

```
1  # 一日目
2  # 正規線形モデリング
3
4  # 2013年8月28最終改訂
5  # 執筆者 馬場真哉（北大水産M2）
6  # ウェブサイト http://logics-of-blue.com/
7  # ミスなどのご連絡は logics.of.blue★gmail.com までお願いします。
8  # (星を@に変更)
9
10 =====
11 # Rによる統計モデリング
12 # この章は全てのプログラムを記述～実行する
13 =====
14
15 # ★★★★★★★★★★★★★★★★★★★
16 # ★ 写経はじめ ★
17 # ★★★★★★★★★★★★★★★★★★★
18
19
20 # こんにちは世界
21 print("Hallow World")
22
23 # printは省略可能
24 "Hallow World"
25
26 # ""を抜かすとエラー
27 Hallow World
28
29 # 簡単な計算
30 1 + 2
31 5*6
32 3 + 4/2 - 1
33 2^10
34
35 # コメント部の説明
36 # 1 + 2
```

```
37
38
39 # 変数の定義
40 # 「<-」の左右にはスペースを必ずあけること。
41 # 見やすいプログラムを書く癖をつけた方がよい
42 x
43 x <- 3
44 x
45 x*2
46
47
48 # データの読み込みと表示
49 # データをコピーしてから実行する
50 data.0.clip <- read.delim("clipboard")
51 data.0.clip
52 head(data.0.clip)
53 pairs(data.0.clip)
54
55 # データの読み込み方法その2
56 # 作業ディレクトリを変更してから以下を実行
57 getwd()
58
59 # Not Run
60 # setwd()
61 # 任意のワーキングディレクトリを指定
62 # コンソールをクリック→ファイル→ディレクトリの変更 から指定できる
63 # getwd()で場所を確認したうえで setwd()に値をコピペで渡すと楽。
64
65 # 読み込みはこの一行でOK
66 data.0 <- read.csv("data0.csv")
67 data.0
68 head(data.0)
69 pairs(data.0)
70
71
72 # データの中身の取出し
```

```
73 # attach() 関数は"絶対に"使わないこと。
74 food          # エラーが出る
75 names(data.0) # data.0 に格納されているモノの名前を調べる
76 data.0$food   # OK
77
78
79 # プロット
80 # 興味の対象を「～」の左側におく
81 # 興味の対象 length はえさの量によってどう変わるか？
82 plot(data.0$length ~ data.0$food)
83
84 # 興味の対象 length は薬の有無によってどう変わるか？
85 plot(data.0$length ~ data.0$medicine)
86
87
88 # ちょっと複雑なプロット
89 # 長いプログラムは分かりやすく改行するとよい
90 # インデント（左端にスペースを空ける）するとなお見やすい。
91 # Google 基準ではインデントはスペース 2 つ分
92 # ただし丸っこで囲まれている場合は、一行目に合わせてインデントを付ける
93 # ↑ただし、これはめんどくさいので私は使っていない。スペース 2 つで十分と思う。
94 # 無理に Google に合わせる必要はないが、改行&インデントは実務上必須。
95 # 私はタブを使うこともよくあるが、Google 基準では使ってはいけないらしい？
96
97 # 見づらいコードを書き続けると、いつか"必ず"後悔することになるので注意。
98 # Google に合わせる必要は皆無だが、自分なりにルールを決める必要はある。
99
100 plot(
101   data.0$length ~ data.0$food,
102   col = c(2,3)[data.0$medicine],
103   pch=16,
104   ylab = "Length",
105   xlab = "food",
106   main = "薬の有無別、体長と餌の量の関係",
107   cex.main = 1.5,
108   font.lab = 2
```

```
109  )
110
111 # 凡例
112 legend(
113   "topleft",
114   legend = c("薬あり", "薬なし"),
115   col = c(2,3),
116   bty = "n",
117   pch = 16
118 )
119
120 # ヘルプ
121 ?legend
122
123 # モデルの作成
124 # 応答変数 : length
125 # 説明変数 : food & medicine
126 # 興味のある対象 (応答変数 length) は説明変数によってどれくらい変化するか? の
127 モデリング
128 lm.model.0 <- lm(length ~ food + medicine, data = data.0)
129 lm.model.0
130
131 # ANOVA による検定結果
132 # 説明は後程
133 anova(lm.model.0)
134
135 # より詳しい結果の表示→最初は飛ばす
136 summary(lm.model.0)
137
138 # 予測
139 predict(
140   lm.model.0,
141   newdata = data.frame(food = 50, medicine = "medicine"))
142
143 # 予測区間つき
144 predict(
```

```
145 lm.model.0,
146 newdata = data.frame(food = 50, medicine = "medicine"),
147 interval = "prediction", level = 0.95)
148
149 # newdata の指定なし
150 predict(lm.model.0)
151
152 # データの型の確認
153 class(data.0)
154
155 # ベクトルデータ
156 vec <- c(1, 2, 3, 4.5, 5, 6, 7.2, 8, 9, 9.9)
157 vec
158 int <- 1:10
159 int
160
161 # 等差数列
162 seq(from = 0.1, to = 1, by = 0.1)
163
164 # データフレームに格納する
165 d <- data.frame(
166   vec = vec,
167   int = int,
168   seq = seq(from = 0.1, to = 1, by = 0.1)
169 )
170
171 # データの取り出し方
172 d
173 d[1,]
174 d[,1]
175 d[2,3]
176 names(d)
177 d$vec
178 d[,c("vec")]
179 d[,c("vec", "seq")]
180
```

```
181 # 予測の図示
182 # newdata の作成
183 newfood <- seq(from = min(data.0$food), to = max(data.0$food), by = 1)
184 newfood
185
186 # 薬があるときのえさと体長の関係の予測のためのデータセット
187 new.1 <- data.frame(
188   food = newfood,
189   medicine = "medicine")
190
191 # 薬がない時のえさと体長の関係の予測のためのデータセット
192 new.2 <- data.frame(
193   food = newfood,
194   medicine = "na")
195
196 # 薬の有無別に予測
197 pred.1 <- predict(
198   lm.model.0, newdata = new.1,
199   interval = "prediction", level = 0.95)
200 pred.2 <- predict(
201   lm.model.0, newdata = new.2,
202   interval = "prediction", level = 0.95)
203
204 # 図示
205 # ちょっと複雑なプロット
206 # これは以前作ったものと同じ。
207 # コピペする方が早いし、正確。
208 # 意地を張ると遅くなるうえに不正確になるので注意。ぜひコピペしてください。
209 plot(data.0$length ~ data.0$food,
210       col = c(2,3)[data.0$medicine],
211       pch=16,
212       ylab = "Length",
213       xlab = "food",
214       main = "薬の有無別、体長と餌の量の関係",
215       cex.main = 1.5,
216       font.lab = 2
```



```

239  #=====
240  # 統計の基本と t 検定
241  #=====
242
243  # 大阪データ
244  osaka <- c(19,19,20,20,20,20,20,21,21,21)
245  # 期待値
246  19*(2/10) + 20*(5/10) + 21*(3/10)
247  mean(osaka)
248
249  # 標本分散
250  ((19-20.1)^2)*(2/10) + ((20-20.1)^2)*(5/10) +
251  ((21-20.1)^2)*(3/10)
252
253  # 普遍分散
254  length(osaka)
255  ((19-20.1)^2)*(2/9) + ((20-20.1)^2)*(5/9) + ((21-20.1)^2)*(3/9)
256  var(osaka)
257
258
259  # 東京データ
260  tokyo <- c(-50, 0, 0, 20, 20, 20, 20, 70, 70)
261
262  # 期待値
263  mean(tokyo)
264
265  # 標本分散
266  (-50-19)^2 * (1/10) + (0-19)^2 * (2/10) +
267  (20-19)^2 * (5/10) + (70-19)^2 * (2/10)
268
269  # 不偏分散
270  (-50-19)^2 * (1/9) + (0-19)^2 * (2/9) +
271  (20-19)^2 * (5/9) + (70-19)^2 * (2/9)
272  var(tokyo)
273
274

```

```
275 #=====
276 # t 検定
277 #=====
278
279 # ★★★★★★★★★★★★★★★★★★
280 # ★ 写経はじめ ★
281 # ★★★★★★★★★★★★★★★★★★
282
283 # サンプルデータの作成
284 d <- c(-1, -1, 0, 0, 1, 3, 5, 6, 7, 7)
285
286 # 平均値
287 mean(d)
288
289 # 標準偏差
290 sd(d)
291
292 # サンプルサイズ
293 length(d)
294
295 # 標準誤差
296 sd(d) / sqrt(length(d))
297
298 # t 検定
299 mean <- mean(d)
300 std.error <- sd(d) / sqrt(length(d))
301
302 t.value <- mean / std.error
303 t.value
304 (1-pt(t.value, df=length(d)-1))*2
305
306 # Rに入っている関数を使う
307 t.test(d)
308
309
310
```



```
336 #=====
337 # 分散分析 ANOVA
338 #=====
339
340 # グラフを描くだけ
341
342 # サンプルデータの作成
343 d2 <- data.frame(
344   Y = c(c(1,2,3,4,5), c(4,5,6,7,8), c(7,8,9,10,11)),
345   option = rep(c("A", "B", "C"), each=5)
346 )
347
348 d2
349
350 # データの可視化
351 plot(d2$Y ~ d2$option)
352
353 # スライドに載せたような散布図形式のグラフを作る
354 par(mar=c(5, 6, 3, 3))
355 plot.default(
356   d2$Y ~ d2$option,
357   ylim=c(0,12), xlim=c(0.5,3.5),
358   ylab="結果", xlab="選択肢",
359   cex=2,
360   cex.lab=2, cex.main=3,
361   xaxt="n")
362 axis(side=1, 1:3, LETTERS[1:3])
363
364
365 # ★★★★★★★★★★★★★★★★★★
366 # ★          写経はじめ          ★
367 # ★★★★★★★★★★★★★★★★★★★★
368
369 # サンプルデータの作成
370 d2 <- data.frame(
371   Y = c(c(1,2,3,4,5), c(4,5,6,7,8), c(7,8,9,10,11)),
```

```
372     option = rep(c("A", "B", "C"), each=5)
373 )
374
375 d2
376
377 # データの可視化
378 plot(d2$Y ~ d2$option)
379
380 # Y の総平均
381 mean(d2$Y)
382
383 # option ごとの Y の期待値の算出
384 tapply(d2$Y, d2$option, mean)
385
386 # R で分散分析を実装する
387 lm.model.anova <- lm(Y ~ option, data=d2)
388
389 # 係数の確認
390 lm.model.anova
391 summary(lm.model.anova)
392
393 # 分散分析による予測
394 predict(lm.model.anova, data.frame(option=c("A", "B", "C")))
395
396 # 予測値は option 別の期待値に等しい
397 tapply(d2$Y, d2$option, mean)
398
399
400 =====
401 # 分散分析モデルにおける検定
402 # 検定スライドを見せてから実施
403 =====
404
405 # 予測値の変化の大きさ
406 d2.effect <- data.frame(
407   Y = c(rep(3,5), rep(6,5), rep(9,5)),
```

```
408     option = rep(c("A", "B", "C"), each=5)
409 )
410
411 d2.effect
412
413 # 偏差平方和 (分散の分子の部分)
414 sum((d2.effect$Y-mean(d2.effect$Y))^2)
415
416 # 偏差平方和を自由度で割る
417 # 自由度 = 自由に動ける値
418 # d2.effect の中身は 3, 6, 9 の 3 種類のみ。
419 # また d2.effect の期待値は元データの期待値と等しくなるという制約がある
420 # 自由に動ける値は 「3 - 1 = 2」
421 # 予測値の差の大きさを表す分散
422 sum((d2.effect$Y-mean(d2.effect$Y))^2) / 2
423
424 effect <- sum((d2.effect$Y-mean(d2.effect$Y))^2) / 2
425 effect
426
427 # ノイズの大きさ
428 # 単に実データから予測値を引いただけ
429 # ↓これが予測誤差(ノイズ)
430 d2$Y - d2.effect$Y
431
432 # データフレームにまとめ
433 d2.noise <- data.frame(
434   Y = d2$Y - d2.effect$Y,
435   option = rep(c("A", "B", "C"), each=5)
436 )
437
438 d2.noise
439
440 # 偏差平方和 (分散の分子の部分)
441 sum( (d2.noise$Y-mean(d2.noise$Y))^2 )
442
443 # 自由度
```

```

444 # 誤差というモノの制約=合計値は 0 になる
445 # そうしないと、単なるノイズ以上の意味を持つてしまうから。
446 # ゆえに、option ごとに合計値が 0 にならなければならない。
447 # サンプルサイズから option の個数（この場合 3）を引く。
448 length(d2.noise$Y)-3
449
450 # ノイズの大きさを表す分散
451 sum( (d2.noise$Y-mean(d2.noise$Y))^2 ) / 12
452
453 noise <- sum( (d2.noise$Y-mean(d2.noise$Y))^2 ) / 12
454 noise
455
456 # F 比
457 # Effect は option の数-1 で割っている
458 # Noise はサンプルサイズ-option 数で割っている
459 # ゆえに、サンプルサイズが多くなるほど Effect と Noise の比は大きくなる。
460 effect
461 noise
462
463 # F 比
464 F.ratio <- effect/noise
465 F.ratio
466
467 # 検定
468 # option によって有意に予測は変わらないことがわかっているデータを集める部分は
469 # 数式を使って省略する。
470 # ここにも自由度が使われている。サンプルサイズが効いているということ。
471 1 - pf(F.ratio, 2, 12)
472
473
474 # 関数を使って計算
475 lm.model.anova <- lm(Y ~ option, data=d2)
476 anova(lm.model.anova)
477 summary(lm.model.anova)
478
479

```



```
516
517 sum.sq.ANOVA
518 sum.sq.NULL
519
520 # 自由度の差
521 3 - 1
522
523 # 分散分析モデルの予測誤差の分散 (再掲)
524 noise.ANOVA <- sum( (d2.noise$Y-mean(d2.noise$Y))^2 ) / 12
525 noise.ANOVA
526
527 # 見やすいように括弧を増やしている。本来は(sum.sq.NULL - sum.sq.ANOVA) / 2
528 / Noise.ANOVA
529 # 残差の差分を自由度の差分で割ったものが分子
530 # option 入りの ANOVA モデルの予測残差の分散が分母
531 F.ratio2 <- ( (sum.sq.NULL - sum.sq.ANOVA) / 2 ) / noise.ANOVA
532 F.ratio2
533
534 1 - pf(F.ratio2, 2, 12)
535
536
537 # ナイーブ予測と分散分析モデル予測との予測誤差の比較
538 anova(lm.model.NULL, lm.model.anova)
539
540 # 普通の ANOVA
541 anova(lm.model.anova)
542
543 # 分散分析とは予測残差の違いを検定するものだ、とご理解ください
544
545
546
547
548
549
```

```
550  #=====
551  # 回帰分析
552  #=====
553
554  # ★★★★★★★★★★★★★★★★★★
555  # ★          写経はじめ      ★
556  # ★★★★★★★★★★★★★★★★★★
557
558  # サンプルデータ
559  d3 <- data.frame(
560    Y = c(3, 5, 4, 6, 7),
561    X = 1:5
562  )
563
564  d3
565  nrow(d3)
566  plot(d3$Y ~ d3$X)
567
568  summary(lm(Y ~ X, data=d3))
569  lm(Y ~ X, data=d3)
570  anova(lm(Y ~ X, data=d3))
571
572
573  # optimで最小二乗法
574
575  # 自作関数の作成
576  # 予測するための関数
577  yosoku <- function(a, b, X) {
578    Yhat <- a*X + b
579    return(Yhat)
580  }
581
582  # 2*3 + 1
583  yosoku(a=2, b=1, X=3)
584
585
```

```

586  # 予測結果の残差平方を出す関数
587  zansa <- function(Y, Yhat) {
588      zansa <- (Y - Yhat)^2
589      return(zansa)
590  }
591
592  # (2-3)^2
593  zansa(2, 3)
594
595
596  # 残差平方和計算
597  # for ループを使用する（本当は使わなくていい。勉強用）
598  zansa2 <- numeric()
599  a <- 2
600  b <- 2
601  for ( i in 1:nrow(d3) ) {
602      Yhat <- yosoku(a, b, d3$X[i])
603      zansa2[i] <- zansa(d3$Y[i], Yhat)
604  }
605  zansa2
606  sum.zansa2 <- sum(zansa2)
607  sum.zansa2
608
609  plot(d3$Y ~ d3$X, ylim=c(0, 15), main=c("a=2, b=2 の時の予測結果"))
610  abline(a=2, b=2)
611
612
613  # 一つの関数にまとめる
614  OLS <- function(para) {
615      zansa2 <- numeric()
616      for ( i in 1:nrow(d3) ) {
617          Yhat <- yosoku(para[1], para[2], d3$X[i])
618          zansa2[i] <- zansa(d3$Y[i], Yhat)
619      }
620      sum.zansa2 <- sum(zansa2)
621      return(sum.zansa2)

```

```
622      }
623
624  OLS(c(2, 2))
625
626  # optim で最適化
627  optim(c(2, 2), OLS)
628
629  # 答え合わせ
630
631  lm.model <- lm(Y ~ X, data=d3)
632  lm.model
633
634  # 検定
635  anova(lm.model)
636
637
638  ######
639  # 実はこの関数を使った方が早い
640  # これを写経する必要はありません
641  # しかし、実務上はこちらのプログラムを使った方が計算速度が速くなります。
642
643  OLS2 <- function(para){
644    Yhat <- yosoku(para[1], para[2], d3$X)
645    zansa2 <- zansa(d3$Y, Yhat)
646    sum.zansa2 <- sum(zansa2)
647    return(sum.zansa2)
648  }
649
650  OLS2(c(2, 2))
651
652  optim(c(2, 2), OLS2)
653
654  # 関数の仕組み
655
656  bai <- function(x) {
657    return(2*x)
```



```
694
695 Effect <- sum(heihou) / 1
696 Effect
697
698 # 予測残差の大きさ = ノイズの大きさ
699 gosa <- zansa(d3$Y, yosokuti)
700 gosa
701
702 Error <- sum(gosa) / 3
703 Error
704
705 # 分散分析
706 F.value <-Effect / Error
707 F.value
708
709 1 - pf(F.value, 1, 3)
710
711 anova(lm.model)
712
713
714 # おまけ NULL モデル(ナイーブ予測)との比較
715 lm.model.NULL <- lm(Y ~ 1, data=d3)
716 predict(lm.model.NULL)
717 mean(d3$Y)
718
719 anova(lm.model.NULL, lm.model)
720 anova(lm.model)
721
722
723
```

```
724 #=====
725 # パラメトリックブートストラップ検定
726 #=====
727
728 # ★★★★★★★★★★★★★★★★★★
729 # ★            写経はじめ    ★
730 # ★★★★★★★★★★★★★★★★★★
731
732 # 回帰分析の時に使ったデータの再掲
733 d3 <- data.frame(
734     Y = c(3, 5, 4, 6, 7),
735     X = 1:5
736 )
737
738 # モデルの作成
739 lm.model.NULL <- lm(Y ~ 1, data=d3)
740 predict(lm.model.NULL)
741
742 lm.model <- lm(Y ~ X, data=d3)
743 anova(lm.model)
744
745
746 # シミュレーションデータの作成
747 # NULL モデルが正しいと仮定してシミュレーションする
748
749 simulate(lm.model.NULL, 1)
750
751 sim.data <- cbind(
752     simulate(lm.model.NULL, 1),
753     1:5)
754
755 colnames(sim.data) <- c("Y", "X")
756 sim.data
757 mean(sim.data$Y)
758 plot(sim.data$Y ~ sim.data$X)
759
```

```
760
761
762 # パラメトリックブートストラップ検定
763 # 本来は N.sim=10000 くらいで行うことが望ましいが
764 # PC のスペックに自信のない方は N.sim=1000 に変えてください
765 # apply 系の関数を使うともっと早くできます。興味がある方は試してみてください。
766
767 set.seed(1)
768 Nsim <- 10000
769 sim <- simulate(lm.model.NULL, Nsim)
770 sim.F.value <- numeric()
771 for(i in 1:Nsim){
772   sim.data <- cbind(sim[i], 1:5)
773   colnames(sim.data) <- c("Y", "X")
774   model <- lm(Y ~ X, data=sim.data)
775   sim.F.value[i] <- summary(model)$fstatistic[1]
776 }
777
778 sim.F.value
779 max(sim.F.value)
780
781 # 見やすいヒストグラムを作る
782 F.value <- subset(sim.F.value, sim.F.value < 20)
783
784 # F 比のヒストグラム
785 hist(F.value, xlim=c(0, 20))
786
787 # 0~1 の範囲内にある F 比は 10000 個中 6000 個くらい
788 # 1~2 の範囲内にある F 比は 10000 個中 1400 個くらい
789 # 2~3 の範囲内にある F 比は 10000 個中 650 個くらい
790
791 # ヒストグラムの数字バージョン
792 stem(F.value)
793
794 # 実際の F 値と比較
795 # この F 比よりも大きな F 比が出た回数を数えればよい
```

```
796 lm.model <- lm(Y ~ X, data=d3)
797 lm.model
798 summary(lm.model)
799 summary(lm.model)$fstatistic[1]
800
801 # 数え方
802 # 論理演算
803 2<3
804 2>3
805
806 # 相手がたくさんあっても大丈夫
807 c(1,2,3,4,5,6,7,8,9,10)>6
808
809 # TRUE の個数がわかる
810 sum(c(1,2,3,4,5,6,7,8,9,10)>6)
811
812
813 # 実データから計算された F 比
814 summary(lm.model)$fstatistic[1]
815
816 # PB 検定の p 値
817 sum(sim.F.value>=summary(lm.model)$fstatistic[1]) / Nsim
818
819 # 理論的に計算された p 値とほぼ等しい事を確認
820 anova(lm.model)
821
822
823 # 偉人の式とシミュレーションの比較
824 # pf は F 比がある値を下回る確率を表す
825 # df は F 比の確率分布を表す
826 # 確率分布とは、面積が 1 になるように標準化されたヒストグラムのこと
827 x <- seq(0, 20, by=0.1)
828 F <- df(x, 1, 3)
829 F
830
831 hist(F.value,xlim=c(0,20),ylim=c(0,1),prob=T,breaks=50)
```



```
868 #=====
869 # 正規線形モデルの再確認
870 #=====
871
872 # ★★★★★★★★★★★★★★★★
873 # ★ 写経はじめ ★
874 # ★★★★★★★★★★★★★★★★
875
876 # データの読み込み
877 # setwd()
878 data0 <- read.csv("data0.csv")
879 head(data0)
880
881 # 分散分析・回帰分析で解説したように、両者はやってることがほぼ同じ
882 # 説明変数がカテゴリデータか定量データかということは気にする必要がない
883 summary(data0)
884 pairs(data0)
885
886 # モデルの作成
887 # 応答変数 : length
888 # 説明変数 : food & medicine
889 # 興味のある対象（応答変数 length）は説明変数によってどれくらい変化するか？
890 model0 <- lm(length ~ food + medicine, data = data0)
891 model0
892
893 # 分散分析結果の再確認
894 # すべての出力の意味が分かるはず
895 # 変数を加えることによって予測残差は“有意に”減ったか？
896 anova(model0)
897
898 # 要約結果の再確認
899 summary(model0)
900
901 # Call: モデルの式
902 # Residuals: 残差
903 # Coefficients: 係数。t検定によるパラメタの有意性の検定結果
```



```
940 # 引っ張られた 100 本の線のうち、95 本の線が中に入る区間
941 sinrai <- predict(
942   model.sim, newdata=new.sim, se.fit=T, interval="confidence")
943 sinrai
944
945 plot(y.sim ~ x.sim)
946 lines(sinrai$fit[,1] ~ new.sim$x.sim, lwd=2)
947 lines(sinrai$fit[,2] ~ new.sim$x.sim, lwd=2, col=4)
948 lines(sinrai$fit[,3] ~ new.sim$x.sim, lwd=2, col=2)
949
950 # 予測区間
951 # 線を引く
952 # その後、同じ確率分布から得られたデータを例えれば 100 個集める。
953 # その 100 個のデータのうち、95 このデータはこの範囲内に収まる
954 # 予測値から左右対称になっていることに注意
955
956 yosoku <- predict(
957   model.sim, newdata=new.sim, se.fit=T, interval="predict")
958 yosoku
959 yosoku$residual.scale
960 # ↑この値が Ysim のノイズの標準偏差 (10) に近いことを確認
961
962 plot(y.sim ~ x.sim)
963 lines(yosoku$fit[,1] ~ new.sim$x.sim, lwd=2)
964 lines(yosoku$fit[,2] ~ new.sim$x.sim, lwd=2, col=4)
965 lines(yosoku$fit[,3] ~ new.sim$x.sim, lwd=2, col=2)
966
967
968 # 両方載せたグラフ
969 plot(y.sim ~ x.sim)
970 lines(sinrai$fit[,1] ~ new.sim$x.sim, lwd=2)
971 lines(sinrai$fit[,2] ~ new.sim$x.sim, lwd=2, col=4)
972 lines(sinrai$fit[,3] ~ new.sim$x.sim, lwd=2, col=2)
973 lines(yosoku$fit[,2] ~ new.sim$x.sim, lwd=2, col=4, lty=2)
974 lines(yosoku$fit[,3] ~ new.sim$x.sim, lwd=2, col=2, lty=2)
975
```

```
976
977  =====
978  # 中心極限定理
979  =====
980
981  # サイコロの合計値は正規分布になるか？
982  N.sample <- 1
983  saikoro <- sample(1:6, size=N.sample, replace = T)
984  saikoro
985
986  # 一回サイコロを投げることを 1000 回実行
987  N.sim <- 1000
988  kekka <- numeric()
989  for(i in 1:N.sim){
990      kekka[i] <- sample(1:6, size=N.sample, replace = T)
991  }
992  kekka
993  hist(kekka, breaks=0:6)
994
995
996  # 五回サイコロを投げることを 1000 回実行
997  N.sample <- 5
998  N.sim <- 1000
999  kekka <- numeric()
1000  for(i in 1:N.sim){
1001      kekka[i] <- sum(sample(1:6, size=N.sample, replace = T))
1002  }
1003  kekka
1004  hist(kekka, breaks=min(kekka):max(kekka) )
1005
1006
1007  # 十回サイコロを投げることを 1000 回実行
1008  N.sample <- 10
1009  N.sim <- 1000
1010  kekka <- numeric()
1011  for(i in 1:N.sim){
```

```
1012     kekka[i] <- sum(sample(1:6, size=N.sample, replace = T))
1013 }
1014 kekka
1015 hist(kekka, breaks=min(kekka):max(kekka))
1016
1017
1018 # 十回サイコロを投げることを 10000 回実行
1019 N.sample <- 10
1020 N.sim <- 10000
1021 kekka <- numeric()
1022 for(i in 1:N.sim){
1023     kekka[i] <- sum(sample(1:6, size=N.sample, replace = T))
1024 }
1025 kekka
1026 hist(kekka, breaks=min(kekka):max(kekka))
1027
1028 # 単体では全く正規分布ぽくないサイコロの目が
1029 # 合計値をとると正規分布によく似た確率分布（ヒストグラム）になっている。
1030
1031
1032
1033
```

```
1034
1035 #=====
1036 # モデル選択と AIC
1037 #=====
1038
1039 # ★★★★★★★★★★★★★★★★★★
1040 # ★ 写経はじめ ★
1041 # ★★★★★★★★★★★★★★★★★★
1042
1043 # データの読み込み
1044 # ディレクトリを変更した後で実行する
1045 data1 <- read.csv("data1.csv")
1046
1047 # データの確認
1048 head(data1)
1049 names(data1)
1050 levels(data1$option1)
1051 levels(data1$option2)
1052 summary(data1)
1053
1054 pairs(data1, panel=panel.smooth)
1055
1056
1057 # 4つの説明変数を全て加えたモデルを作成
1058 lm.model.1 <- lm(Y ~ ., data=data1)
1059 summary(lm.model.1)
1060 anova(lm.model.1)
1061
1062 # option2が不要なようなので、これを削除したモデルを作成
1063 lm.model.2 <- update(lm.model.1, ~ .-option2)
1064 lm.model.2
1065
1066 # 変数を切った時、予測誤差は有意に増加したか？
1067 anova(lm.model.1, lm.model.2)
1068
1069 # option2を切ったモデルで再度分散分析
```

```
1070  anova(lm.model.2)
1071
1072
1073  # x1 もいらなきそうだったので削除する
1074  lm.model.3 <- update(lm.model.2, ~ .-x1)
1075
1076  # 変数を切った時、予測誤差は有意に増加したか？
1077  anova(lm.model.2, lm.model.3)
1078
1079  # 二つの変数を削除した後のモデルでもう一度分散分析
1080  anova(lm.model.3)
1081
1082
1083  # 念のため x2 も削除してみる
1084  lm.model.4 <- update(lm.model.3, ~ .-x2)
1085
1086  # x2 という変数がなくなった時、予測残差は有意に増えたか？
1087  anova(lm.model.3, lm.model.4)
1088
1089  # Best なモデルは lm.model.3 に決定
1090  best.model.anova <- lm.model.3
1091  summary(best.model.anova)
1092
1093  # ベストじゃない変数全部入りモデルとの比較
1094  summary(lm.model.1)
1095
1096
1097  # AIC によるモデル選択
1098  # install.packages("MuMIn") をした後で実行
1099  library(MuMIn)
1100
1101  # 変数を抜き差しして、すべての変数の組み合わせにおいて AIC を算出して比較する
1102  list <- dredge(lm.model.1, rank="AIC")
1103  list
1104
1105  all.model <- get.models(list)
```



```
1119 #=====
1120 # 8 章 Type II ANOVA
1121 #=====
1122
1123 # ★★★★★★★★★★★★★★★★★★
1124 # ★ 写経はじめ ★
1125 # ★★★★★★★★★★★★★★★★
1126
1127
1128 # データの読み込み
1129 # ディレクトリを変更した後で実行する
1130 data2 <- read.csv("data2.csv")
1131
1132 # データの確認
1133 # 卖上 (sell) を上げるためにどうすればよいかを知りたい、という設定
1134 # experience : チラシを配る店員さんの経験年数
1135 # n.sheets : 配れたチラシの枚数
1136 # time : チラシを配った時間帯 (昼か夜か)
1137 # sex : チラシを配った店員さんの性別
1138 head(data2)
1139 pairs(data2, panel=panel.smooth)
1140 summary(data2)
1141
1142 # やってはいけない解析方法
1143 # 統計モデルを使わずに、一つ一つ検定していく
1144
1145 # 経験・チラシ配布枚数は売り上げに貢献しているか? @t 検定
1146 summary(lm(sell ~ experience, data=data2))
1147 summary(lm(sell ~ n.sheets, data=data2))
1148
1149 # 経験・チラシ配布枚数は売り上げに貢献しているか? @分散分析
1150 anova(lm(sell ~ experience, data=data2))
1151 anova(lm(sell ~ n.sheets, data=data2))
1152
1153 # チラシを配った人の性別と、配った時間は売り上げに貢献しているか? @t 検定
1154 t.test(data2$sell ~ data2$sex)
```

```
1155 t.test(data2$sell ~ data2$time)
1156
1157
1158 # 4つの変数を一つにまとめて統計モデリングする
1159 sell.model1 <- lm(sell ~., data = data2)
1160
1161 # 係数が変数を一つずつ入れた時と大きく異なっていることに注意
1162 # 特に experience の値がおかしい
1163 sell.model1$coef
1164
1165 # 分散分析する。
1166 anova(sell.model1)
1167
1168 # 各係数を t 検定する
1169 summary(sell.model1)
1170
1171 # 説明変数の関係性
1172 par(mfrow=c(1,2))
1173 plot(data2$n.sheets ~ data2$experience)
1174 plot(data2$sex ~ data2$time)
1175 par(mfrow=c(1,1))
1176
1177
1178
```

```
1179 #=====
1180 # Type II ANOVA
1181 #=====
1182
1183 # install.packages("car") をしたあとに実行
1184 library(car)
1185
1186 # 普通の分散分析(Type I ANOVA)
1187 anova(sell.model1)
1188
1189 # Type II ANOVA
1190 Anova(sell.model1,type=c("II"))
1191
1192 # 性別はいらなさそうだったので切る
1193 sell.model2 <- update(sell.model1, ~.-sex)
1194 anova(sell.model1, sell.model2)
1195
1196 Anova(sell.model2)
1197
1198 # さらに経験もいらなさそうだったので切る
1199 sell.model3 <- update(sell.model2, ~.-experience)
1200 anova(sell.model2, sell.model3)
1201
1202 Anova(sell.model3)
1203
1204 # これがベストモデル
1205 summary(sell.model3)
1206
1207
1208 # AIC
1209
1210 library(MuMIn)
1211 sell.model.list <- dredge(sell.model1, rank="AIC")
1212 sell.model.list
1213
1214 all.model.sell <- get.models(sell.model.list)
```



```
1224 #=====
1225 # 交互作用
1226 #=====
1227
1228 # ★★★★★★★★★★★★★★★★
1229 # ★          写経はじめ      ★
1230 # ★★★★★★★★★★★★★★★★
1231
1232
1233 # データの読み込み
1234 # ディレクトリを変更した後で実行する
1235 data3 <- read.csv("data3.csv")
1236
1237
1238 # データの確認
1239 head(data3)
1240 summary(data3)
1241 pairs(data3, panel=panel.smooth)
1242
1243 # 気温や天気がビールの利益に与える影響は？
1244 # beer:ビールによる利益
1245 # temperature:気温
1246 # weather:天気
1247
1248 # ダメな分析手法、変数を別個に入れて検定
1249 summary(lm(beer ~ temperature, data=data3))
1250 t.test(data3$beer ~ data3$weather)
1251
1252 # 変数を一緒にモデリング
1253 model.beer0 <- lm(beer~, data=data3)
1254 library(car)
1255 Anova(model.beer0)
1256
1257 # 晴れの日、雨の日でデータを分割
1258 fine <- subset(data3, weather=="fine")
1259 rain <- subset(data3, weather=="rain")
```

```
1260
1261 # お天気別気温の影響
1262 par(mfrow=c(1,2))
1263 plot(fine$beer~fine$temperature, main="晴れ")
1264 plot(rain$beer~rain$temperature, main="雨")
1265 par(mfrow=c(1,1))
1266
1267 # 交互作用を入れる
1268 # model の formula
1269 # どれを使っても、今回は、結果は同じ
1270 # + で区切られたのが主効果（交互作用じゃないもの）
1271 # : でつながったものが交互作用
1272 lm(beer ~ temperature + weather + temperature:weather, data=data3)
1273 lm(beer ~ temperature * weather , data = data3)
1274 lm(beer ~ (.)^2, data = data3)
1275
1276 # 交互作用入りのモデルを使って再度検定
1277 model.beer <- lm(beer ~ (.)^2, data = data3)
1278 summary(model.beer)
1279
1280 # anova
1281 anova(model.beer)
1282
1283 # Type II ANOVA
1284 Anova(model.beer, type="II")
1285
1286 # weather の部分が Type I ANOVA と全くいっしょ
1287 # 主効果の検定をする際に、交互作用を無視してしまっている。
1288
1289 # Type III ANOVA
1290 # 交互作用が入っている時はこっちを使った方がよいかも
1291 # 交互作用に着目する場合は Type III ANOVA を使うべき
1292 # 主効果の影響を主に見ているのならば、Type II ANOVA のままでよいという説もある。
1293 Anova(model.beer, type="III")
1294
1295 # 3つのANOVAの比較
```

```
1296 anova(model.beer)
1297 Anova(model.beer, type="II")
1298 Anova(model.beer, type="III")
1299
1300 # 交互作用入りモデルの係数
1301 model.beer$coef
1302
1303 # AICでモデル選択
1304 library(MuMIn)
1305 dredge(model.beer, rank="AIC")
1306
1307
1308 # 結果の図示
1309 plot(data3$beer ~ data3$temperature,
1310       col=c(1,2)[data3$weather],
1311       pch = 16,
1312       xlab = "気温",
1313       ylab = "ビールによる利益",
1314       main = "お天気別ビールによる利益",
1315       font.lab = 2
1316 )
1317 legend("bottomright", pch = 16, col = c(1,2), legend = c("fine", "rain"))
1318
1319 # 回帰直線を引く
1320 # 各係数の使い方に注意
1321 abline(
1322   model.beer$coef[1],
1323   model.beer$coef[2],
1324   lwd = 2)
1325 abline(
1326   model.beer$coef[1] + model.beer$coef[3],
1327   model.beer$coef[2] + model.beer$coef[4],
1328   lwd = 2, col = 2)
1329
1330 # 晴れの日の傾きと切片
1331 model.beer$coef[1]
```



```
1368 # subset の意味
1369 data
1370 subset(data, data$sisaku_B == "a.not")
1371 subset(data, data$sisaku_B == "act")
1372
1373 #####
1374 # もっとわかりやすい（が、ほんとはダメな）例
1375
1376 data <- data.frame(Y = c(1,5,3,10),
1377   sisaku_A = rep(c("a.not","act","a.not","act"),1),
1378   sisaku_B = rep(c("a.not","a.not","act","act"),1))
1379
1380 data
1381
1382 model <- lm(Y ~ (.)^2, data=data)
1383 model$coef
1384
1385 model.2 <- lm(Y ~ sisaku_A, data=data, subset=(sisaku_B == "a.not"))
1386 model.2$coef
1387
1388 model.3 <- lm(Y ~ sisaku_A, data=data, subset=(sisaku_B == "act"))
1389 model.3$coef
1390
1391
1392 model$coef
1393 model.2$coef
1394 model.3$coef
1395
1396 subset(data, data$sisaku_B == "a.not")
1397 subset(data, data$sisaku_B == "act")
1398
1399 # 本当はこんなデータ使ってはいけない。あくまで、わかりやすさを優先した解説
1400 anova(model)
1401
1402 # お疲れ様でした
```