

負の二項分布とガンマ分布の関係

プログラマたん bot

2019/8/3

負の二項分布

複数の異なる定義があるので、順に説明する。

[Wikipedia](#)に従い、負の二項分布を、「1回の試行に確率 p で成功するとき、 r 回目の失敗までの、成功回数 y の分布」と定義する。この分布は以下の通りである。

$$\text{NegativeBinomial}(y|r, p) = \binom{y+r-1}{y} (1-p)^r p^y \quad (\text{式 1})$$

これは r 回目に失敗するまでに、 y 回成功し $r-1$ 回失敗する試行の組み合わせと、その確率である。

「[高校数学の美しい物語](#)」の定義は、成功回数 y を試行回数 $z := y + r$ に置き換えて、

$$\binom{y+r-1}{x} = \binom{y+r-1}{r-1} \quad (\text{式 2})$$

であることを用い、

$$\text{NegativeBinomial}(z|r, p) = \binom{z-1}{r-1} (1-p)^r p^{z-r} \quad (\text{式 3})$$

としたものである。

“[Stan Functions Reference](#)” による、別の定義を与える。

$$\text{NegativeBinomial}(y|\alpha, \beta) = \binom{y+\alpha-1}{\alpha-1} \left(\frac{\beta}{1+\beta}\right)^\alpha \left(\frac{1}{1+\beta}\right)^y \quad (\text{式 4})$$

[R の dnbinom 関数の説明](#)によると、

$$\text{NegativeBinomial}(y|size, prob) = \frac{\Gamma(y+size)}{\Gamma(size)y!} p^{size} (1-p)^y \quad (\text{式 5})$$

である。ただしこの定義は、これまで挙げた負の二項分布の定義とは逆で、「1回の試行に確率 p で成功するとき、 r 回目の成功までの、失敗回数 y の分布」である。

そのため本文書では成功と失敗を逆に読み、 p を $1 - p$ に置き換え、以下の式に読み替える。

$$\text{NegativeBinomial}(y|size, prob) = \frac{\Gamma(y + size)}{\Gamma(size)y!} (1 - p)^{size} p^y \quad (\text{式 6})$$

これは式 4 の階乗をガンマ関数におきかえ、パラメータを以下のように置き換えることで得られる。

$$\Gamma(k + 1) = k!, size := \alpha, \beta := 1/p - 1 \quad (\text{式 7})$$

ガンマ分布

[Wikipedia](#)によると、ガンマ分布の確率密度関数は

$$\text{GammaDistribution}(x|shape, rate) = \frac{rate^{shape}}{\Gamma(shape)} x^{shape-1} e^{-rate*x} \quad (\text{式 8})$$

である。慣習的に、 $\alpha := shape$, $\beta := rate$, $\theta := scale = 1/\beta$ と表記する。文献によっては $\beta := scale$ としているものがあり、注意が必要である (例えばC++ の `std::gamma_distribution`)。平均 $\mathbb{E}[x]$ と分散 $Var[x]$ は、以下の通りである。

$$\mathbb{E}[x] = \frac{\alpha}{\beta} = \alpha\theta, Var[x] = \frac{\alpha}{\beta^2} = \frac{\mathbb{E}[x]}{\beta} = \alpha\theta^2 = \mathbb{E}[x]\theta \quad (\text{式 9})$$

平均と分散が既知なら (観測値があれば)、負の二項分布のパラメータ α, β を求めることができる。

$$\beta = \frac{\mathbb{E}[x]}{Var[x]}, \alpha = \mathbb{E}[x]\beta = \frac{\mathbb{E}[x]}{\theta} \quad (\text{式 10})$$

改めて表記上の注意

これまで見たように、負の二項分布とガンマ分布のパラメータ表記には、複数の方法がある。本文書では、以下のように定義する。

- 負の二項分布: 1 回の試行に確率 $prob$ で成功するとき、 r 回目の失敗までの、成功回数 y の分布
- ガンマ分布のパラメータ: $\alpha := shape$, $\beta := rate$, $\theta := scale = 1/\beta$

これ以外にも以下の定義があるため、既存のソフトウェアを使う場合は注意が必要である。

- $prob$ の意味が逆 ($prob := 1 - prob$) で、「1 回の試行に確率 p で成功するとき、 r 回目の成功までの、失敗回数 y の分布」に基づく
- $rate = \beta$ ではなく $scale = \theta$ を与える。もしくは、 β の意味が逆数 ($\beta = scale$) になっている。

表1: 負の二項分布の y の定義と $prob$ の与え方

| y の定義 | ソフトウェア |
|---------|---|
| 成功回数の分布 | Stan |
| 失敗回数の分布 | R dnbinom , C++ , Boost C++ Libraries |

表2: ガンマ分布への $beta$ の与え方

| beta,theta の与え方 | ソフトウェア |
|-----------------|---|
| rate | Stan |
| scale | C++ , Boost C++ Libraries |
| rate または scale | R dgamma |

ガンマ-ポアソン分布

観測値が過分散であるときに、標本には個体差があり、応答変数は個体差に従った分布である、というモデルを適用することがある。例えば発芽率は、 n 個の種の発芽率 $p_i, i \in 1..n$ がある分布に従い、それぞれの種は確率 p_i で発芽する、と仮定することである。詳しくは、“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 (確率と情報の科学)”, 久保拓弥著, 岩波書店, 2012/5 を参照すること。

ここで個体差をガンマ分布で表現し、それぞれの個体差を説明変数とするポアソン分布にした値が観測される、というモデルを考える。

$$y \sim \text{PoissonDistribution}(x), x \sim \text{GammaDistribution}(\text{shape}, \text{rate}) \quad (\text{式 } 11)$$

ガンマ-ポアソン分布を一まとめにして、負の二項分布で表現することができる。証明は[StackExchange](#)にあるが、長いのでここには載せない。

平均と分散

これまでに挙げた文献から、平均と分散を引用する。

NB(α, β) の平均と分散

負の二項分布の定義 [式 4](#) において、平均 $\mathbb{E}[y]$ と分散 $\text{Var}[y]$ は、以下の通りである。

$$\mathbb{E}[y] = \frac{\alpha}{\beta}, \text{Var}[y] = \frac{\alpha}{\beta^2}(\beta + 1) = \mathbb{E}[y](1 + \frac{1}{\beta}) \quad (\text{式 } 12)$$

であり、平均と分散が既知なら (観測値があれば)、負の二項分布のパラメータ α, β を求めることができる。

$$\beta = \frac{\mathbb{E}[y]}{Var[y] - \mathbb{E}[y]}, \alpha = \mathbb{E}[y]\beta \quad (\text{式 13})$$

負の二項分布の代わりに、ガンマ-ポアソン分布 式 11 におけるガンマ分布 式 8 のパラメータで表現することを考える。ポアソン分布 $y \sim \text{poisson}(\lambda)$ の平均 $\mathbb{E}[y]$ と分散 $Var[y]$ は

$$\mathbb{E}[y] = \lambda, Var[y] = \lambda \quad (\text{式 14})$$

であり、StackExchangeによると、独立な期待値と分散の加法性 (共分散が 0) から、 x がガンマ分布 $Gamma(\alpha_{gamma}, \beta_{gamma})$ に従うとき、

$$\mathbb{E}[y] = \lambda = \mathbb{E}[x], Var[y] = Var[Gamma(\alpha_{gamma}, \beta_{gamma})] + \lambda \quad (\text{式 15})$$

$$\mathbb{E}[x] = \mathbb{E}[Gamma(\alpha_{gamma}, \beta_{gamma})] = \frac{\alpha_{gamma}}{\beta_{gamma}} \quad (\text{式 16})$$

$$Var[y] = \mathbb{E}[x] + Var[Gamma(\alpha_{gamma}, \beta_{gamma})] = \frac{\alpha_{gamma}}{\beta_{gamma}} \left(1 + \frac{1}{\beta_{gamma}}\right) = \mathbb{E}[x] \left(1 + \frac{1}{\beta_{gamma}}\right) \quad (\text{式 17})$$

である。つまり、 $\alpha = \alpha_{gamma}, \beta = \beta_{gamma}$ が成り立つ。

NB(size, prob) の平均と分散

負の二項分布の定義 式 6 において、平均 $\mathbb{E}[y]$ と分散 $Var[y]$ は、以下の通りである。R の `dnbinom` 関数の説明とは、 $prob$ の意味が逆であることに注意する ($prob := 1 - prob$)。

$$\mathbb{E}[y] = size \frac{prob}{1 - prob}, Var[y] = size \frac{prob}{(1 - prob)^2} = \frac{\mathbb{E}[y]}{1 - prob} \quad (\text{式 18})$$

であり、平均と分散が既知なら (観測値があれば)、負の二項分布のパラメータ $size, prob$ を求めることができる。

$$prob = 1 - \frac{\mathbb{E}[y]}{Var[y]}, size = Var[y] \frac{(1 - prob)^2}{prob} = \frac{\mathbb{E}[y]^2}{Var[y] - \mathbb{E}[y]} \quad (\text{式 19})$$

負の二項分布の代わりに、ガンマ-ポアソン分布 式 11 におけるガンマ分布 式 8 のパラメータで表現することを考える。先ほどとは逆に、 x がガンマ分布 $Gamma(\alpha_{gamma}, \beta_{gamma})$ に従うとしてボトムアップに考える。

$$\beta_{gamma} = \frac{\mathbb{E}(x)}{Var(x)} = \frac{\mathbb{E}(x)}{Var(y) - \mathbb{E}(x)} = \frac{\mathbb{E}(y)}{Var(y) - \mathbb{E}(y)} = \frac{1}{prob} - 1 \quad (式 20)$$

式 10, 式 19 および 式 20 を見比べると、

$$\alpha_{gamma} = \beta_{gamma} \mathbb{E}[x] = \beta_{gamma} \mathbb{E}[y] = size \quad (式 21)$$

つまり、 $\alpha_{gamma} = size, \beta_{gamma} = 1/p - 1$ が成り立つ。これは 式 4 と 式 6 の関係 (式 7) であった。

再生性

NB(size, prob) の再生性

負の二項分布 $NB(size, prob)$ は $size$ に対して再生性がある。つまり

$$NB(size_A + size_B, prob) = NB(size_A, prob) + NB(size_B, prob) \quad (式 22)$$

$$NB(y|size_A + size_B, prob) = \sum_{i=0}^y NB(i|size_A, prob) * NB(y-i|size_B, prob) \quad (式 23)$$

が成り立つ。これは負の二項分布のモーメント母関数の定義 ([データ科学便覧](#))

$$M_x(t) = \left(\frac{p}{1 - (1-p) * exp(t)} \right)^r \quad (式 24)$$

で $r = size, p = prob$ と置くと明らかと言える。

ただしこの定義の基になる負の二項分布の定義は、本文書で一貫して用いているものとは逆の「1 回の試行に確率 p で成功するとき、 r 回目の成功までの、失敗回数 y の分布」であるため (R の `dnbinom` 関数と同じ)、本文書の定義に合わせると以下ようになる。

$$M_x(t) = \left(\frac{1 - prob}{1 - prob * exp(t)} \right)^{size} \quad (式 25)$$

ここではガンマ関数の引数が整数に限るつまり階乗に置き換えられるものとして、sum 記号の中を幾何的に解く。

言葉で書くと、 r 回目に失敗するまでに、 y 回成功するのは、 i 回成功してから $y-i$ 成功することと同じ、という意味である。式 6 を用いて以下のように展開できる。

$$NB(y|size_A+size_B, prob) = \sum_{i=0}^y \frac{\Gamma(i+size_A)}{\Gamma(size_A)i!} (1-prob)^{size_A} prob^i \frac{\Gamma(y-i+size_B)}{\Gamma(size_B)(y-i)!} (1-prob)^{size_B} prob^{y-i} \quad (式 26)$$

$$= (1-prob)^{size_A+size_B} prob^y \sum_{i=0}^y \frac{\Gamma(i+size_A)}{\Gamma(size_A)i!} \frac{\Gamma(y-i+size_B)}{\Gamma(size_B)(y-i)!} \quad (式 27)$$

碁盤の目状の街があり、縦 m 区間、横 n 区間あるとする。このとき街の左上の端から右下の端まで、右または下方向に移動することを考える。つまり下、下、右、下、... と移動する。このとき左上の端から右下の端までの経路は、 $(m+n)!/(m!n!)$ 通りある。

なぜなら下同士、右同士の移動を移動時の座標に基づいて、下 1、下 2、右 1、下 3、と区別したときの経路数は $(m+n)!$ であるが、実は下同士、右同士の移動を区別しないので、 $m!n!$ で割ると、上記の通りになるからである。これは $m+n$ 個から m 個取り出す組み合わせの数である。

ここで縦 y 区間、横 $size_A + size_B - 1$ 区間ある街を考える。ここで、

1. 縦 i 区間、横 $size_A - 1$ 区間移動して
2. 右に 1 区間移動して
3. 縦 $y - i$ 区間、横 $size_B - 1$ 区間移動する

経路の数を考え、取りうるすべての $i \in 0..y$ について足す。これは i 回成功してから $y - i$ 成功する経路を漏れなく重複なく数えることと同じである。よって、

$$\binom{y+size_A+size_B-1}{y} = \sum_{i=0}^y \binom{y+size_A-1}{y} \binom{y+size_B-1}{y} \quad (式 28)$$

である。

ガンマ分布の再生性

ガンマ分布 $GammaDistribution(\alpha, \beta)$ は $\alpha = shape$ に対して再生性がある。つまり

$$x_A \sim GammaDistribution(\alpha_A, \beta), x_B \sim GammaDistribution(\alpha_B, \beta) \quad (式 29)$$

なら

$$x_A + x_B \sim GammaDistribution(\alpha_A + \alpha_B, \beta) \quad (式 30)$$

である。これはガンマ分布のモーメント母関数の定義

$$M_x(t) = (1 - t/\beta)^\alpha \quad (\text{式 31})$$

から明らかと言えるが、証明は[こちら](#)。

NB(size, prob) の再生性を確認する
再生性の定義より、

$$\sum_{i=1}^n NB(\text{size}, \text{prob}) = NB(\text{size} * n, \text{prob}) \text{ or } \sum_{i=1}^n NB(\text{size}/n, \text{prob}) = NB(\text{size}, \text{prob}) \quad (\text{式 32})$$

である。これは n 人がそれぞれ、確率 prob で当たるくじを size 回当たるまでくじを引く ($NB(\text{size}, \text{prob})$) ときの合計試行回数と、一人で $\text{size} * n$ 回当たるまでくじを引く ($NB(\text{size} * n, \text{prob})$) ときの試行回数が、同じ分布に従うということである。 n 人が同一人物と考えれば、そうなるのが自然である。くじに外れる確率を $\text{prob} := 1 - \text{prob}$ としても意味は同じである。

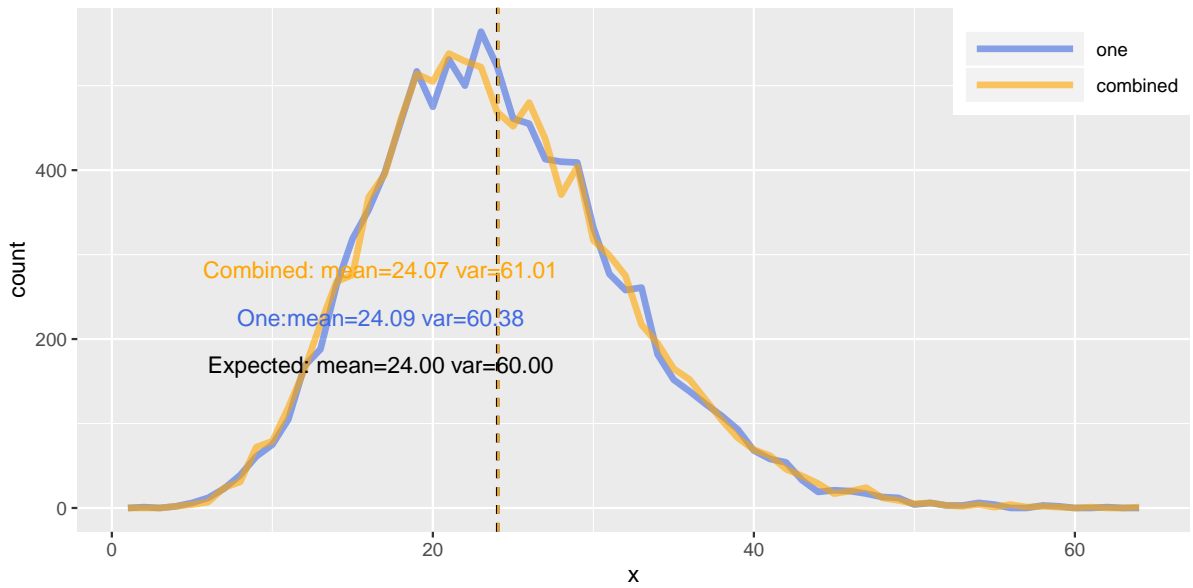


図1 negative binomial(size=16.00, prob=0.60) alpha combined

図1は $\text{size} = 16, \text{prob} = 0.6$ の負の二項分布に基づいて生成した乱数のヒストグラムである。R の `rnbinom` は本文書の prob を $1.0 - \text{prob}$ としているので、`rnbinom` の引数は 0.4 であることに注意。横軸は x 、縦軸は各 x の出現回数である。Expected は式 19 による平均と分散、one は単一の負の二項分布、combined は size を 16 分割した複数の負の二項分布の和である。

式 9 より、ガンマ分布の期待値は $\alpha := \text{shape}$ にも $\theta := \text{scale} = 1/\beta$ にも比例するが、式 12 よりガンマ-ポアソン分布の期待値と分散は α に比例しても、分散は $\theta = 1/\beta$ に比例しない。なので、 θ を n 分割した負の二項分布の和は、元の負の二項分布とは異なる (図2)。

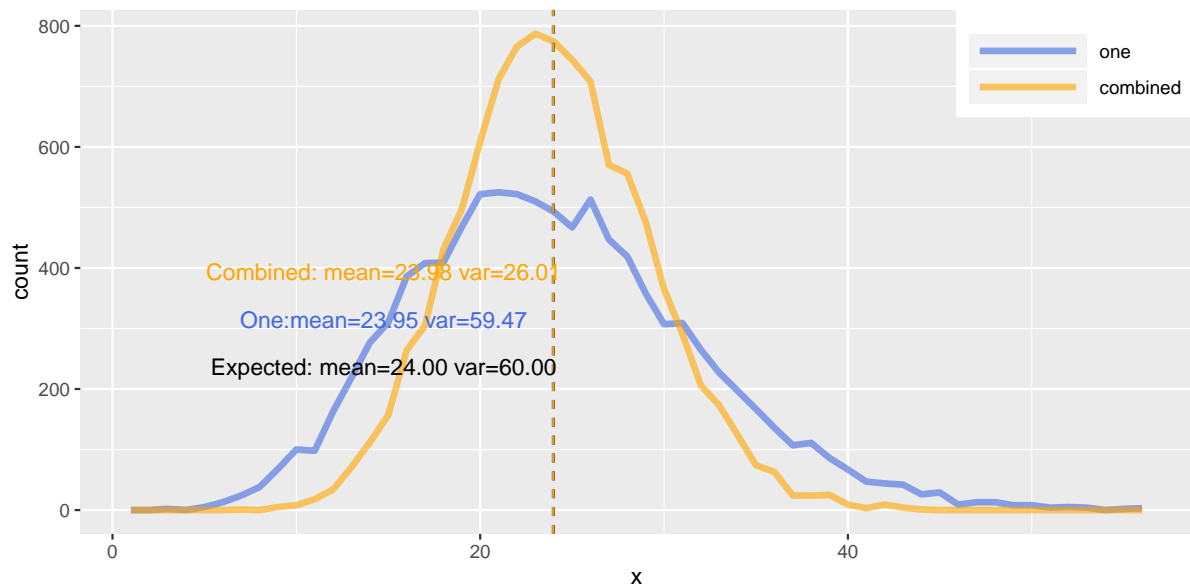


図2 negative binomial(size=16.00, prob=0.60) theta combined

ガンマ-ポアソン分布 (α, θ) の再生性を確認する

図1と同じ分布になるようにガンマ-ポアソン分布を作る。図3が示すように *shape* には再生性があり、図4が示すように *shape* には再生性がない。

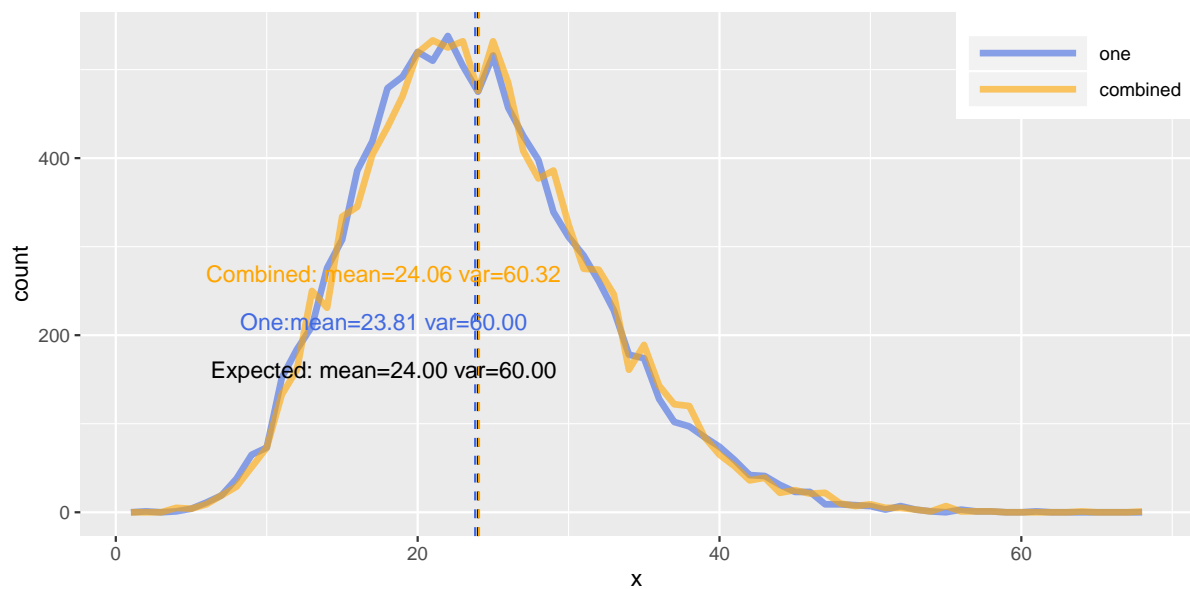


図3 Gamma poisson(shape=16.00, rate=0.67) alpha combined

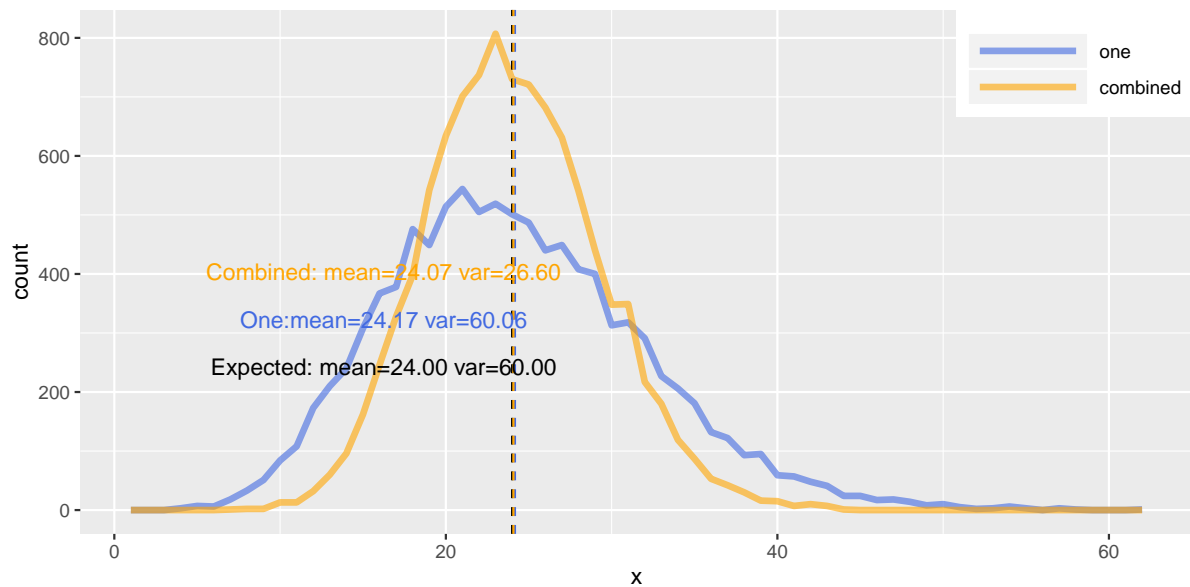


図4 Gamma poisson(shape=16.00, rate=0.67) theta combined

Gamma(α, θ) の再生性を確認する

負の二項分布と同様、 $\alpha = size$ の分割は再生性がある (図5)。

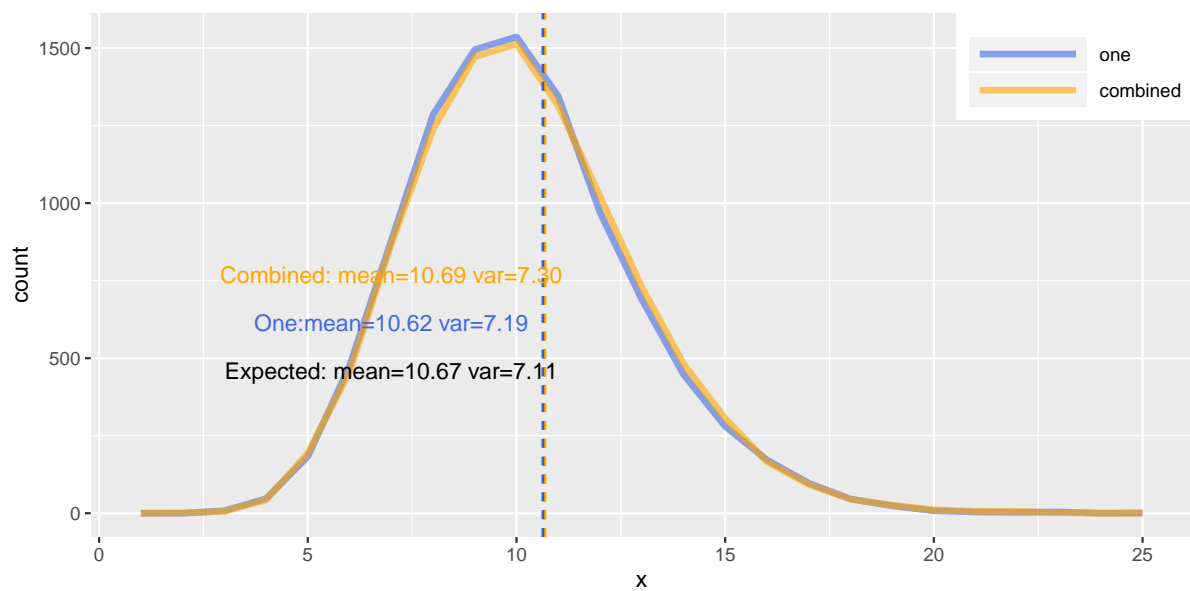


図5 Gamma distribution alpha combined

やはり負の二項分布と同様、 θ の分割は再生性がない (期待値は同じだが分散は異なる: 図6)。

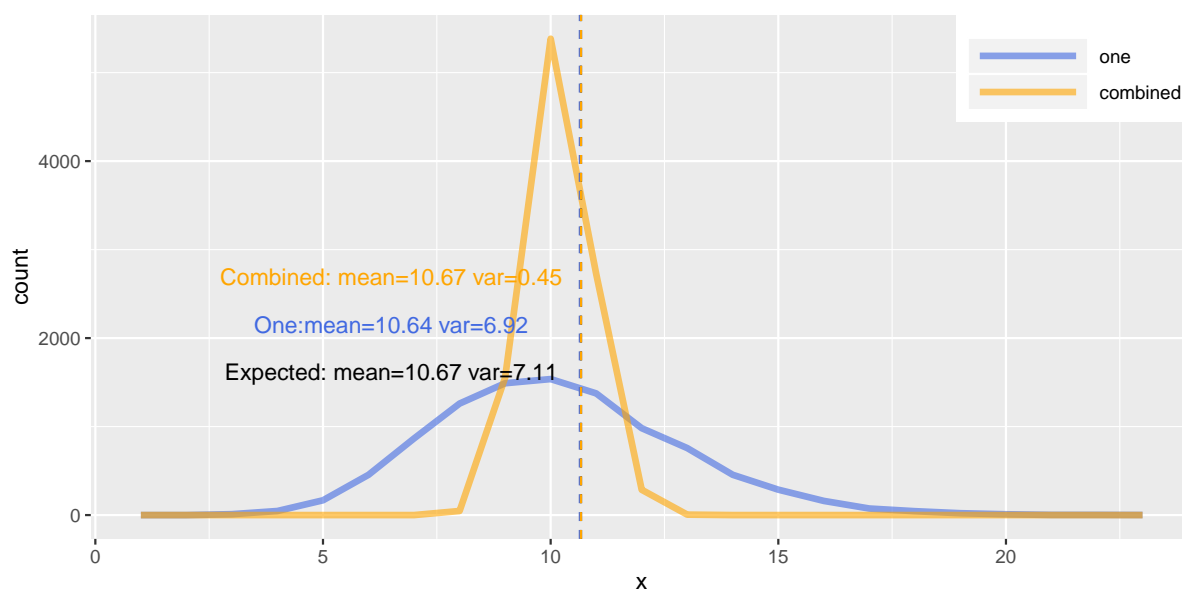


图6 Gamma distribution theta combined