

データ解析のための統計モデリング (2008 年 10-11 月)

全 5 (+2) 回中の第 2 回 (2008-10-30)

# さまざまな確率分布と最尤推定

久保拓弥 kubo@ees.hokudai.ac.jp

<http://hosho.ees.hokudai.ac.jp/~kubo/ce/EesLecture2008.html>

この講義のーとが「データ解析のための統計モデリング入門」として出版されました!

<http://hosho.ees.hokudai.ac.jp/~kubo/ce/IwanamiBook.html>

まちがいを修正し, 詳しい解説・新しい内容をたくさん追加したものです

## 今日のもくじ

1. いきなり R & ポアソン分布で統計モデリング ?!	2
2. ポアソン分布とは何か?	6
3. 統計モデルの中核概念: 乱数発生と推定	9
4. ポアソン分布のパラメーター $\lambda$ を最尤推定	11
5. 二項分布で「あり・なし」データをモデル化	14
6. 正規分布についてもちょっとばかり	15
7. 統計モデルにどの確率分布を使えそうか	17

「生態学の統計モデリング」第 2 回目です。R のインストールはうまくできましたか? 今日はこの R を動かしながらハナシをすすめていきたいと思えます。

前回は,

- データ解析とは, 観測データをうまく説明できるような統計モデルの構築 (統計モデリング) にほかならない
- 統計モデルの基本部品は確率分布である

といったことを述べてみました。今日はその基本部品たる確率分布とやらを説明してみたいと考えています。

確率分布, さてさてどう説明したものでしょうか.....うーむ, ハナシをたいへん乱暴にすすめて恐縮なんですけれど, 「確率とは何か? 確率分布の厳

## データ解析は統計モデリング

- 統計モデルは観測データのパターンをうまく説明できるようなモデル
- 基本的部品: 確率分布 (とそのパラメーター)
- データにもとづくパラメーター推定, あてはまりの良さを定量的に評価できる

2008-10-30

4 / 7

## 図 1: 統計モデルとは何か?

密な定義は?」ははぶいてしましましょう.....というのも、ホントにゼロに近いところから出発して話していくと、おそらく聴かされてるほうは面倒で退屈でわかりにくい、と思うわけですよ。<sup>1</sup>

こういうふうに乱暴に始めるのはいいんだけど、どういう確率分布をもちだすのが説明に都合よろしいでしょうかねえ.....それではポアソン分布の説明から始めてみましょうか。前回も述べましたが、この講義ではポアソン分布・二項分布・正規分布という三種類の確率分布だけをおつかうことにしよう、と計画しています。<sup>2</sup>

## 1. いきなり R & ポアソン分布で統計モデリング?!

いきなり統計ソフトウェア R なんかを起動してそこに何かデータが入っていたとしましょう<sup>3</sup>.....そうですね、私が誰かから R のデータ<sup>4</sup>をもらってしまって、それは「data と名づけられた<sup>5</sup> vector object」になっていてこれは「ある人が毎日うけとったメール (e-mail) の数、50 日ぶん」である、と。何だかよくわかりませんが、とりあえず R でそのデータを開いてみましょう。<sup>6</sup>

```
> data
[1] 2 2 4 6 4 5 2 3 1 2 0 4 3 3 3 3 4 2 7 2 4 3 3 3 4
[26] 3 7 5 3 1 7 6 4 6 5 2 4 7 2 2 6 2 4 5 4 5 1 3 2 3
```

これで意味不明ぎみだった「data と名づけられた vector object」なるものの正体が少しわかったと思います。これは 0 個, 1 個, 2 個, ... などと数えられるデータ、つまり カウントデータ (ここでは一日に受けとったメールの本数) になっていて、それが vector なる一次元の構造をもつデータ構造

1. とはいえ、このあたり、前回示したような統計学教科書などあとからながめてみて、自分でおぎなえるところはおぎなってみてください。

2. 今回は全 7 回から全 5 回になったので、正規分布を使った統計モデリングはほとんど登場しません。

3. さすがにわれながら不条理な状況設定だと認識しております。

4. これはホントの観測データではなく、講義用の架空データです。どうやって作るかは後述。

5. データを入力したヒトが勝手につけた名前、と考えてください

6. 表示の左側の [1] だの [26] だのはそれぞれ「このすぐ右にあるデータは 1 番目のデータです」「このすぐ右にあるデータは 26 番目のデータです」を示しています。

に格納されていたらしい, その vector object は data と名づけられている, という状況におかれているわけです.

上の R のユーザーインターフェイス上で示されていることは, ユーザーが data と R のコマンドプロンプト<sup>7</sup> 上で入力すると<sup>8</sup>, 「data の内容はこうです」が示されています. たしかに 50 個の整数からなっているように見えますね.

この data をデータ解析しなければならぬ状況におかれている, とします. さっそく R を使って調べてみましょう.<sup>9</sup>

```
> length(data) # data にはいくつのデータが含まれるのか?
[1] 50

> summary(data)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 0.00   2.00   3.00   3.56   4.75   7.00
```

length() 関数によってこの data に含まれるデータ数は 50 個だと確認ができ, summary() 関数によって, 平均値や最小値・最大値・四分位数などがわかります.

上の summary(data) の読みかたを少し説明してみましょう. Min. と Max. はそれぞれ data 中の最小値・最大値です. また 1st Qu., Median, 3rd Qu. はそれぞれ data を小さい順にならべたときの 25%, 50%, 75% 点の値です.

<sup>10</sup> Median は中央値・中位値と呼ばれる推定値ですね. そして (標本) 平均値 (Mean) が 3.56 である, と. ついでながら (標本) 分散 (variance) はこのように計算できます.

```
> var(data)
[1] 2.9861
```

この標本分散の推定値は data の「ばらつき具合」を表現しています. なお (標本) 標準偏差 (standard deviation; SD) とは分散の平方根です. R ではこんなふうに計算できます.

```
> sqrt(var(data))
[1] 1.7280
> sd(data)
[1] 1.7280
```

ここで重要なのは平均値なんかだけをながめてわかったような気分になら

7. R 上の > に命令を入力すると R がそれを実行してくれる, というインターフェイスのことです.

8. 自動的に print(data) が実行されています.

9. R では # のうしろから行末までに書かれてることは読み飛ばされます. # はコメントマークと呼ばれています.

10. この値は観測値そのものではなく, たいていの場合, 補間された値になっています.

ないことです。解析しなければならないデータはもっと「見えるカタチ」にしましょう。たとえば「5 通うけとったのは 50 日中の何日だったのか?」といったことも調べる必要があります。このような度数分布を R で調べる方法はたくさんあるのですが、たとえば一番簡単には、<sup>11</sup>

```
> summary(as.factor(data))
 0  1  2  3  4  5  6  7
 1  3 11 12 10  5  4  4
```

11. `summary()` 関数の挙動が先ほどと異なってますよね。「どうしてこうなってしまうんだろう」と苦悩してくださいね :-)

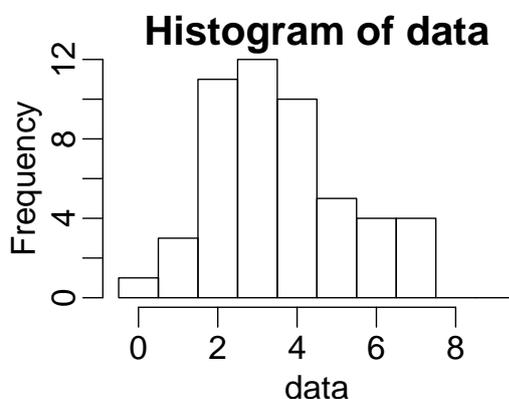
とすればよいだけです。あるいはもっと簡単に

```
> table(data)
 0  1  2  3  4  5  6  7
 1  3 11 12 10  5  4  4
```

としてもいいです。

前回に皆さんに「実践してもらいたいこと」としてあげた「観測データの図をたくさん作ろう」にしたがって、この `data` を度数分布図 (histogram) として図示してみましょう。<sup>12</sup>

```
> hist(data, breaks = seq(-0.5, 8.5, 1))
```



12. `seq(...)` の意味がわからないヒトは自分の R で `seq(-0.5, 8.5, 1)` を入力して実行してみてください。

R はホントに便利ですね。

ここまででわかったことは、某氏が毎日うけとるメールの本数 (データは 50 日ぶん) は

- 一日にうけとるメール数の平均値は 3.56 通
- しかしそれより多い日も少ない日もある

ということで、何か確率論的な (probablistic) モデルでこの現象を説明できそうだと、つまり統計モデリング可能だろう、とハナシをつなげたいわけです。

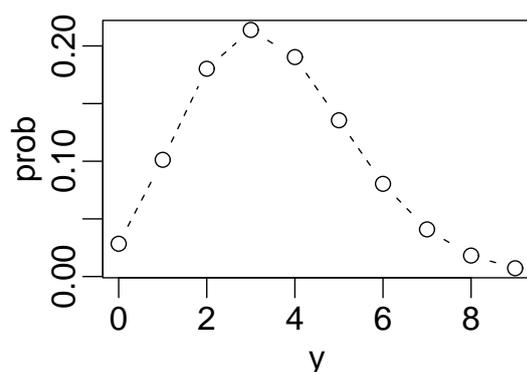
この観測されたパターンはどういった確率分布によって表現できるのでしょうか?

こういうカウントデータの統計モデリング, その基本部品となる確率分布としては(とりあえず) ポアソン分布 (Poisson distribution) を使ってみる<sup>13</sup> というのがよくある統計モデリングの手口です. いささか強引ですが, ここでは「そういうものだ」というコトにしてハナシを進めていくことにします.

ポアソン分布っていうのはいったいぜんたい何なのか? いきなり図でも示して説明してみましようか.

たとえば, 平均 3.56 のポアソン分布<sup>14</sup> とはこういうものです.

```
> y <- 0:9
> prob <- dpois(y, lambda = 3.56)
> plot(y, prob, type = "b", lty = 2)
```



```
> data.frame(y, prob) # y と prob を表にしてみる
```

```
  y  prob
1  0 0.02843882
2  1 0.10124222
3  2 0.18021114
4  3 0.21385056
5  4 0.19032700
6  5 0.13551282
7  6 0.08040427
8  7 0.04089132
9  8 0.01819664
10 9 0.00719778
```

上の図表で示そうとしているコトは何なのかを説明してみると,

- ポアソン分布は, 一日にくるメール数  $y$  (この  $y$  は  $y \in \{0, 1, 2, 3, \dots\}$ ) と,  $y$  それぞれに対応する確率 (ここでは  $\text{prob}$ ) の組みあわせである
- R では「一日  $y$  通くる確率」である  $\text{prob}$  は  $\text{dpois}(y, \text{lambda} = \text{平$

13. すでにお忘れかもしれませんが、「世の中にはポアソン分布と呼ばれる確率分布がある」なるとうとつなる前提でハナシが展開中, なんです.

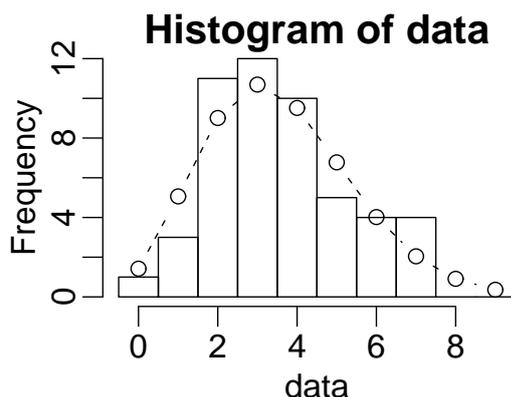
14. これは  $\text{data}$  の平均値を (今のところ特に理由もなく) 使っています.

均値) で計算できる

といったことなのです。だいたいわかるでしょうか? この平均が 3.56 であるポアソン分布だと、一日のメール数がゼロである確率は 0.03 ぐらい、一番確率が高くなるのは一日に 3 通ぐらいくる場合で、その確率は 0.21 ぐらい、ということをお知らせしています。平均 3.56 のポアソン分布なるものが「なんとなく」わかりましたか? <sup>15</sup>

ここで重要なのは、統計モデリングにおいてはこの確率分布 (平均 3.56 のポアソン分布) を使って観測データ (data, つまり一日にくるメール数, 50 日ぶん) にみられるパターンが「説明」できる、と考えていることです。それは図で示すようになります。 <sup>16</sup>

```
> # data, y, prob はすでに定義済みだとします
> hist(data, breaks = seq(-0.5, 8.5, 1))
> lines(y, prob * 50, type = "b", lty = 2) # 50 日ぶん
```



ヒストグラムで示されている観測データと、ポアソン分布による予測がかなり「近い」というかんじですよね? こういう関係を見て、 <sup>17</sup>

- メール来信数の観測データがポアソン分布 (平均 3.56 の場合) の予測とだいたい「あって」いるようにみえる
- このポアソン分布を使った統計モデルによって観測された現象が「説明」(あるいは「表現」) されてるよね

とみなしてしまうのが統計モデリングなのです。

15. なぜ「平均 3.56」という数値を使うのか? それはこのあとの最尤推定の節で説明します。

16. ここでは hist() でヒストグラムを表示させてから、それを消さずに lines() によってポアソン分布による 50 日ぶんの予測値を上げきしています。この「上げき」技法は R による作図のたいへん重要な基本わざです。

17. もちろん実際の統計モデリングはもう少し「客観性」ごときものが要求されます。「ポアソン分布はいいのかもしれないけど、どうして平均 3.56 なの?」「他の分布はありえないわけ?」「『あってる』って何をもってそういうわけ?」 こういった疑問に対応できる方向で、この講義をすすめていきたい、と考えています。

## 2. ポアソン分布とは何か?

ここまでで、観測データと確率分布を使った統計モデルの関係がわかってきた、ということにしましょう<sup>18</sup>。上の「メール数データ」の統計モデリングに使った確率分布、すなわちポアソン分布の性質をもう少し詳しく調べてみましょう。

ポアソン分布がどんなものであるかを規定している確率密度関数 (probabilistic density function) は

$$p(y | \lambda) = \frac{\lambda^y \exp(-\lambda)}{y!}$$

と定義され、左辺の  $p(y | \lambda)$  は<sup>19</sup> (平均値をあらわす) パラメーター  $\lambda$ <sup>ラムダ</sup> がある値のときに、確率変数が  $y$  (例: メール数が 2 のときは  $y = 2$  と考える) である確率、という意味です。その確率というのが右辺で定義されていて、 $y$  とパラメーター  $\lambda$  だけで定義されています。<sup>20</sup>

ついでに説明しとくと  $\exp(-y) = e^{-y}$  で  $e$  は自然対数の底 ( $e = 2.7182\dots$ ) です。 $y!$  は  $y$  の階乗で、たとえば  $4!$  は  $1 \times 2 \times 3 \times 4$  をあらわしています。

さて、ポアソン分布の性質をいくつかあげてみると

- $y = 0, 1, 2, 3, \dots, \infty$  (無限大) の値をとる,  $\sum_{y=1}^{\infty} p(y | \lambda) = 1$
- 分布の平均は  $\lambda$  である ( $\lambda \geq 0$ )
- 分布のカタチを決めるパラメーターは  $\lambda$  だけである
- 分散と平均値は等しい:  $\lambda = \text{平均} = \text{分散}$

最後に「分散」なる用語がでてきました。とりあえず、これは「値のちらばりぐあい」とでも考えてください。<sup>21</sup>

ポアソン分布のパラメーター  $\lambda$  を変化させると、分布のカタチはこのように変化していきます。<sup>22 23</sup>

```
> y <- 0:20
> plot(y, dpois(y, lambda = 3.5), type = "b", lty = 2, pch = 21,
+       ylab = "prob")
> lines(y, dpois(y, lambda = 7.7), type = "b", lty = 2, pch = 23)
> lines(y, dpois(y, lambda = 15.1), type = "b", lty = 2, pch = 24)
> legend("topright", legend = c(3.5, 7.7, 15.1), pch = c(21, 23, 24),
+       title = "lambda", cex = 0.7)
```

(次のページにつづく)

18. だいじょうぶかな?

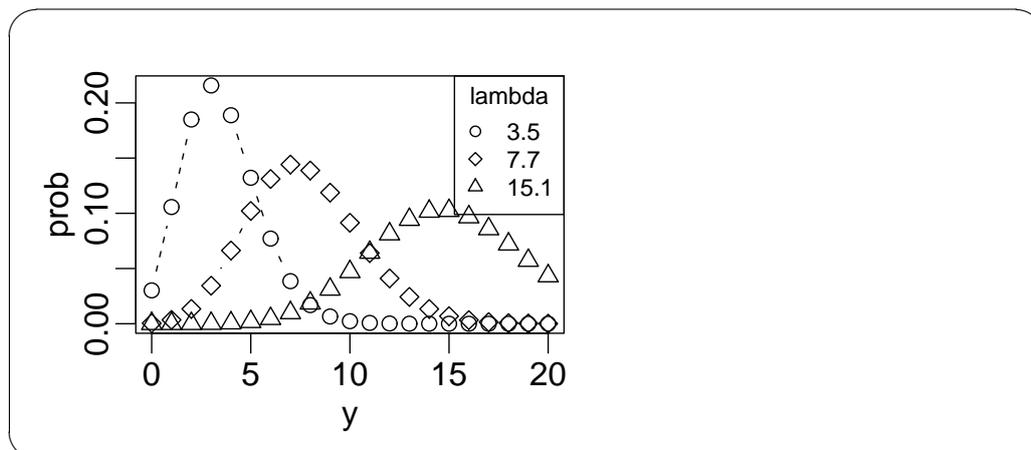
19. めんどくなときは単に  $p(y)$  と書かれたりします。

20. なんでこういう式になるのか? 興味があるヒトはぜひ統計学本などで調べてみてください。

21. またあとで詳しく説明する機会があるかもしれませんが、じつはこの講義では平均とは何かについてきちんと説明してませんね!

22. ここでは `plot()` で「わく」を作らせて、さらに `lines()` による上がり、というわざが使われています

23. R のコマンドプロンプトでカッコを閉じないで改行すると、次の行頭に `>` ではなく `+` があらわれて、前の行のつづきを書くことができます。



ここで統計モデルの部品としてこのポアソン分布が選ばれた理由は

- data に含まれてる値 ( $y$  としましょうか) が  $0, 1, 2, \dots$  といった非負の整数である
- $y$  に下限 (ゼロ) はあるみたいだけど上限はよくわからない
- 標本平均と標本分散<sup>24</sup> がだいたい等しい

24. 標本平均・標本分散とは、それぞれ手もちの観測データの平均値と分散のことです。

```
> mean(data) # data の標本平均を計算
[1] 3.56
> var(data) # data の標本分散を計算
[1] 2.986122
```

といった理由によるものです。

ある確率論的な事象 (event) がポアソン分布になる、というのはどういう条件のもとでおこることなのでしょう? たとえば、ここであつかつてるような「一日にくるメール数」だとすると、メール数がポアソン分布になる必要条件の例としては

- 毎日くるメール数は一定— 曜日・季節による変化があつてはいけない
- 毎日くるメール数は独立— たとえば昨日きたメール数に依存してはいけない
- 一日の中でも各メールの来信は独立— ある日はメールやりとりが頻繁だから増えた、となつてはいけない

といったことがあげられます。これらの条件が満たされていないときはどうしたらよいのでしょうか? 対策はふたつあります。

1. 上の条件を満たしてないみたいだけど，現象の近似的な表現手段としてポアソン分布を使う
2. ポアソン分布では観測された現象が説明できないから，ポアソン分布を少し修正した確率分布を使った統計モデリングにする

対策 2. の例としては，たとえばポアソン分布の平均値が日によって変わっているとしましょう．そうするとメール数の標本分散は標本平均よりかなり大きくなります．こういう場合は負の二項分布 (negative binomial distribution) を使った統計モデルや一般化線形混合モデル<sup>25</sup> を使ってメール数を表現する，という方法もあります．

対策 1. と 2. のどちらをすればよいのでしょうか？ これは状況によるのであり，判断するためにはそのデータについていろいろと調べなければなりません．データがどのようにして得られたのか，どのようにして生成されていると考えられるデータなのか，データの標本平均と標本分散はどういった関係にあるのか，といったことをよく考え R でデータをよく調べる必要があります．ここはデータ解析者の統計モデリングの技術とセンスが問われるところです．<sup>26</sup>

### 3. 統計モデルの中核概念: 乱数発生と推定

ここまでで「ポアソン分布ってのは何となくこんなものらしい」というのがわかったコトにして，つぎにこういった確率分布が統計モデリング・データ解析の中で果たす役割について考えてみましょう．

私たちが何か観測データ，たとえば先ほどから使ってる data

```
> data
[1] 2 2 4 6 4 5 2 3 1 2 0 4 3 3 3 3 4 2 7 2 4 3 3 3 4
[26] 3 7 5 3 1 7 6 4 6 5 2 4 7 2 2 6 2 4 5 4 5 1 3 2 3
```

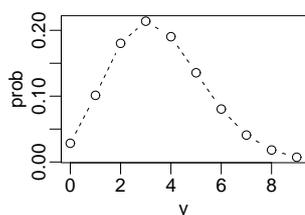
こういう数字の羅列を見たときに，「ああ，これは何か確率論的な現象で統計モデリングによって記述可能だな」と考えることが統計モデリングの第一歩となります．このときにデータ解析者のアタマの中では，図 2 のように考えている，ということになります．研究対象である現象 (生態学であつかうのはおもに自然現象ですが) の実体は確率分布である，<sup>27</sup> と．私たちはこれを観測することによって，乱数 (random number) の集まりを標本 (sample) として得ている，と考えます．確率分布から乱数をとることを sampling とよびます．

25. この講義の第 6 回であつきます 予定でしたが第 6 回の講義はどこかに消えてしまいました．

26. そしてこのあとの講義で述べる「モデル選択」を使えば，どちらの対策が良いのかある程度は「客観的」に判断できるかもしれません．

27. まあ，人間にはなぜだかよくわからない理由で

確率論的な現象 = 確率分布



観測データ (標本, 乱数)

```
2 2 4 6 4 5 2 3 1 2
0 4 3 3 3 3 4 2 7 2
4 3 3 3 4 3 7 5 3 1
7 6 4 6 5 2 4 7 2 2
6 2 4 5 4 5 1 3 2 3
```

図 2: 確率分布から乱数発生, というモデル

逆に, 観測データ (標本) から「そもそもこれを生み出した確率分布はどんなものなんだ?」を決めてみようとするのが推定 (estimation) です (図 3)。

観測データ

```
2 2 4 6 4 5 2 3 1 2
0 4 3 3 3 3 4 2 7 2
4 3 3 3 4 3 7 5 3 1
7 6 4 6 5 2 4 7 2 2
6 2 4 5 4 5 1 3 2 3
```



推定された確率分布

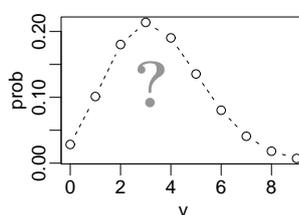


図 3: 観測データから確率分布を推定する

このように

1. ばらつきのあるデータ (標本) はある確率分布  $X$  から生成されたのだろう (自然現象ってのは確率分布  $X$  なのだろう)
2. 私たちは確率分布  $X$  がどのようなカタチをしているのか知らないけれど, データから確率分布  $X$  のカタチが推定できるのではないかな

と考えることが統計モデリングの中核的な考えかた, となります. 確率分布のカタチはパラメーター (parameter) に依存しており, <sup>28</sup> その値をデータから決めてやることをパラメーター推定ともいいます. このことから

- 現象がどのような確率分布で説明されそうか
- データからどのようにパラメーターの値を推定するか

この点をはっきりさせることが自然科学のデータ解析にとっては重要になります. <sup>29</sup> しかしながら, 多くの自然科学の論文ではこの点がかなりいいかげんにあつかわれていています. つまりきちんと説明されていませんし, おそらくよく考えられていないのでしょう.

この節の最後に, 今回とりあげている架空データである data つまり「毎日くるメール数, 50 日ぶん」の生成方法について説明してみましよう. こ

28. とこで世の中では「ばらめとりつく検定」といえば正規分布を仮定した統計モデルの検定のこと, てな意味で使われてしまっていたりしますが, これはまったく不適切なものです. どんな確率分布だってパラメーターによってカタチを決められているわけですから.

29. 皆さんが大好きな「検定」なんてのは, これらに比べると瑣末なことだと私は考えています. 検定なるものは推定されたパラメーターをどう比較するかという検討なので, 統計モデリングや推定をきちんとはじめて議論可能となるからです.

これはホントの観測データではなく、私が R の乱数生成機能をつかって作ったものです。R にはじつにさまざまな乱数生成機能がついています。たとえばポアソン分布にしたがう乱数を生成するのが `rpois()` 関数で、`data` はこれを使って生成されました。

```
> data <- rpois(50, lambda = 3.5)
> data
 [1] 2 2 4 6 4 5 2 3 1 2 0 4 3 3 3 3 4 2 7 2 4 3 3 3 4
[26] 3 7 5 3 1 7 6 4 6 5 2 4 7 2 2 6 2 4 5 4 5 1 3 2 3
```

じつは「真の平均値」は 3.56 ではなく 3.5 だったのです。このように観測されたデータから推定された平均値と「真の平均値」は異なります。

乱数生成のおもしろいところは結果が毎回毎回かわるところです。たとえば、

```
> y <- rpois(50, lambda = 3.5); print(y); print(mean(y))
 [1] 5 6 2 5 2 4 3 4 4 6 4 1 3 5 5 1 3 1 5 2 0 2 0 4 3
[26] 2 6 0 2 3 6 3 2 3 1 2 6 4 2 2 3 4 0 6 5 5 4 1 5 5
 [1] 3.24
```

ここでは標本平均値もいっしょに計算していて、上の「データ」だと標本平均値は 3.24 になりました。あと二回ほど同じことをやってみましょうか。

```
> y <- rpois(50, lambda = 3.5); print(y); print(mean(y))
 [1] 1 4 2 2 4 1 3 1 2 5 4 3 5 7 1 1 2 2 3 9 0 3 5 3 5
[26] 3 3 3 4 4 4 4 3 0 4 9 3 1 7 3 2 1 2 4 3 0 6 2 2 2
 [1] 3.14

> y <- rpois(50, lambda = 3.5); print(y); print(mean(y))
 [1] 5 7 2 2 4 3 5 4 4 3 3 7 5 3 5 2 5
[18] 6 2 4 4 2 3 9 1 6 0 3 3 11 1 2 2 3
[35] 3 3 2 2 3 4 3 4 6 5 4 4 4 4 4 2
 [1] 3.76
```

2 回目・3 回目の標本平均値はそれぞれ 3.14 と 3.76 になっていますね。

#### 4. ポアソン分布のパラメーター $\lambda$ を最尤推定

ここまでメール来信数の架空データである `data` の解析として、

1. ばらつきのあるデータ data はポアソン分布から生成されたのだろう
2. 私たちはそのポアソン分布がどういうカタチをしているのか知らないけれど, data の標本平均値 3.56 を使うと (データを生成したもともとの) ポアソン分布のカタチを決めているパラメーター  $\lambda$ <sup>30</sup> が推定できるのではないかな

30. ポアソン分布の  $\lambda$  は平均・分散をあらわすパラメーターでしたね

というふうに統計モデリングしてきました. 1. でポアソン分布を選んだ理由は上で述べてみました. それでは 2. のように「標本平均値 3.56 がなんとなくホントの平均値に『近い』のではないかな?」と考えた理由は何でしょう? このあたりは最尤推定法 (maximum likelihood estimation) の考えかたにもとづいています.<sup>31</sup>

31. 最尤法という場合もあります.

メール来信数データである data の第  $i$  日目の来信数を  $y_i$  としましょう, つまり  $\{y_1 = 2, y_2 = 2, y_3 = 4, \dots, y_{49} = 2, y_{50} = 3\}$  ということです. このときに「平均  $\lambda$  のポアソン分布から  $\{y_i\} = \{y_1, y_2, y_3, \dots, y_{50}\}$  が得られる確率」のことを尤度 (likelihood) とよびます. ここでポアソン分布を

$$p(y | \lambda) = \frac{\lambda^y \exp(-\lambda)}{y!}$$

とすると尤度 (尤度関数) は<sup>32 33</sup>

$$\begin{aligned} L(\lambda | \{y_i\}) &= (y_1 \text{ が } 2 \text{ である確率}) \times (y_2 \text{ が } 2 \text{ である確率}) \\ &\quad \times \dots \times (y_{50} \text{ が } 3 \text{ である確率}) \\ &= p(y_1 | \lambda) \times p(y_2 | \lambda) \times p(y_3 | \lambda) \times \dots \times p(y_{50} | \lambda) \\ &= \prod_{i=1}^{50} p(y_i | \lambda) \\ &= \prod_{i=1}^{50} \frac{\lambda^{y_i} \exp(-\lambda)}{y_i!} \end{aligned}$$

32. この式は「平均  $\lambda$  のポアソン分布のもとで観測データが得られる確率」の意味ですが, 形式的には  $L(\lambda | \{y_i\})$ , つまり「ある観測データのもとでパラメーター  $\lambda$  が得られる確率」というふうに尤度はパラメーターの関数であるように定義しています.

33.  $\prod$  は積の記号です.

となります. このように定義された「データが得られる確率」こと尤度, つまり尤らしさを最大化する  $\lambda$  が「一番よい」 $\lambda$  と考えることにしましょう.

尤度  $L(\lambda | \{y_i\})$  を最大化する  $\lambda$  とはどのようなものなのかを調べてみましょう. 尤度がかけ算のカタチだとたいへん扱いづらいので, 対数変換した対数尤度 (log likelihood) あるいは対数尤度関数になおして計算します.<sup>34</sup>

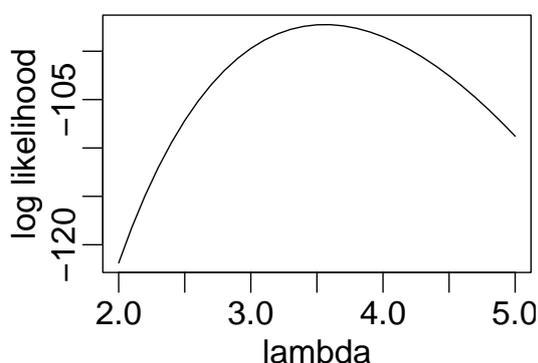
$$\log L(\lambda | \{y_i\}) = \sum_{i=1}^{50} \left\{ y_i \log(\lambda) - \lambda - \sum_{k=1}^{y_i} \log(k) \right\}$$

ここで  $\lambda$  を変化させていって対数尤度  $\log L(\lambda | \{y_i\})$  が最大になる点があるだろう, と考えます. このように対数尤度を最大化する  $\lambda$  を  $\hat{\lambda}$ <sup>ラムダハット</sup> とし

34.  $\log(L)$  は  $L$  の単調増加関数なので,  $L$  の最大化  $\Leftrightarrow \log(L)$  の最大化, となっています.

ます．ここで data が決まっているときの， $\lambda$  と対数尤度の関係を図にしてみましよう．

```
> lambda <- seq(2, 5, 0.1)
> likelihood <- function(lambda) sum(dpois(data, lambda, log = TRUE))
> plot(
+ lambda,
+ sapply(lambda, likelihood),
+ type = "l",
+ xlab = "lambda",
+ ylab = "log likelihood"
+ )
```



この図をみると，対数尤度（縦軸）が最大になるふきんでは対数尤度関数の「傾き」はゼロになっていることがわかります．つまり，対数尤度関数を横軸  $\lambda$  で偏微分したときに

$$\frac{\partial \log L(\lambda \mid \{y_i\})}{\partial \lambda} = \sum_{i=1}^{50} \left\{ \frac{y_i}{\lambda} - 1 \right\} = \frac{\sum_{i=1}^{50} y_i}{\lambda} - 50$$

この  $\log L(\lambda \mid \{y_i\})$  がゼロになるような  $\lambda$  が  $\hat{\lambda}$  なのです．実際に計算してみると

$$\frac{\sum_{i=1}^{50} y_i}{\hat{\lambda}} - 50 = 0$$

なので

$$\hat{\lambda} = \frac{\sum_{i=1}^{50} y_i}{50} = \frac{\text{全部の } y_i \text{ の和}}{\text{データ数}} = \text{データの平均値}$$

となっています．<sup>35</sup> このように  $\{y_i\}$  に具体的な値をいれない状態の  $\hat{\lambda}$  のことを最尤推定量 (maximum likelihood estimator)，さらに  $\{y_1 = 2, y_2 = 2, y_3 = 4, \dots, y_{49} = 2, y_{50} = 3\}$  というふうに具体的に  $y_i$  に値をいれて計算して得た  $\hat{\lambda} = 3.56$  のことを最尤推定値 (maximum likelihood estimate) と言います．

つまりこのような単純な状況でポアソン分布を仮定したときに，data という観測データのもとでは，最尤推定すると平均値 3.56 が最も尤もらしい

35. 残念ながら，パラメータまわりがもっと複雑になってしまう実際のデータ解析ではこんなに簡単に最尤推定量は導出できません．そこで計算機を使った対数尤度最大化によって最尤推定値を数値的にもとめます．次回以降で紹介する一般化線形モデルの推定関数などではこの数値的な最尤法によってパラメータを推定しています．

値である，ということがわかりました。<sup>36</sup>

## 5. 二項分布で「あり・なし」データをモデル化

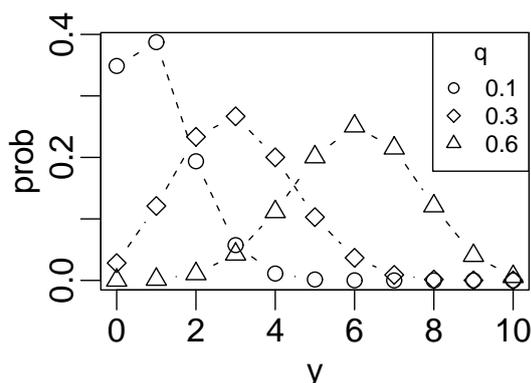
データ数が 0 個, 1 個, 2 個, ... というふうに数えられるカウントデータの統計モデリングにつかえる確率分布はポアソン分布のほかにもいろいろあります。

ポアソン分布はデータが 0 以上のどんな値でもとりうることをあらわしている確率分布ですが，そうではないカウントデータもあります。たとえば  $N$  個のコインを投げたときに  $y$  枚がおもてになった，という観測データでは  $y$  が  $\{0, 1, 2, \dots, N\}$  の値しかとることができません。このときには二項分布 (binomial distribution)

$$p(y | N, q) = \binom{N}{y} q^y (1 - q)^{N-y}$$

を使うと  $y$  が  $\{0, 1, 2, \dots, N\}$  になる確率を計算できます。

```
> y <- 0:10
> plot(y, dbinom(y, 10, prob = 0.1), type = "b", lty = 2, pch = 21,
+       ylab = "prob")
> lines(y, dbinom(y, 10, prob = 0.3), type = "b", lty = 2, pch = 23)
> lines(y, dbinom(y, 10, prob = 0.6), type = "b", lty = 2, pch = 24)
> legend("topright", legend = c(0.1, 0.3, 0.6), pch = c(21, 23, 24),
+       title = "q", cex = 0.7)
```



この二項分布を使った統計モデルは，たとえば生態学のデータ解析では，たとえばある実験処理をして  $N$  個体中  $y$  個体が応答したといった，反応の「あり・なし」データ解析に使います。

二項分布については次の次の講義「第 4 回: 11/10 (月) 一般化線形モデル (GLM) 2」でくわしく紹介します。

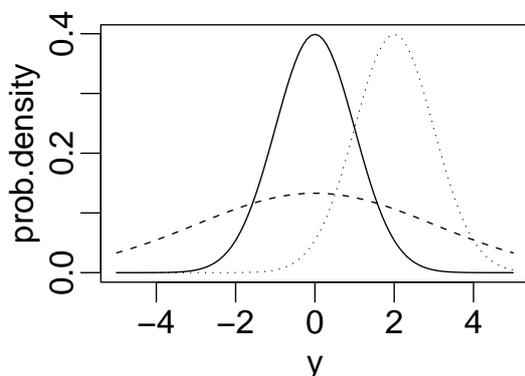
36. 尤度と似て非なるものがあります。準尤度 (quasi likelihood) は overdispersion を表現するために作られ，擬似尤度 (pseudo likelihood) は空間相関など単純な最尤法ではあつかいにくい現象のモデリングのために作られました。どちらも現在はあまり使われていません。Bayesian な方法が普及したので，これらが不要になったからです。

## 6. 正規分布についてもちょっとばかり

よく理解されないままよく使われてしまっている正規分布 (normal distribution; 別名: ガウス分布, Gaussian distribution) <sup>37</sup> はポアソン分布や二項分布とは異なり,  $\{1.5, -3.2, 7.7, \dots\}$  といった連続値のデータをあつかう確率分布です. 正規分布で 0 個, 1 個, 2 個,  $\dots$  といったカウントデータをあつかっているヒトがいますけれど, これはまったく不適切な使いかたです.

正規分布はポアソン分布と同じく平均値のパラメーターをもちます. 平均をあらわすパラメーター  $\mu$  は  $\pm\infty$  の範囲で自由に平均値を変えることができます. <sup>38</sup> . またポアソン分布とは異なり, 値のばらつきをあらわす標準偏差 <sup>39</sup> をあらわすパラメーター  $\sigma$  はゼロ以上なら自由な値を設定できます. つまり, 確率分布の平均だけでなく, ばらつき具合も指定できます.

```
> y <- seq(-5, 5, 0.1)
> plot(y, dnorm(y, mean = 0, sd = 1), type = "l",
+ ylab = "prob.density") # 実線
> lines(y, dnorm(y, mean = 0, sd = 3), lty = 2) # 破線
> lines(y, dnorm(y, mean = 3, sd = 1), lty = 3) # 点線
```



正規分布みたいな連続値の確率分布にはちょっと面倒なところがあって, 上のグラフは確率ではなく確率密度をあらわしています. つまり正規分布の確率密度関数は平均パラメーター  $\mu$  と標準偏差パラメーター  $\sigma$  を使って

$$p(y | \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \mu)^2}{2\sigma^2}\right)$$

というふうに定義されるのですが, これは確率ではなく確率密度なのです.

どういうことかということ たとえば成人男子の平均身長が 170 cm での標準偏差が 5 cm だとします. <sup>40</sup> このときに「ある人の身長が 177 cm から 178 cm にある確率を計算せよ」という問題には

37. というか、「なんでもかんでも正規分布」なヒトたちはほとんど「確率分布とは何か」といったことを考えたことがないのでは? と思えてしまいます.

38. ポアソン分布では平均値パラメーター  $\lambda$  は非負でなければならない, という制約があります.

39. 標準偏差<sup>2</sup> = 分散

40. データにもとづく値ではありません.

```
> pnorm(178, mean = 170, sd = 5) - pnorm(177, mean = 170, sd = 5)
[1] 0.02595737

> dnorm(177.5, mean = 170, sd = 5) * 1.0 # 近似的な確率の簡単計算法
[1] 0.02590352
```

と答えることができます。しかし「ある人の身長が 178 cm ぴったりである確率を計算せよ」という問題だと

```
> pnorm(178, mean = 170, sd = 5) - pnorm(178, mean = 170, sd = 5)
[1] 0
```

確率はゼロになります。ただし確率密度はゼロではありません。

```
> dnorm(178, mean = 170, sd = 5)
[1] 0.02218417
```

正規分布の最尤推定について少しだけ言及してみます。前回の授業で「最小二乗法は最尤推定の一部である」とのべました。最小二乗法は統計モデルの部品として正規分布を使っているときに使えるパラメーター推定法です。この関係について少し調べてみましょう。

たとえばある人間の集団の身長データを  $\{y_i\}$  ( $i = 1, 2, \dots, N$ ) とすると、正規分布をつかった統計モデルの尤度関数は

$$\begin{aligned} L(y_i | \mu, \sigma) &= \prod_{i=1}^N p(y_i | \mu, \sigma) \Delta y \\ &= \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \mu)^2}{2\sigma^2}\right) \Delta y \end{aligned}$$

となります。ここで  $\Delta y$  は「身長に関する一定の微小な幅 (例: 0.1 cm など)」<sup>41</sup> です。というのも正規分布の確率密度  $p(y | \mu, \sigma)$  はそれだけでは「確率」ではなく「(身長の) 幅」をかけて「面積」にしてやらないと確率になりません。<sup>42</sup>

対数尤度関数は

$$\log L(y_i | \mu, \sigma) = -0.5N \log(2\pi\sigma^2) - \frac{\sum_{i=1}^N (y_i - \mu)^2}{2\sigma^2} + N \log \Delta y$$

いま標準偏差パラメーターである  $\sigma$  が  $\mu$  とは無関係ななにか定数だとすると、<sup>43</sup> 対数尤度  $\log L(y_i | \mu, \sigma)$  を最大化するにはうしろの項つまり  $-\frac{\sum_{i=1}^N (y_i - \mu)^2}{2\sigma^2}$  の分子  $\sum_{i=1}^N (y_i - \mu)^2$  を最小に（「二乗誤差の和」を最小

41. 小さければなんでもいいです。

42. 先ほどの `dnorm()` などを使った正規分布の確率計算を参照してください。

43. もちろん  $N \Delta y$  も定数です。

化) するような平均のパラメーター  $\hat{\mu}$  を推定してやればよい, ということがわかると思います.

このことから正規分布を部品とする統計モデルにおける最尤推定法は最小二乗法, <sup>44</sup> が成立しています.

44. ただし  $\sigma$  が一定, といった仮定が必要です.  
ね.

## 7. 統計モデルにどの確率分布を使えそうか

観測データを説明する統計モデル, その基本部品である確率分布はどう選ばよいのでしょうか? とりあえずデータをみたら次の点に注意してみてください.

1. 説明したい量は離散か連続か?
  - 離散: { 生きてる, 死んでる }, カウントデータ, ...
  - 連続: { 0.56, 1.33, 12.4, 9.84, ... }, ...
2. 説明した量の範囲は?
  - $\{0, 1, \dots, N\}$ ,  $\{0, 1, \dots, \infty\}$ ,  $[y_{\min}, y_{\max}]$ ,  $[-\infty, \infty]$ , ...
3. 説明したい量の標本分散 (ばらつき) と標本平均の関係は?
  - 分散  $\approx$  定数, 分散  $\approx$  平均, 分散  $\propto$  平均, 分散  $\propto$  平均<sup>n</sup>, ...

いくつかの確率分布について, それが統計モデルの部品として使いそうな必要条件をざっとみると,

- ポアソン分布: データが離散値, ゼロ以上の範囲, 上限とくになし, 標本平均  $\approx$  標本分散
- 負の二項分布: データが離散値, ゼロ以上の範囲, 上限とくになし, 標本平均  $<$  標本分散
- 二項分布: データが離散値, ゼロ以上で有限の範囲 ( $\{0, 1, 2, \dots, N\}$ )
- 正規分布: データが連続値, 範囲が  $[-\infty, +\infty]$
- 対数正規分布: データが連続値, 範囲が  $[0, +\infty]$
- ガンマ分布: データが連続値, 範囲が  $[0, +\infty]$
- ベータ分布: データが連続値, 範囲が  $[0, 1]$

また R にはさまざまな確率分布が使えるように準備されています。いくつかの関数について、それぞれの乱数生成関数とデータからパラメータを推定する関数を示してみましょう。<sup>45</sup>

	確率分布	乱数生成	パラメータ推定
(離散)	ベルヌーイ分布	<code>rbinom()</code>	<code>glm(family = binomial)</code>
	二項分布	<code>rbinom()</code>	<code>glm(family = binomial)</code>
	ポアソン分布	<code>rpois()</code>	<code>glm(family = poisson)</code>
	負の二項分布	<code>rnbinom()</code>	<code>glm.nb()</code>
(連続)	ガンマ分布	<code>rgamma()</code>	<code>glm(family = gamma)</code>
	正規分布	<code>rnorm()</code>	<code>glm(family = gaussian)</code>

一般化線形モデルの推定関数である `glm()` の使いかたは次回以降の講義であつかうことにします。

今回は「第 3 回: 一般化線形モデル (GLM) 1」, ポアソン分布を仮定した統計モデルを一般化線形モデルとしてあつかう方法について考えてみましょう。

45. ベルヌーイ分布と二項分布はほとんど同じものです。二項分布は  $N$  個の独立した事象がどれも確率  $p$  で生じるときに、いくつかの事象が生じるかをあらかず確率分布で、確率変数は  $\{0, 1, 2, \dots, N\}$  の値をとります。ベルヌーイ分布は二項分布の  $N = 1$  の場合で、確率変数は  $\{0, 1\}$  の値をとります。

## 確率分布と最尤推定

- 確率論的な現象は **確率分布** でうまく表現できそう
- 確率分布のカタチは **パラメーター** できまる
- パラメーターの値は、観測データにもとづいて **最尤推定法** で得られる

この講義で使う確率分布: ポアソン分布, 二項分布, 正規分布

図 4: 本日の要点