

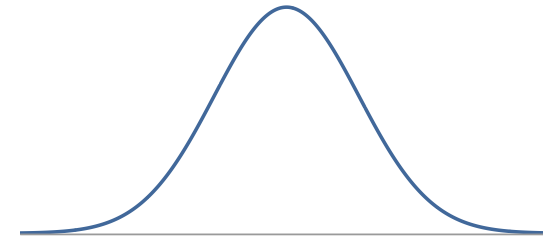
第4章：確率変数の分布

OpenIntro Statistics 第4版（日本語版）

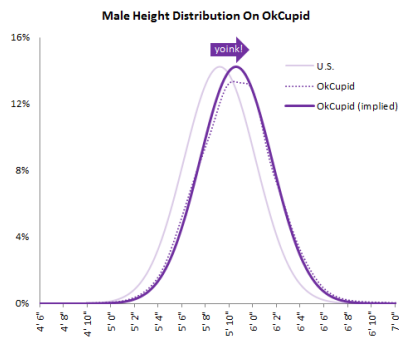
原著スライド：Mine Çetinkaya-Rundel（OpenIntro）
 CC BY-SA ライセンスのもと使用・翻訳。
 一部の画像はフェアユース（教育目的）に基づき使用。

正規分布

- 単峰で対称な、ベル型の曲線
- 多くの変数はほぼ正規分布に従うが、完全に正規分布に従うものはない
- $N(\mu, \sigma)$ と表す → 平均 μ 、標準偏差 σ の正規分布



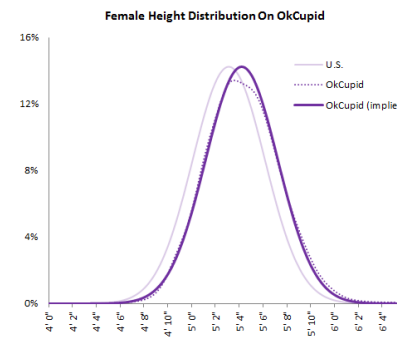
男性の身長



「OkCupid の男性の身長は、期待される正規分布にほぼ従っているが、全体的に右にずれている。男性はほぼ普遍的に数センチ多めに申告する傾向がある。」
 「約5フィート8インチ付近から、点線の曲線の上部がさらに右に傾いていることがわかる。これは、6フィートに近づくとつれて男性が少し多めに切り上げていることを意味し、その羨望の心理的な基準値を目指している。」

<http://blog.okcupid.com/index.php/the-biggest-lies-in-online-dating/>

女性の身長



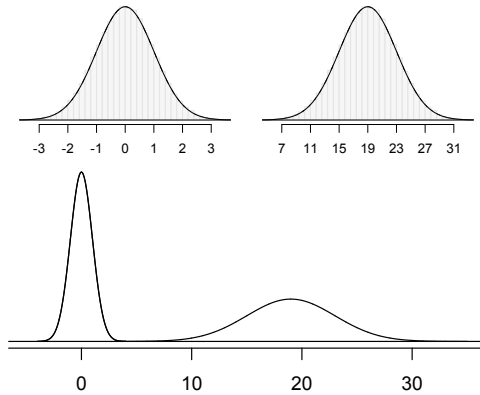
「女性のデータを調べたところ、身長の誇張は同様に広まっていたが、特定の基準値に向けた急激な変化はなかった。」

<http://blog.okcupid.com/index.php/the-biggest-lies-in-online-dating/>

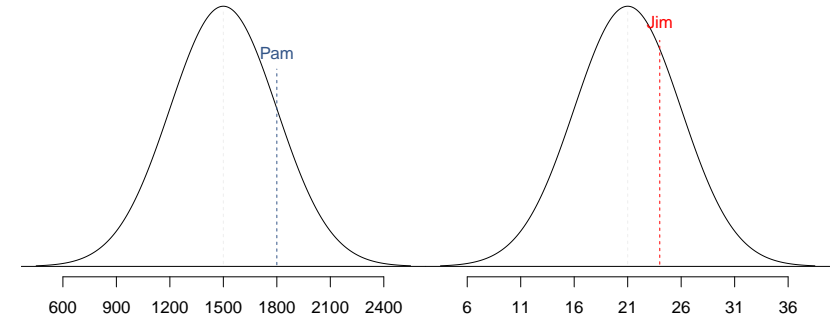
パラメータの異なる正規分布

μ : 平均、 σ : 標準偏差

$N(\mu = 0, \sigma = 1)$ $N(\mu = 19, \sigma = 4)$



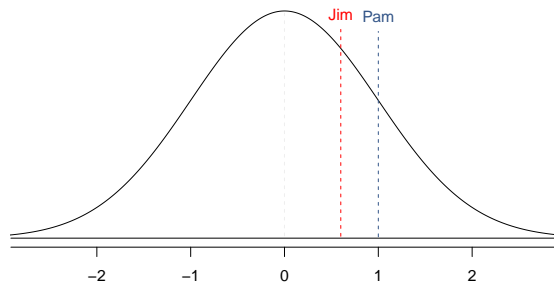
SAT のスコアは平均 1500、標準偏差 300 の正規分布にほぼ従う。ACT のスコアは平均 21、標準偏差 5 の正規分布にほぼ従う。大学の入学担当者が 2 人の志願者のうち、他の受験者に対してどちらが標準化テストでより良い成績を収めたかを判断したい。SAT で 1800 点を獲得したパムと、ACT で 24 点を取ったジムのどちらが優れているか？



Zスコアによる標準化

これら 2 つの生の点数を直接比較することはできないため、代わりに各観測値が平均から何標準偏差離れているかを比較する。

- パムのスコアは平均より $\frac{1800-1500}{300} = 1$ 標準偏差上にある。
- ジムのスコアは平均より $\frac{24-21}{5} = 0.6$ 標準偏差上にある。



Zスコアによる標準化 (続き)

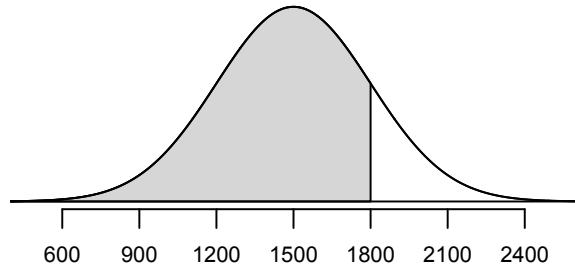
- これらは標準化スコア、または Zスコアと呼ばれる。
- 観測値の Zスコアは、それが平均から何標準偏差上または下にあるかを示す。

$$Z = \frac{\text{観測値} - \text{平均}}{SD}$$

- Zスコアはどんな形の分布に対しても定義できるが、分布が正規分布の場合にのみ Zスコアを使ってパーセンタイルを計算できる。
- 平均から 2SD 以上離れている観測値 ($|Z| > 2$) は通常、異常値と見なされる。

パーセンタイル

- **パーセンタイル**とは、ある特定のデータ点より下に位置する観測値の割合のことである。
- グラフ的には、パーセンタイルはその観測値の左側の確率分布曲線の下面積である。



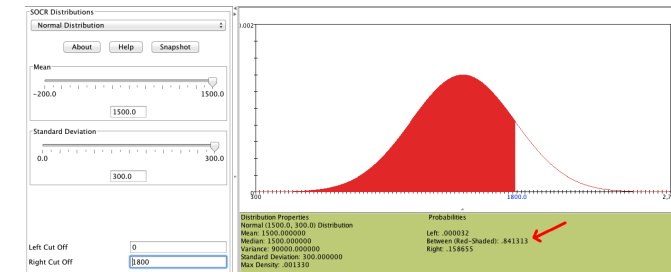
パーセンタイルの計算 — 計算機を使う

曲線の下面積（パーセンタイル）を計算するには様々な方法がある：

- R :

```
> pnorm(1800, mean = 1500, sd = 300)
[1] 0.8413447
```

- アプレット：https://gallery.shinyapps.io/dist_calc/



パーセンタイルの計算 — 表を使う

Z	Z の小数第 2 位									
	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.5000	0.5040	0.5080	0.5120	0.5160	0.5199	0.5239	0.5279	0.5319	0.5359
0.1	0.5398	0.5438	0.5478	0.5517	0.5557	0.5596	0.5636	0.5675	0.5714	0.5753
0.2	0.5793	0.5832	0.5871	0.5910	0.5948	0.5987	0.6026	0.6064	0.6103	0.6141
0.3	0.6179	0.6217	0.6255	0.6293	0.6331	0.6368	0.6406	0.6443	0.6480	0.6517
0.4	0.6554	0.6591	0.6628	0.6664	0.6700	0.6736	0.6772	0.6808	0.6844	0.6879
0.5	0.6915	0.6950	0.6985	0.7019	0.7054	0.7088	0.7123	0.7157	0.7190	0.7224
0.6	0.7257	0.7291	0.7324	0.7357	0.7389	0.7422	0.7454	0.7486	0.7517	0.7549
0.7	0.7580	0.7611	0.7642	0.7673	0.7704	0.7734	0.7764	0.7794	0.7823	0.7852
0.8	0.7881	0.7910	0.7939	0.7967	0.7995	0.8023	0.8051	0.8078	0.8106	0.8133
0.9	0.8159	0.8186	0.8212	0.8238	0.8264	0.8289	0.8315	0.8340	0.8365	0.8389
1.0	0.8413	0.8438	0.8461	0.8485	0.8508	0.8531	0.8554	0.8577	0.8599	0.8621
1.1	0.8643	0.8665	0.8686	0.8708	0.8729	0.8749	0.8770	0.8790	0.8810	0.8830
1.2	0.8849	0.8869	0.8888	0.8907	0.8925	0.8944	0.8962	0.8980	0.8997	0.9015

シックスシグマ

「シックスシグマプロセスという言葉は、グラフに示すように、プロセス平均と最も近い規格限界との間に6つの標準偏差がある場合、事実上すべての製品が規格を満たすという考え方に由来する。」

6σ

品質管理

ハイソックスのケチャップ工場では、ケチャップのボトルに入る量が平均 36 オンス、標準偏差 0.11 オンスの正規分布に従うとされている。30 分ごとに生産ラインからボトルが 1 本選ばれ、その内容量が正確に記録される。ケチャップの量が 35.8 オンス未満または 36.2 オンス超の場合、そのボトルは品質管理検査に不合格となる。ケチャップが 35.8 オンス未満のボトルは何パーセントか？

正確な確率の求め方 — R を使う

```
> pnorm(-1.82, mean = 0, sd = 1)
[1] 0.0344
```

または

```
> pnorm(35.8, mean = 36, sd = 0.11)
[1] 0.0345
```

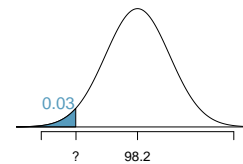
練習問題

品質管理検査に合格するボトルは何パーセントか？

- (a) 1.82%
- (b) 3.44%
- (c) 6.88%
- (d) 93.12%
- (e) 96.56%

カットオフ点の求め方

健康な人間の体温は平均 98.2°F、標準偏差 0.73°F の正規分布にほぼ従う。人間の体温の下位 3% のカットオフは何°F か？



$$\begin{aligned}
 P(X < x) &= 0.03 \\
 \rightarrow P(Z < -1.88) &= 0.03 \\
 Z &= \frac{\text{観測値} - \text{平均}}{SD} \\
 \rightarrow \frac{x - 98.2}{0.73} &= -1.88 \\
 x &= (-1.88 \times 0.73) + 98.2 = 96.8^\circ F
 \end{aligned}$$

```
> qnorm(0.03)
[1] -1.880794
```

Mackowiak, Wasserman, and Levine (1992). A Critical Appraisal of 98.6 Degrees F, the Upper Limit of the Normal Body Temperature, and Other Legacies of Carl Reinhold August Wunderlick.

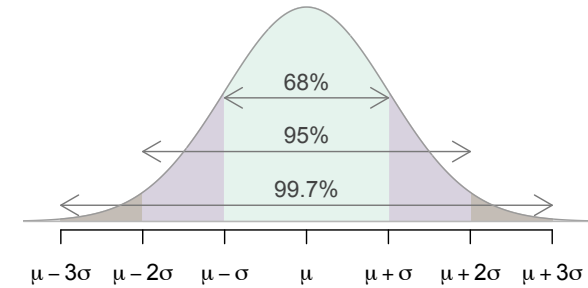
練習問題

健康な人間の体温は平均 98.2°F 、標準偏差 0.73°F の正規分布にほぼ従う。
人間の体温の上位 10% のカットオフは何 $^{\circ}\text{F}$ か？

- (a) 97.3°F
- (b) 99.1°F
- (c) 99.4°F
- (d) 99.6°F

68-95-99.7 則

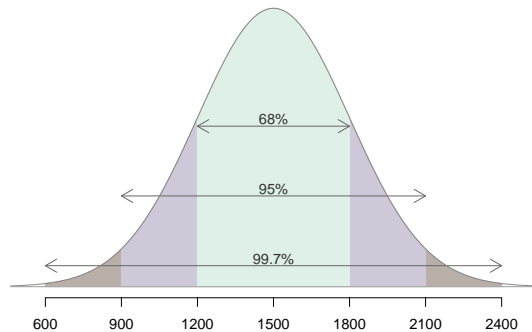
- ほぼ正規分布に従うデータでは、
 - 約 68% が平均から 1SD 以内に収まる、
 - 約 95% が平均から 2SD 以内に収まる、
 - 約 99.7% が平均から 3SD 以内に収まる。
- 平均から 4SD、5SD、またはそれ以上離れた観測値が存在することもあるが、データがほぼ正規分布に従う場合、このような出来事は非常にまれである。



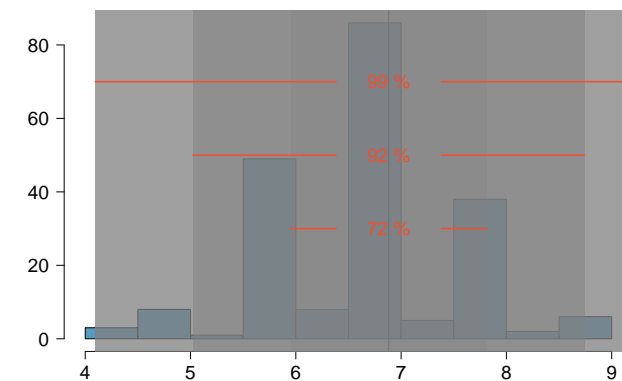
68-95-99.7 則を使った変動性の記述

SAT のスコアは平均 1500、標準偏差 300 の正規分布にほぼ従う。

- SAT で 1200~1800 点の間に収まる学生は約 68%。
- SAT で 900~2100 点の間に収まる学生は約 95%。
- SAT で 600~2400 点の間に収まる学生は約 99.7%。



平日の睡眠時間



- 平均 = 6.88 時間、SD = 0.92 時間
- データの 72% は平均から 1SD 以内： 6.88 ± 0.93
- データの 92% は平均から 2SD 以内： $6.88 \pm 2 \times 0.93$
- データの 99% は平均から 3SD 以内： $6.88 \pm 3 \times 0.93$

幾何分布

スミス博士はミルグラムの実験を繰り返したいが、強烈なショックを与えない人が見つかるまでサンプリングしたいと考えている。最初の人でやめる確率は何か？

$$P(1人目が拒否) = 0.35$$

……3人目の人でやめる確率は？

$$\begin{aligned} P(1人目と2人目がショック、3人目が拒否) &= \frac{S}{0.65} \times \frac{S}{0.65} \times \frac{R}{0.35} \\ &= 0.65^2 \times 0.35 \approx 0.15 \end{aligned}$$

……10人目の人でやめる確率は？

幾何分布（続き）

幾何分布は、独立同一分布 (iid) のベルヌーイ確率変数において、最初の成功までの待ち時間を記述する。

- 独立性：各試行の結果は互いに影響しない
- 同一性：成功確率は各試行で同じ

幾何確率

成功確率を p 、失敗確率を $(1 - p)$ 、独立試行回数を n とすると

$$P(n回目の試行で初めて成功) = (1 - p)^{n-1} p$$

サイコロを6回振って初めて6の目が出る確率を幾何分布で計算できるか？ 成功（6の目が出る）と失敗（6の目が出ない）は明確に定義されており、各試行でどちらかが必ず起こることに注意せよ。

- (a) いいえ、サイコロを振ると2つ以上の結果がある
- (b) はい、計算できる

期待値

スミス博士は、ショックを与えることを拒否する最初の人を見つけるまで、何人の人を試験すると予想されるか？

幾何分布の期待値、すなわち平均は $\frac{1}{p}$ と定義される。

$$\mu = \frac{1}{p} = \frac{1}{0.35} = 2.86$$

スミス博士は、ショックを与えることを拒否する最初の人を見つけるまで、2.86人の人を試験すると予想される。

しかし、整数でない人数をどのように試験するのか？

シナリオ数の計算

二項係数

二項係数 (choose 関数) は、 n 回の試行で k 回成功する方法の数を計算するのに役立つ。

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

- $k = 1, n = 4$: $\binom{4}{1} = \frac{4!}{1!(4-1)!} = \frac{4 \times 3 \times 2 \times 1}{1 \times (3 \times 2 \times 1)} = 4$
- $k = 2, n = 9$: $\binom{9}{2} = \frac{9!}{2!(9-2)!} = \frac{9 \times 8 \times 7!}{2 \times 1 \times 7!} = \frac{72}{2} = 36$

注: R を使っても計算できる:

```
> choose(9,2)
[1] 36
```

二項係数の性質

次のうち誤りはどれか?

- n 回の試行で 1 回成功する方法は n 通りある、 $\binom{n}{1} = n$ 。
- n 回の試行で n 回成功する方法は 1 通りしかない、 $\binom{n}{n} = 1$ 。
- n 回の試行で n 回失敗する方法は 1 通りしかない、 $\binom{n}{0} = 1$ 。
- n 回の試行で $n - 1$ 回成功する方法は $n - 1$ 通りある、 $\binom{n}{n-1} = n - 1$ 。

二項分布 (続き)

二項確率

p が成功確率、 $(1 - p)$ が失敗確率、 n が独立試行回数、 k が成功回数を表すとき

$$P(n \text{ 回中 } k \text{ 回の成功}) = \binom{n}{k} p^k (1 - p)^{(n-k)}$$

二項分布を適用するために満たす必要のない条件はどれか?

- 試行は独立でなければならない
- 試行回数 n は固定されていなければならない
- 各試行の結果は成功または失敗に分類されなければならない
- 望ましい成功回数 k は試行回数より大きくななければならない
- 成功確率 p は各試行で同じでなければならない

2012年のギャラップ調査によると、アメリカ人の26.2%が肥満である。10人のアメリカ人の無作為標本の中で、ちょうど8人が肥満である確率は何か？

- (a) かなり高い
- (b) かなり低い

Gallup: <http://www.gallup.com/poll/160061/obesity-rate-stable-2012.aspx>, January 23, 2013.

2012年のギャラップ調査によると、アメリカ人の26.2%が肥満である。10人のアメリカ人の無作為標本の中で、ちょうど8人が肥満である確率は何か？

- (a) $0.262^8 \times 0.738^2$
- (b) $\binom{8}{10} \times 0.262^8 \times 0.738^2$
- (c) $\binom{10}{8} \times 0.262^8 \times 0.738^2 = 45 \times 0.262^8 \times 0.738^2 = 0.0005$
- (d) $\binom{10}{8} \times 0.262^2 \times 0.738^8$

誕生日問題

無作為に選んだ2人が同じ誕生日を持つ確率は何か？

かなり低く、 $\frac{1}{365} \approx 0.0027$ である。

366人の中で、少なくとも2人が同じ誕生日を持つ確率は何か？

確実に1である！（閏年の誕生日の可能性を除く。）

誕生日問題（続き）

121人の中で、少なくとも2人（1組のペア）が同じ誕生日を持つ確率は何か？

計算はやや複雑だが、121人の中で一致がない確率の余事象として考えることができる。

$$\begin{aligned}
 P(\text{一致なし}) &= 1 \times \left(1 - \frac{1}{365}\right) \times \left(1 - \frac{2}{365}\right) \\
 &\quad \times \cdots \times \left(1 - \frac{120}{365}\right) \\
 &= \frac{365 \times 364 \times \cdots \times 245}{365^{121}} \\
 &= \frac{365!}{365^{121} \times (365 - 121)!} \\
 &= \frac{121! \times \binom{365}{121}}{365^{121}} \approx 0
 \end{aligned}$$

$$P(\text{少なくとも1組のペア}) \approx 1$$

期待値

2012年のギャラップ調査によると、アメリカ人の26.2%が肥満である。

100人のアメリカ人の無作為標本の中で、何人が肥満であると予想されるか？

- 簡単に計算すると、 $100 \times 0.262 = 26.2$ である。
- より形式的には、 $\mu = np = 100 \times 0.262 = 26.2$ となる。
- しかし、これはすべての100人の無作為標本でちょうど26.2人が肥満であることを意味するわけではない。実際、それは不可能である。標本によってこの値は小さかったり大きかったりする。この値がどれだけ変動すると予想されるか？

異常な観測値

平均から2標準偏差以上離れている観測値は異常と見なすという考えと、先ほど計算した平均と標準偏差を使って、100人の無作為標本における肥満者数の妥当な範囲を計算できる。

$$26.2 \pm (2 \times 4.4) = (17.4, 35)$$

期待値とその変動性

二項分布の平均と標準偏差

$$\mu = np \quad \sigma = \sqrt{np(1-p)}$$

- 肥満率に戻ると：

$$\sigma = \sqrt{np(1-p)} = \sqrt{100 \times 0.262 \times 0.738} \approx 4.4$$

- 無作為に選んだ100人のアメリカ人のうち26.2人が肥満であると予想され、標準偏差は4.4である。

注：二項分布の平均と標準偏差が常に整数であるとは限らないが、それで構わない。これらの値は平均的に期待されることを表している。

2012年8月のギャラップ調査によると、アメリカ人の13%がホームスクーリングは子どもに優れた教育を提供すると考えている。1,000人のアメリカ人の無作為標本でこの意見を持つ人が100人のみだった場合、それは異常と見なされるか？

- (a) いいえ (b) はい

<http://www.gallup.com/poll/156974/private-schools-top-marks-educating-children.aspx>

二項分布の形状

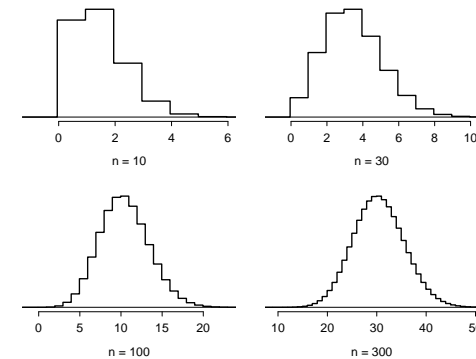
この活動にはウェブアプリを使用する。

https://gallery.shinyapps.io/dist_calc/ にアクセスし、左のドロップダウンメニューで「Binomial coin experiment」を選ぶ。

- 試行回数を 20、成功確率を 0.15 に設定する。成功回数の分布の形状を説明せよ。
- p を 0.15 に固定したまま、成功回数の分布が単峰かつ対称になるのに必要な最小の標本サイズを求めよ。
- さらに検討：
 - n を一定に保ちながら p が変化すると分布の形状はどうなるか？
 - p を一定に保ちながら n が変化すると分布の形状はどうなるか？

成功回数の分布

$p = 0.10$ 、 $n = 10, 30, 100, 300$ の二項モデルから得られたサンプルのヒストグラム。 n が増加するとどうなるか？



十分大きいとはどれくらいか？

期待される成功回数と失敗回数がともに少なくとも 10 であれば、標本サイズは十分大きいと見なされる。

$$np \geq 10 \quad \text{かつ} \quad n(1-p) \geq 10$$

以下の 4 組の二項分布パラメータのうち、正規分布で近似できる分布はどれか？

- (a) $n = 100, p = 0.95$
- (b) $n = 25, p = 0.45 \rightarrow 25 \times 0.45 = 11.25; 25 \times 0.55 = 13.75$
- (c) $n = 150, p = 0.05$
- (d) $n = 500, p = 0.015$

Facebook ユーザーの分析

最近の研究によると「Facebook ユーザーは与えるよりも多くを受け取っている」。例えば：

- 標本内の Facebook ユーザーの 40% が友達リクエストを送ったが、63% が少なくとも 1 件のリクエストを受け取った
- ユーザーは友達のコンテンツに平均 14 回「いいね」を押したが、自分のコンテンツには平均 20 回「いいね」がついた
- ユーザーは 9 件の個人メッセージを送ったが、12 件を受け取った
- ユーザーの 12% が写真で友達をタグ付けしたが、35% が自分自身を写真でタグ付けされた

このパターンをどのように説明できるか？

<http://www.pewinternet.org/Reports/2012/Facebook-users/Summary.aspx>

この研究では、Facebook ユーザーの約 25% がパワーユーザーと見なされることもわかった。同じ研究で、平均的な Facebook ユーザーは 245 人の友達を持つことが明らかになった。245 人の友達を持つ平均的な Facebook ユーザーが、パワーユーザーと見なされる友達が 70 人以上いる確率は何か？ 必要な仮定を述べよ。

$n = 245$ 、 $p = 0.25$ が与えられ、確率 $P(K \geq 70)$ を求める。独立性が必要であり、ここでは仮定する (Facebook のデータにアクセスできれば確認できる)。

$$P(X \geq 70) = P(K = 70 \text{ または } K = 71 \text{ または } \dots \text{ または } K = 245) \\ = P(K = 70) + P(K = 71) + \dots + P(K = 245)$$

これは膨大な計算になりそうだ……

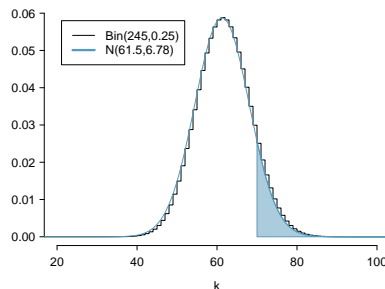
二項分布の正規近似

標本サイズが十分大きい場合、パラメータ n と p の二項分布は、パラメータ $\mu = np$ 、 $\sigma = \sqrt{np(1-p)}$ の正規モデルで近似できる。

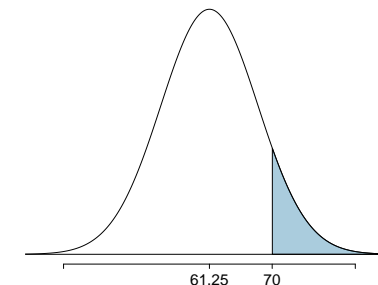
- Facebook のパワーユーザーの場合、 $n = 245$ 、 $p = 0.25$ である。

$$\mu = 245 \times 0.25 = 61.25 \quad \sigma = \sqrt{245 \times 0.25 \times 0.75} = 6.78$$

- $Bin(n = 245, p = 0.25) \approx N(\mu = 61.25, \sigma = 6.78)$ 。



245 人の友達を持つ平均的な Facebook ユーザーが、パワーユーザーと見なされる友達が 70 人以上いる確率は何か？



$$Z = \frac{\text{観測値} - \text{平均}}{SD} \\ = \frac{70 - 61.25}{6.78} = 1.29$$

$$P(Z > 1.29) = 1 - 0.9015 = 0.0985$$

```
> pnorm(1.29)
[1] 0.9014747
```

小さい区間での正規近似の限界

- 二項分布に対する正規近似は、条件が満たされていても、少数の値の範囲の確率を推定する際に精度が低下する傾向がある。
- 値の区間に対するこの近似は、カットオフ値を両方向に0.5ずつ拡張することで通常改善される。
- 正規近似を適用する際にこの修正を加えるヒントは、観測値の範囲を調べる際に最も役立つ。裾確率を計算する際にもこの修正を適用することは可能だが、区間が通常かなり広いため修正の利点は消えることが多い。

負の二項分布

- **負の二項分布**は、 n 回目の試行で k 回目の成功が起こる確率を記述する。
 - 負の二項分布の状況を識別するのに役立つ4つの条件：
 1. 試行は独立である。
 2. 各試行の結果は成功または失敗に分類される。
 3. 成功確率 (p) は各試行で同じである。
 4. 最後の試行は成功でなければならない。
- 最初の3つの条件は二項分布と共通している。

負の二項分布

$P(n \text{ 回目の試行で } k \text{ 回目の成功}) = \binom{n-1}{k-1} p^k (1-p)^{n-k}$ 、
ここで p は個々の試行が成功する確率である。すべての試行は独立と仮定する。

心理学研究室でアルバイトをしている大学生が、研究に参加するカップルを10組募集するよう頼まれた。彼女は学生センターの前に立ち、建物を出る5人おきに1人に対して、交際中かどうか、もし交際中であれば重要な他者と一緒に研究に参加したいかを尋ねることにした。このような人を見つける確率が10%だとする。目標を達成するまでに30人に尋ねる必要がある確率は何か？

与えられた情報： $p = 0.10$ 、 $k = 10$ 、 $n = 30$ 。30回目の試行で10回目の成功の確率を求めるため、負の二項分布を使用する。

$$\begin{aligned}
 P(30 \text{ 回目で } 10 \text{ 回目の成功}) &= \binom{29}{9} \times 0.10^{10} \times 0.90^{20} \\
 &= 10,015,005 \times 0.10^{10} \times 0.90^{20} \\
 &= 0.00012
 \end{aligned}$$

二項分布 vs. 負の二項分布

負の二項分布は二項分布とどのように異なるか？

- 二項分布の場合、通常は固定された試行回数があり、代わりに成功の回数を考える。
- 負の二項分布の場合、固定された回数の成功を観測するまでに何回の試行が必要かを調べ、最後の観測が成功であることが必要である。

