.CPUE <- data.frame(Year=2001:2010, Effort=rep(1,10), SST =
722 mean(data5$SST))
723 new.CPUE
724
725 s.CPUE <- predict(model.CPUE, new.CPUE, type="respon")
726 s.CPUE
727 Mean.CPUE <- as.numeric( tapply(data5$CPUE, data5$Year, mean) )
728 Mean.CPUE
729
730 Y <- 2001:2010
731 Y
732 plot(Mean.CPUE ~ Y, type ="l", lwd=2,
733      xlim = c(2001,2010), ylim=c(0, 2.5), main="CPUE 標準化")
734 lines(s.CPUE ~ Y, col=2, lwd=2)
735 legend("topleft", lwd=2, col=c(1,2),
736       legend=c("年平均 CPUE", "標準化された CPUE"))
737
738
739 # predict 関数を使わないで実装する
740 model.CPUE$coef
741 exp(model.CPUE$coef[1])
742
743 model.CPUE$coef
744
745
746 # CPUE の年変動の「トレンド」を見るだけなら、SST の係数は使わなくても OK
747 s.CPUE.2 <- numeric()
748 s.CPUE.2[1] <- exp(model.CPUE$coef[1])
749 for ( i in 1:9){
750     s.CPUE.2[i+1] <- exp(model.CPUE$coef[1] + model.CPUE$coef[i+1])
751 }
752 s.CPUE.2
753
754 # 全く同じトレンドになっていることに注目
755 par(mfrow=c(2,1))
756 plot(s.CPUE ~ Y, type="l")

```

```

757 plot(s.CPUE.2 ~ Y, type="l")
758 par(mfrow=c(1,1))
759
760
761 # ♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪
762 # ♪          写経終わり          ♪
763 # ♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪
764
765
766 # SAS (別の統計ソフト) のような出力を出すなら options で設定する
767 # R に慣れるため、あまり使わないほうが良い?
768 options("contrasts"=c("contr.SAS", "contr.SAS"))
769 model.CPUE.SAS <- glm(
770   Catch ~ as.factor(Year) + offset(log(Effort)) + SST,
771   family="poisson",
772   data=data5)
773 summary(model.CPUE.SAS)$coefficient
774
775 # 両者の比較
776 summary(model.CPUE)$coef
777 summary(model.CPUE.SAS)$coef
778
779 # 係数の見方
780 # 無理に SAS の表記を使う必要はないと思うが、水産ではなぜかつかわれているみたい
781
782 # 2001 年
783 model.CPUE$coef[1]
784 model.CPUE.SAS$coef[1] + model.CPUE.SAS$coef[2]
785 # 2002 年
786 model.CPUE$coef[1] + model.CPUE$coef[2]
787 model.CPUE.SAS$coef[1] + model.CPUE.SAS$coef[3]
788
789 # 2010 年 (最後の年) の SAS の出力はなぜか切片の値になっている。
790 model.CPUE$coef[1] + model.CPUE$coef[10]
791 model.CPUE.SAS$coef[1]
792

```

```

793
794 # predict 関数を使わないで、標準化する
795 exp(model.CPUE$coef[1] + mean(data5$SST) * model.CPUE$coef[11])
796 exp(model.CPUE$coef[1] + model.CPUE$coef[2] +
797     mean(data5$SST) * model.CPUE$coef[11])
798 exp(model.CPUE$coef[1] + model.CPUE$coef[3] +
799     mean(data5$SST) * model.CPUE$coef[11])
800 exp(model.CPUE$coef[1] + model.CPUE$coef[4] +
801     mean(data5$SST) * model.CPUE$coef[11])
802 exp(model.CPUE$coef[1] + model.CPUE$coef[5] +
803     mean(data5$SST) * model.CPUE$coef[11])
804 exp(model.CPUE$coef[1] + model.CPUE$coef[6] +
805     mean(data5$SST) * model.CPUE$coef[11])
806
807 s.CPUE
808
809 # CPUE の年変動の「トレンド」を見るだけなら、SST の係数は使わなくても OK
810 s.CPUE.2 <- numeric()
811 s.CPUE.2[1] <- exp(model.CPUE$coef[1])
812 for ( i in 1:9){
813     s.CPUE.2[i+1] <- exp(model.CPUE$coef[1] + model.CPUE$coef[i+1])
814 }
815 s.CPUE.2
816
817 # 全く同じトレンドになっていることに注目
818 par(mfrow=c(2,1))
819 plot(s.CPUE ~ Y, type="l")
820 plot(s.CPUE.2 ~ Y, type="l")
821 par(mfrow=c(1,1))
822
823 # 単に SST の分をかけているだけ。
824 # 一般に「標準化 CPUE」を出す際に、predict の結果を使っているのか、係数の年変動
825 を使っているのかは不明
826 s.CPUE / s.CPUE.2
827 exp(mean(data5$SST) * model.CPUE$coef[11])
828

```



```

829 # SAS の出力を用いて CPUE 標準化
830 model.CPUE.SAS$coef
831 s.CPUE.2.SAS <- numeric()
832 for ( i in 1:9){
833   s.CPUE.2.SAS[i] <- exp(model.CPUE.SAS$coef[i+1])
834 }
835 # 最終年を 0 に合わせているのでこうなる。
836 # ややわかりにくい。
837 s.CPUE.2.SAS[10] <- exp(0)
838
839 Mean.CPUE
840 s.CPUE
841 s.CPUE.2
842 s.CPUE.2.SAS
843
844 s.CPUE / s.CPUE.2
845 s.CPUE / s.CPUE.2.SAS
846
847 # 倍数の意味
848 exp(mean(data5$SST) * model.CPUE$coef[11])
849 exp(mean(data5$SST) * model.CPUE.SAS$coef[11] +
850   model.CPUE.SAS$coef[1])
851 exp(model.CPUE.SAS$coef[1])
852
853 mean(data5$Catch)
854 head(data5)
855
856 # トrendはどれを使っても同じ
857 par(mfrow=c(2,2))
858 plot(Mean.CPUE ~ Y, type="l")
859 plot(s.CPUE ~ Y, type="l")
860 plot(s.CPUE.2 ~ Y, type="l")
861 plot(s.CPUE.2.SAS ~ Y, type="l")
862 par(mfrow=c(1,1))
863

```

```

864 #=====
865 # ロジスティック回帰
866 # Rにもともと入っているタイタニックデータを使用
867 #=====
868
869
870 # ★★★★★★★★★★★★★★★★★★
871 # ★          写経はじめ          ★
872 # ★★★★★★★★★★★★★★★★★★
873
874
875 Titanic
876
877 # データの整形。使いにくいので
878 Titanic2 <- data.frame(Titanic)
879 Titanic2
880
881 data.Titanic <- data.frame(
882   cbind(Titanic2[1:16,-4], Survive = Titanic2[17:32,5]))
883 names(data.Titanic)[4] <- "Death"
884 data.Titanic
885
886 # 二項分布のモデルの作成
887 # cbind(成功事例, 失敗事例)
888 model.binom_1 <- glm(
889   cbind(Survive, Death) ~ .,
890   data=data.Titanic,
891   family="binomial")
892 summary(model.binom_1)
893
894 # オッズの変化率
895 exp(coef(model.binom_1))[-1]
896
897 # モデルの診断
898 par(mfrow=c(2,2))
899 plot(model.binom_1, which=1:4)

```



```

935 #=====
936 # Gamma 分布
937 #=====
938
939 # ★★★★★★★★★★★★★★★★★★
940 # ★          写経はじめ          ★
941 # ★★★★★★★★★★★★★★★★★★
942
943 cars
944 plot(cars)
945
946 # ガンマ回帰
947 model.Gamma <- glm(dist ~ speed, data=cars, family="Gamma"(link=log))
948 summary(model.Gamma)
949
950 # モデルの診断
951 par(mfrow=c(2,2))
952 plot(model.Gamma, which=1:4)
953 par(mfrow=c(1,1))
954
955 # ♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪
956 # ♪          写経終わり          ♪
957 # ♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪
958
959 # 予測区間を出す
960
961 new <- data.frame(
962   speed = seq(min(cars$speed), max(cars$speed), by = 0.1))
963 new
964
965 pred.Gamma.response <- predict(model.Gamma, new, type="response")
966 pred.Gamma.response
967
968
969 # 必要なパラメタの準備
970 library("MASS")

```

```
971 shape <- gamma.shape(model.Gamma)$alpha
972 shape
973 scale <- pred.Gamma.response/shape # scale パラメータ
974 lower <- qgamma(0.025, shape = shape, scale = scale)
975 upper <- qgamma(0.975, shape = shape, scale = scale)
976
977
978 # 予測区間の図示
979 par(mfrow=c(1,2))
980 plot(cars, ylim =c(-50,120), main="Gamma 分布", cex.main=2)
981 lines(new$speed, pred.Gamma.response, lwd=2)
982 lines(new$speed, lower, lty = 2, lwd=2, col=4)
983 lines(new$speed, upper, lty = 2, lwd=2, col=2)
984 abline(h=0)
985
986 # 比較用の普通の正規線形モデル (単回帰)
987 model.lm <- lm(dist ~ speed, data=cars)
988 pred.lm <- predict(model.lm, new, interval="prediction")
989 head(pred.lm)
990 plot(cars, ylim =c(-50,120), main="正規分布", cex.main=2)
991 lines(new$speed, pred.lm[,1], lwd=2)
992 lines(new$speed, pred.lm[,2], lty = 2, lwd=2, col=4)
993 lines(new$speed, pred.lm[,3], lty = 2, lwd=2, col=2)
994 abline(h=0)
995 par(mfrow=c(1,1))
996
997
```

```

998 #=====
999 # 対数線形モデル
1000 #=====
1001
1002 # ★★★★★★★★★★★★★★★★★★
1003 # ★          写経はじめ          ★
1004 # ★★★★★★★★★★★★★★★★★★
1005
1006
1007 sample.table <- data.frame(
1008     type = rep(c("type_A", "type_B"), 2),
1009     Result = rep(c("Good", "Bad"), each=2),
1010     Freq = c(35, 45, 50, 30)
1011 )
1012 sample.table
1013
1014 sample.table2 <- xtabs(Freq ~., sample.table)
1015 sample.table2
1016
1017 #  $\chi$  二乗検定
1018 # イエーツの補正をされているらしい?
1019 chisq.test(sample.table2)
1020 # 普通の  $\chi$  二乗検定
1021 chisq.test(sample.table2, correct=F)
1022 # フィッシャーの正確確率検定
1023 fisher.test(sample.table2)
1024
1025 # 対数線形モデル
1026 # 交互作用の有無を検定する
1027 glm.loglin <- glm(Freq ~ (.)^2, data=sample.table, family="poisson")
1028 summary(glm.loglin)
1029 library(car)
1030 Anova(glm.loglin, type="III")
1031
1032 # loglin 関数を使う
1033 m.loglin <- loglin(sample.table2, margin=list(c(1), c(2)), param=T,

```

```
1034 fit=F)
1035 summary(m.loglin)
1036
1037 # 対数尤度比統計量 (deviance の差)
1038 m.loglin$lrt
1039
1040 # 尤度比検定
1041 1 - pchisq(m.loglin$lrt, m.loglin$df)
1042
1043 #  $\chi^2$  乗検定
1044 m.loglin$pearson
1045 1 - pchisq(m.loglin$pearson, m.loglin$df)
1046 chisq.test(sample.table2, correct=F)
1047
1048
1049
1050 # もっと複雑なデータ
1051 HairEyeColor
1052 # ポアソン回帰できるようにデータを整形
1053 Hair <- data.frame(HairEyeColor)
1054 Hair
1055
1056 #対数線形モデル
1057 model.loglin.full <- glm(Freq ~ (.)^3, data=Hair, family="poisson")
1058
1059 # 検定
1060 summary(model.loglin.full)
1061 Anova(model.loglin.full, type="III")
1062
1063 # loglin を使っても一緒の結果。
1064 # margin の設定が難しいので、glm 関数のほうが楽
1065 fm <- loglin(HairEyeColor, list(c(1, 2), c(1, 3), c(2, 3)))
1066 fm
1067 1 - pchisq(fm$lrt, fm$df)
1068
1069
```

```
1070
1071 # 3 次の交互作用を切った
1072 model.loglin.2 <- update(model.loglin.full, ~. - Hair:Eye:Sex)
1073 Anova(model.loglin.2, type="III")
1074
1075 # 2 次の交互作用を一つ切った
1076 model.loglin.3 <- update(model.loglin.2, ~. - Eye:Sex)
1077 Anova(model.loglin.3, type="III")
1078 summary(model.loglin.3)
1079
1080 # モデルの診断
1081 par(mfrow=c(2,2))
1082 plot(model.loglin.3, which=1:4)
1083 par(mfrow=c(1,1))
1084
1085
1086 # AIC でモデル選択
1087 library(MuMIn)
1088 dredge(model.loglin.full, rank="AIC")
1089 plot(dredge(model.loglin.full, rank="AIC"))
1090
1091
1092 # ♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪
1093 # ♪                写経終わり                ♪
1094 # ♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪♪
1095
```



```
1096 #=====
1097 # ゼロ切断とゼロ過剰モデル
1098 #=====
1099
1100 # ★★★★★★★★★★★★★★★★★★
1101 # ★          写経はじめ          ★
1102 # ★★★★★★★★★★★★★★★★★★
1103
1104 # install.packages("VGAM")
1105 # install.packages("pscl")
1106 library(VGAM)
1107 library(pscl)
1108
1109 # Zero Truncated
1110
1111 # サンプルデータの作成
1112 # 期待値を表すパラメタλの政界値は「1」
1113 set.seed(1)
1114 Y <- rpois(lambda=1, n=30)
1115 Y
1116 mean(Y)
1117
1118 # 問題なし
1119 exp(glm(Y ~ 1, family="poisson")$coef)
1120
1121
1122 # ゼロを切断してみる
1123 Y.cut <- subset(Y, Y>0)
1124 Y.cut
1125
1126 # 期待値が高くなっている
1127 # こちらはウソの期待値。これに騙されてはいけない。
1128 mean(Y.cut)
1129
1130 # 普通の GLM
1131 model.pois <- glm(Y.cut ~ 1, family="poisson")
```



```
1168
1169
1170 # ゼロを切った
1171 data.cut <- subset(data, data$Y2>0)
1172 data.cut
1173
1174 # 0を切る前と切った後で比較
1175 par(mfrow=c(1,2))
1176 plot(data$Y2 ~ data$x, ylim=c(0,25))
1177 plot(data.cut$Y2 ~ data.cut$x, ylim=c(0,25))
1178 par(mfrow=c(1,1))
1179
1180 # 普通のGLM
1181 model.pois_2 <- glm(Y2 ~ x, data=data.cut, family="poisson")
1182 summary(model.pois_2)
1183
1184 # ゼロ切断モデル
1185 model.pos.pois_2 <- vglm(Y2 ~ x, data=data.cut, family="pospoisson")
1186 summary(model.pos.pois_2)
1187
1188 # 係数の確認
1189 coef(model.pois_2)
1190 coef(model.pos.pois_2)
1191 # 正しい係数
1192 coef(glm(Y2 ~ x, family="poisson"))
1193
1194
1195 # 残念ながら尤度比検定はちょっとやりにくい
1196 library(car)
1197 Anova(model.pos.pois_2)
1198
1199 # lrtestを使う
1200 lrtest(model.pos.pois_2)
1201
1202 # 残差の確認もやや面倒
1203 plot(model.pos.pois_2)
```

```

1204 # 一応出せるが、Working の残差。deviance ではない
1205 plot(resid(model.pos.pois_2) ~ exp(predict(model.pos.pois_2)))
1206
1207
1208 # 予測値の確認
1209 plot(data.cut$Y2 ~ data.cut$x, ylim=c(0,25))
1210
1211 new <- data.frame(x = seq(-2,3, by=0.1))
1212
1213 lines (
1214   predict(glm(Y2 ~ x, family="poisson", data=data),
1215     new, type="respons") ~
1216     new$x, lwd=3)
1217 lines (
1218   predict(model.pos.pois_2, new, type="respons") ~ new$x, lwd=2, col=2)
1219 lines (
1220   predict(model.pois_2, new, type="respons") ~ new$x, lwd=2, col=4)
1221
1222 legend("topleft", lwd=c(3,2,2), col=c(1,2,4),
1223   legend=c("正しい","ゼロ切断","ダメな普通の GLM"))
1224
1225
1226
1227 #=====
1228 # ゼロ過剰モデル
1229 #=====
1230
1231 # ★★★★★★★★★★★★★★★★★★★★★★
1232 # ★          写経はじめ          ★
1233 # ★★★★★★★★★★★★★★★★★★★★★★
1234
1235 library(pscl)
1236 library(lmtest)
1237
1238
1239 # データ

```



```

1276 # art ~ . の[.]は全ての変数を入れるという指示。「1」を入れると、定数項だけにな
1277 る。
1278
1279 model_zip <- zeroinfl(art ~ . | 1, data = bioChemists)
1280 model_zinb <- zeroinfl(art ~ . | 1, data = bioChemists, dist = "negbin")
1281 summary(model_zip)
1282 summary(model_zinb)
1283
1284 AIC(model_zip,model_zinb)          # 負の二項分布を支持
1285
1286 # しかし、これだとゼロインフルモデルにした意味がない
1287 summary(model_zinb)
1288 summary(model_nb)
1289
1290 # MuMInによるモデル選択にも対応している
1291 library(MuMIn)
1292 list.1 <- dredge(model_zinb, rank="AIC")
1293 list.1
1294
1295 plot(list.1)
1296
1297 all.model <- get.models(list.1)
1298 best.model<-all.model[1]
1299 best.model
1300
1301
1302
1303 # ★★★★★★★★★★★★★★★★★★★★★★
1304 # ★          写経はじめ          ★
1305 # ★★★★★★★★★★★★★★★★★★★★★★
1306
1307
1308 # 全部入れたモデル
1309 # 理由は分からないが、略式で formula を書くと、MuMIn を使った時にバグが出る
1310 model_zip2 <- zeroinfl(
1311   art ~ fem + mar + kid5 + phd + ment |

```

```
1312         fem + mar + kid5 + phd + ment, data = bioChemists)
1313 model_zinb2 <- zeroinfl(
1314   art ~ fem + mar + kid5 + phd + ment |
1315     fem + mar + kid5 + phd + ment, data = bioChemists, dist = "negbin")
1316
1317 AIC(model_zip2,model_zinb2)
1318
1319 # ZINB は二項分布における正負の符号が違っていることに注意
1320 # なぜならば、ZINB は「偽物のゼロ」であるかどうかを二項分布で見分けているから。
1321 # プラスの影響ならば「偽物のゼロ」が生じやすいということ
1322
1323 summary(model_zinb2)
1324
1325
1326 # 大変時間がかかるので注意
1327 # やらない方がよいかも
1328 # 7分くらいかかった
1329 # また、分離が起こって、計算できないモデルも多かった
1330
1331 library(MuMIn)
1332 # Not Run
1333 # system.time( list.2 <- dredge(model_zinb2, rank="AIC"))
1334
1335 plot(list.2)
1336 head(list.2, 10)
1337
1338 # all.model.2 <- get.models(list.2)
1339 # best.model.2 <- all.model.2[1]
1340 best.model.2
1341
1342 # 先生がたくさん論文を書いていると、学生が優秀なのとうっかりで
1343 # 論文を出さなくなっちゃうことは少ない。ゆえに影響力は負。
1344 # 結婚していると、学生が優秀なのとうっかりで論文を出さなくなっちゃうことは少ない。
1345 # ゆえに影響力は負。
1346 # やる気が増すのかもしれない？
1347 # 変数の解釈が難しいので要注意。
```





```
1384     data=bioChemists.bin,
1385     family="binomial")
1386 summary(model_bin)
1387
1388
1389 # 書いた人はどれだけたくさん書くか。
1390 bioChemists.pos <- subset(bioChemists, art != 0)
1391 bioChemists.pos
1392
1393 library(VGAM)
1394 model.pos <- vglm(
1395     art ~ fem + mar + kid5 + phd + ment,
1396     family="posnegbinomial",
1397     data=bioChemists.pos)
1398 summary(model.pos)
1399
1400
1401 # ZINB は二項分布における正負の符号が違っていることに注意
1402 # なぜならば、ZINB は「偽物のゼロ」であるかどうかを二項分布で見分けているから。
1403 # プラスの影響ならば「偽物のゼロ」が生じやすいということ
1404 # 先生がたくさん論文を書いていると、学生が優秀なのによりっきり論文を出さなくなっ
1405 ちゃうことは少ない。ゆえに影響力は負。
1406
1407 # 通常の GLM と比較すると hurdle は、二項分布+ゼロ切断だということがよくわかる
1408 # 二項分布にプラスの影響ならば純粹に「論文を投稿しやすい」という意味になる
1409 # 先生がたくさん論文を投稿していたら、学生も論文を投稿しやすいという意味
1410 summary(model_zinb2)
1411 summary(model_hurdle)
1412 summary(model_bin)
1413
1414 summary(model.pos)
1415 summary(model_hurdle)
1416
1417
```

```
1418 #=====
1419 # GLMM
1420 #=====
1421
1422 # ★★★★★★★★★★★★★★★★★★
1423 # ★          写経はじめ          ★
1424 # ★★★★★★★★★★★★★★★★★★
1425
1426
1427 #=====
1428 # Step1 過分散への対応
1429 #=====
1430 data6 <- read.csv("data6.csv")
1431
1432 head(data6)
1433 pairs(data6)
1434
1435 # 通常の GLM
1436 model.glm.pois <- glm(
1437   y ~ x1 + x2,
1438   family="poisson",
1439   data=data6)
1440 summary(model.glm.pois)
1441
1442 # 尤度比検定 Type II ANOVA
1443 library(car)
1444 Anova(model.glm.pois)
1445
1446 # 擬似尤度
1447 model.glm.quasi.pois <- glm(
1448   y ~ x1 + x2,
1449   family="quasipoisson",
1450   data=data6)
1451 summary(model.glm.quasi.pois)
1452
1453
```

```
1454 # 尤度比検定
1455 Anova(model.glm.quasi.pois)
1456
1457 # GLMM
1458 library(lme4)
1459
1460 # ランダムエフェクトを入れる準備
1461 ID <- 1:nrow(data6)
1462 ID
1463
1464 data6$ID <- ID
1465 head(data6)
1466
1467 # GLMM
1468 model.glmm <- glmer(
1469   y ~ x1 + x2 + (1 | ID),
1470   data=data6,
1471   family=poisson)
1472 summary(model.glmm)
1473 Anova(model.glmm)
1474
1475 # 比較
1476 summary(model.glm.pois)
1477 summary(model.glm.quasi.pois)
1478 summary(model.glmm)
1479
1480
1481
1482 #=====
1483 # Step2.擬似反復を攻略する
1484 #=====
1485 head(sleepstudy)
1486 pairs(sleepstudy)
1487
1488 # 睡眠時間が減ることによる悪影響を調べた実験データ
1489 # 初日は普通に寝る
```

```
1490 # 二日目以降は一日 3 時間しか寝ない
1491 # Reaction はあるテストに対する反応速度。早いほうが良い評価である。
1492 # 寝不足は反応速度にどれくらい影響するか?
1493
1494 # このデータは一人の被験者につき 10 日間調べている
1495 # 本当のサンプルサイズは被験者の数であるはず
1496 # しかし、GLM をした時のサンプルサイズは実質 10 倍されている
1497 # サンプルサイズが多いように見せかけられると有意になりやすくなってしまう。
1498
1499 glm.sleep <- glm(Reaction ~ Days, data=sleepstudy)
1500 summary(glm.sleep)
1501 Anova(glm.sleep)
1502
1503 # GLMM で擬似反復を攻略する
1504
1505 # 被験者ごとに反応速度が異なるというランダムエフェクトを入れる
1506 glmm.sleep_1 <- lmer(Reaction ~ Days + (1|Subject), sleepstudy)
1507
1508 # さらに追加で、被験者ごとに日数による影響が異なることもモデルに組み込む
1509 # 0+がないと切片に 2 重にランダムエフェクトが組み込まれてしまうので注意
1510 glmm.sleep_2 <- lmer(
1511   Reaction ~ Days + (1|Subject) + (0+Days|Subject), sleepstudy)
1512
1513 summary(glmm.sleep_1)
1514 summary(glmm.sleep_2)
1515
1516 Anova(glmm.sleep_1)
1517 Anova(glmm.sleep_2)
1518
1519 anova(glmm.sleep_1, glmm.sleep_2)
1520
1521
1522 # 擬似反復と過分散を両方退治する
1523
1524 head(cbpp)
1525 pairs(cbpp)
```

