

大数の法則

大数の法則とは、観測数が増えるにつれて、ある特定の結果が起こる割合 \hat{p}_n が、その結果の確率 p に収束することを述べたものである。

大数の法則 (続き)

公正なコインを投げたとき、最初の 10 回すべてで表が出た場合、次の投げで表が出る確率はいくつだと思うか？ 0.5、0.5 より小さい、それとも 0.5 より大きい？

$H H H H H H H H H H ?$

- 確率はやはり 0.5 であり、次の投げで表が出る確率は 50% である。

$$P(H \text{ on } 11^{\text{th}} \text{ toss}) = P(T \text{ on } 11^{\text{th}} \text{ toss}) = 0.5$$

- コインは「裏を出すべき」ではない。
- 大数の法則についての一般的な誤解は、確率的过程は過去に起こったことを補正するはずだというものであるが、これは誤りであり、**ギャンブラーの誤謬** (または**平均の法則**) とも呼ばれる。

互いに排反な結果とそうでない結果

互いに排反 (*mutually exclusive*) な結果：同時には起こり得ない。

- 1 回のコイン投げの結果が表と裏の両方になることはない。
- 学生が同じ科目で合格と不合格の両方になることはない。
- 1 枚のカードがエースとクイーンの両方であることはない。

互いに排反でない結果：同時に起こり得る。

- 学生が同じ学期に統計学で A、経済学で A をとることがある。

互いに排反でない事象の和事象

よくシャッフルされた 1 組のトランプから、ジャックまたは赤いカードを引く確率は何か？

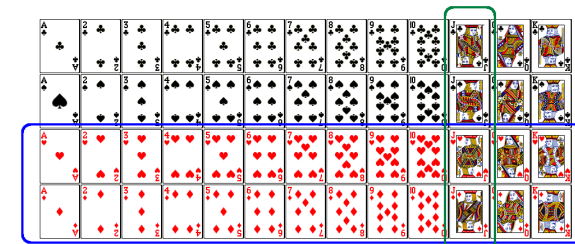


Figure from <http://www.milefoot.com/math/discrete/counting/cardfreq.htm>.

練習問題

無作為に選ばれた学生が、大麻を合法化すべきと考えているか、または親の政治的見解に同意している確率は何か？

大麻合法化	親の政治観に同意		合計
	なし	あり	
反対	11	40	51
賛成	36	78	114
合計	47	118	165

- (a) $\frac{40+36-78}{165}$
- (b) $\frac{114+118-78}{165}$
- (c) $\frac{78}{165}$
- (d) $\frac{78}{188}$
- (e) $\frac{11}{47}$

まとめ

一般加法定則

$$P(A \text{ または } B) = P(A) + P(B) - P(A \text{ かつ } B)$$

注：互いに排反な事象では $P(A \text{ かつ } B) = 0$ なので、上式は $P(A \text{ または } B) = P(A) + P(B)$ に簡略化される。

確率分布

確率分布は、起こりうるすべての事象とそれらが起こる確率を一覧にしたものである。

- 1人の子どもの性別の確率分布：

事象	男	女
確率	0.5	0.5

- 確率分布の規則：
 1. 列挙された事象は互いに排反でなければならない
 2. 各確率は0以上1以下でなければならない
 3. 確率の合計は1でなければならない
- 2人の子どもの性別の確率分布：

練習問題

調査で、回答者の52%が民主党員だと答えた。このサンプルから無作為に選ばれた回答者が共和党員である確率はいくつか？

- (a) 0.48
- (b) 0.48より大きい
- (c) 0.48より小さい
- (d) 与えられた情報だけでは計算できない

唯一の2つの政党が共和党と民主党である場合は (a) が可能である。しかし、政党に所属しない人やこれら2つ以外の政党に所属する人もいる可能性があるため、(c) も可能である。(b) は合計確率が1を超えることになるため、絶対に不可能である。

標本空間と余事象

標本空間は、試行で起こりうるすべての結果の集合である。

- ある夫婦に1人の子どもがいる場合、その性別の標本空間は？ $S = \{M, F\}$
- ある夫婦に2人の子どもがいる場合、その性別の標本空間は？

余事象は、確率の和が1になる2つの互いに排反な事象である。

- ある夫婦に1人の子どもがいる。その子が男の子でないとわかった場合、その性別は？ $\{M, F\}$ → 男の子と女の子は余事象である。
- ある夫婦に2人の子どもがいる。両方が女の子でないとわかった場合、可能な性別の組み合わせは？

独立

2つのプロセスが独立であるとは、一方の結果を知っても他方の結果についての有用な情報が得られない場合をいう。

- 最初の投げでコインが表になったことを知っても、2回目の投げでコインがどちらになるかについての有用な情報は得られない。→ 2回のコイン投げの結果は独立である。
- 最初に引いたカードがエースであることを知ると、2回目の引きでエースを引く確率についての有用な情報が得られる。→ デッキからの2回の引き（非復元）の結果は従属である。

練習問題

2013年1月9日～12日、SurveyUSAがノースカロライナ州の居住者500人の無作為標本を対象に、銃の広範な所有は法を守る市民を犯罪から守るか、社会をより危険にするかを質問した。回答者全体の58%が「市民を守る」と回答した。白人回答者の67%、黒人回答者の28%、ヒスパニック回答者の64%がこの見解を共有していた。次のうち正しいものはどれか？

銃所有に関する意見と人種・民族は最も可能性が高いのは

- 余事象
- 互いに排反
- 独立
- 従属
- 互いに排反

独立性の確認

$P(B \text{ が真であるとき } A \text{ が起こる}) = P(A | B) = P(A)$ ならば、A と B は独立である。

標本データに基づく従属性の判定

- 標本データから計算された条件付き確率が2つの変数間の従属性を示している場合、次のステップは、確率の観測された差が偶然に起こる可能性があるかを判定する仮説検定を行うことである。
- 条件付き確率間の観測された差が大きいほど、その差が真であるという証拠が強くなる。
- 標本が大きい場合、小さな差でも真の差の強い証拠となりうる。

$P(\text{市民を守る} \mid \text{白人}) = 0.67$, $P(\text{市民を守る} \mid \text{ヒスパニック}) = 0.64$ であった。
白人とヒスパニックで銃の広範な所有が市民を守ると思う割合に真の差があると確信するのはどちらの条件下か？ $n = 500$ か $n = 50,000$ か？

独立事象の積の法則

$$P(A \text{ かつ } B) = P(A) \times P(B)$$

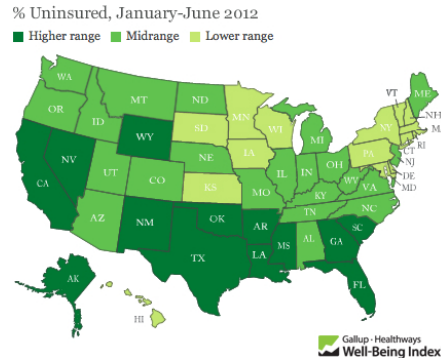
より一般的には、 $P(A_1 \text{ かつ } \dots \text{ かつ } A_k) = P(A_1) \times \dots \times P(A_k)$

コインを2回投げるとき、2回続けて裏が出る確率は何か？

$$P(\text{1回目に裏}) \times P(\text{2回目に裏}) = \frac{1}{2} \times \frac{1}{2} = \frac{1}{4}$$

練習問題

最近のギャラップ調査によると、2012年6月時点でテキサス州民の25.5%が健康保険を持っていない。非加入率が一定であると仮定して、無作為に選んだ2人のテキサス州民が両方とも非加入である確率は何か？



- (a) 25.5^2
- (b) 0.255^2
- (c) 0.255×2
- (d) $(1 - 0.255)^2$

<http://www.gallup.com/poll/156851/uninsured-rate-stable-across-states-far-2012.aspx>

互いに排反 vs. 余事象

2つの互いに排反な事象の確率の和は常に1になるか？

2つの余事象の確率の和は常に1になるか？

すべてをまとめると…

テキサス州民を5人無作為に選ぶとき、少なくとも1人が保険未加入である確率は何か？

- テキサス州民を5人無作為に選んだ場合、保険未加入者の人数の標本空間は：

$$S = \{0, 1, 2, 3, 4, 5\}$$

- 少なくとも1人が保険未加入である場合に注目する：

$$S = \{0, 1, 2, 3, 4, 5\}$$

- 標本空間を2つのカテゴリに分割できる：

$$S = \{0, \text{少なくとも1人}\}$$

すべてをまとめると…

標本空間の確率の和は1でなければならないため：

$$\begin{aligned} P(\text{少なくとも1人が未加入}) &= 1 - P(\text{未加入者なし}) \\ &= 1 - [(1 - 0.25)^5] \\ &= 1 - 0.745^5 \\ &= 1 - 0.23 \\ &= 0.77 \end{aligned}$$

少なくとも1

$$P(\text{少なくとも1}) = 1 - P(\text{1つもない})$$

練習問題

ある大学の学部生の約20%がベジタリアンまたはビーガンである。学部生3人の無作為標本から、少なくとも1人がベジタリアンまたはビーガンである確率は何か？

- (a) $1 - 0.2 \times 3$
- (b) $1 - 0.2^3$
- (c) 0.8^3
- (d) $1 - 0.8 \times 3$
- (e) $1 - 0.8^3$

再発

研究者たちは、コカイン慢性使用者72人を3つのグループに無作為に割り当てた：デシプラミン（抗うつ薬）、リチウム（コカインの標準治療）、プラセボ。研究結果を以下にまとめる。

	再発あり	再発なし	合計
デシプラミン	10	14	24
リチウム	18	6	24
プラセボ	20	4	24
合計	48	24	72

条件付き確率（続き）

患者が抗うつ薬（デシプラミン）を投与されたとわかっている場合、再発した確率は何か？

	再発あり	再発なし	合計
デシプラミン	10	14	24
リチウム	18	6	24
プラセボ	20	4	24
合計	48	24	72

$$P(\text{再発} \mid \text{デシプラミン}) = \frac{10}{24} \approx 0.42$$

$$P(\text{再発} \mid \text{リチウム}) = \frac{18}{24} \approx 0.75$$

$$P(\text{再発} \mid \text{プラセボ}) = \frac{20}{24} \approx 0.83$$

条件付き確率（続き）

患者が再発したとわかっている場合、抗うつ薬（デシプラミン）を投与されていた確率は何か？

	再発あり	再発なし	合計
デシプラミン	10	14	24
リチウム	18	6	24
プラセボ	20	4	24
合計	48	24	72

$$P(\text{デシプラミン} \mid \text{再発}) = \frac{10}{48} \approx 0.21$$

$$P(\text{リチウム} \mid \text{再発}) = \frac{18}{48} \approx 0.375$$

$$P(\text{プラセボ} \mid \text{再発}) = \frac{20}{48} \approx 0.42$$

一般乗法定則

- 2つの事象が独立である場合、その結合確率は単純にそれぞれの確率の積であることを見た。事象が独立でないと考えられる場合、結合確率は少し異なる方法で計算される。
- A と B が2つの結果または事象を表す場合、

$$P(A \text{ かつ } B) = P(A|B) \times P(B)$$

この式は、条件付き確率の公式を変形したものに過ぎない。

- A を目的の結果、B を条件として考えると便利である。

独立性と条件付き確率

以下は、入門統計学クラスの学生の性別と専攻の（仮定の）分布である：

	社会科学	非社会科学	合計
女性	30	20	50
男性	30	20	50
合計	60	40	100

- 無作為に選ばれた学生が社会科学専攻である確率は $\frac{60}{100} = 0.6$ 。
- 無作為に選ばれた学生が女性であるわかっている場合、社会科学専攻である確率は $\frac{30}{50} = 0.6$ 。
- $P(SS|M)$ も 0.6 に等しいため、このクラスの学生の専攻は性別に依存しない： $P(SS | F) = P(SS)$ 。

独立性と条件付き確率 (続き)

一般的に、 $P(A|B) = P(A)$ であれば、事象 A と B は独立であるといわれる。

- 概念的に： B を与えても A について何も教えてくれない。
- 数学的に： 事象 A と B が独立であれば $P(A \text{ かつ } B) = P(A) \times P(B)$ であることがわかっている。すると、

$$P(A|B) = \frac{P(A \text{ かつ } B)}{P(B)} = \frac{P(A) \times P(B)}{P(B)} = P(A)$$

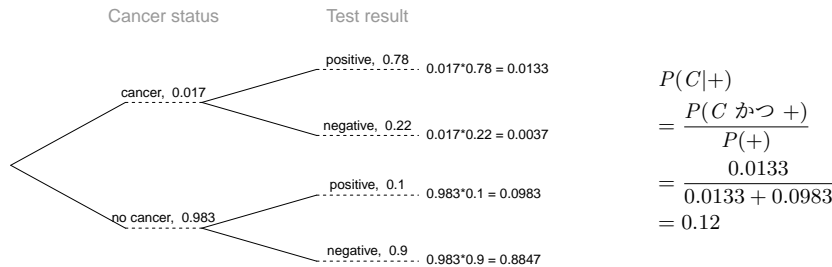
乳がん検診

- 米国がん協会の推定によれば、女性の約 1.7% が乳がんを患っている。
<http://www.cancer.org/cancer/cancerbasics/cancer-prevalence>
- Susan G. Komen For The Cure Foundation によれば、マンモグラフィは実際に乳がんを患っている女性の約 78% を正確に識別する。
<http://www5.komen.org/BreastCancer/AccuracyofMammograms.html>
- 2003 年に発表された論文によれば、すべてのマンモグラフィのうち最大 10% ががんを患っていない患者に対して偽陽性を示す。
<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC1360940>

注：これらの割合は近似値であり、推定が非常に難しい。

確率の逆転

患者が乳がん検診を受けるとき、2つの競合する主張がある：患者はがんである、患者はがんでない。マンモグラフィが陽性の結果を示した場合、患者が実際にかんである確率はいくつ？



$$\begin{aligned}
 P(C|+) &= \frac{P(C \text{ かつ } +)}{P(+)} \\
 &= \frac{0.0133}{0.0133 + 0.0983} \\
 &= 0.12
 \end{aligned}$$

注：樹形図は確率を逆転させるのに有用である： $P(+|C)$ が与えられ、 $P(C|+)$ を求める。

練習問題

1 回検査を受けて陽性の結果を得た女性が再検査を受けたいとする。2 回目の検査では、この特定の女性ががんである確率はいくつと仮定すべきか？

- (a) 0.017
- (b) 0.12
- (c) 0.0133
- (d) 0.88

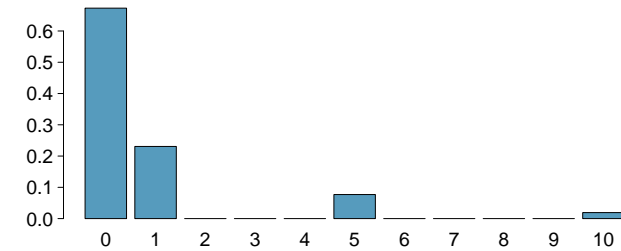
離散確率変数の期待値

カードゲームで、ハートを引いたら1ドル、エース（ハートのエースを含む）を引いたら5ドル、スペードのキングを引いたら10ドル、それ以外のカードを引いたら何も得られない。勝利金の確率モデルを書き、期待される勝利金を計算しなさい。

事象	X	$P(X)$	$X P(X)$
ハート（エース以外）	1	$\frac{12}{52}$	$\frac{12}{52}$
エース	5	$\frac{4}{52}$	$\frac{20}{52}$
スペードのキング	10	$\frac{1}{52}$	$\frac{10}{52}$
その他	0	$\frac{35}{52}$	0
合計			$E(X) = \frac{42}{52} \approx 0.81$

離散確率変数の期待値（続き）

以下はこのゲームの勝利金の確率分布を視覚的に表したものである：



変動性

確率変数の値の変動性にもしばしば関心がある。

$$\sigma^2 = \text{Var}(X) = \sum_{i=1}^k (x_i - E(X))^2 P(X = x_i)$$

$$\sigma = \text{SD}(X) = \sqrt{\text{Var}(X)}$$

離散確率変数の変動性

先のカードゲームの例で、ゲームごとに勝利金がどのくらい変動すると予想されるか？

X	$P(X)$	$X P(X)$	$(X - E(X))^2$	$P(X) (X - E(X))^2$
1	$\frac{12}{52}$	$1 \times \frac{12}{52} = \frac{12}{52}$	$(1 - 0.81)^2 = 0.0361$	$\frac{12}{52} \times 0.0361 = 0.0083$
5	$\frac{4}{52}$	$5 \times \frac{4}{52} = \frac{20}{52}$	$(5 - 0.81)^2 = 17.5561$	$\frac{4}{52} \times 17.5561 = 1.3505$
10	$\frac{1}{52}$	$10 \times \frac{1}{52} = \frac{10}{52}$	$(10 - 0.81)^2 = 84.4561$	$\frac{1}{52} \times 84.4561 = 1.6242$
0	$\frac{35}{52}$	$0 \times \frac{35}{52} = 0$	$(0 - 0.81)^2 = 0.6561$	$\frac{35}{52} \times 0.6561 = 0.4416$
		$E(X) = 0.81$		

線形結合

- 確率変数 X と Y の線形結合は次のように表される

$$aX + bY$$

ここで a と b はある固定された数である。

- 確率変数の線形結合の平均値は次のように計算される

$$E(aX + bY) = a \times E(X) + b \times E(Y)$$

線形結合の期待値の計算

統計学の宿題問題1つには平均10分、化学の宿題問題1つには平均15分かかる。今週は統計学の宿題5問と化学の宿題4問が出された。今週の統計学と化学の宿題に費やすと予想される合計時間は何か？

線形結合

- 2つの独立な確率変数の線形結合の変動性は次のように計算される

$$V(aX + bY) = a^2 \times V(X) + b^2 \times V(Y)$$

- 線形結合の標準偏差は分散の平方根である。

注： 確率変数が独立でない場合、分散の計算はやや複雑になり、本講義の範囲を超える。

線形結合の分散の計算

統計学の宿題問題1つにかかる時間の標準偏差は1.5分、化学の宿題問題1つには2分である。統計学の宿題5問と化学の宿題4問が出された今週に統計学と物理学の宿題に費やすと予想される時間の標準偏差は何か？ 各問題にかかる時間は互いに独立であると仮定する。

