

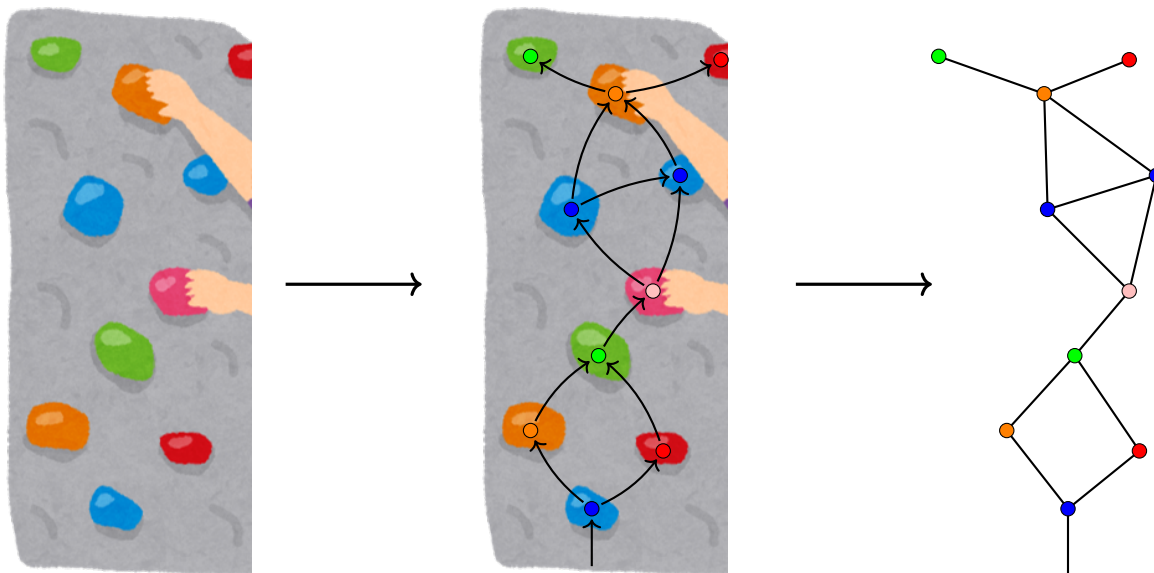
MAT120：数量的推論 第11講ハンドアウト

逐次確率

前回の講義では、不確実性を数量化する方法として確率を導入し、結果を数え上げることで古典的確率を計算する方法を学びました。今回の講義では、条件付き確率を発展させます。これは、追加の情報を得たあとの確率を意味します。また、多段階の実験を整理するための実用的な道具である、決定木についても学びます。最後に、コイン投げにおける結果の数え上げを始めます。これは次回の講義で扱う二項実験へとつながっていきます。

1 導入：岩を登る人

岩場を登っている人の手の動きを追いかけてみましょう。この状況は、色のついた岩どうしがどのようにつながっているかを表す一種の「グラフ」を作ることで理想化できます。



この図では、登る人が岩を上へ移動するにつれて、その手は下の青い岩から出発して、上の緑の岩または赤い岩までの経路をたどります。この過程のある段階では、登る人はどのように進むかを選ばなければなりません。たとえば岩場のいちばん最初では、青い岩からオレンジの岩へ進むこともできますし、あるいは赤い岩へ進むこともできます。他方で、ある地点では何をするかを選ぶことはできず、進み方は一意に決まっています。たとえば最初のオレンジの岩からは、次の緑の岩へ進むしかありません。

登る人が選択をするたびに、左の岩を選ぶ確率が50%、右の岩を選ぶ確率も50%であるとしましょう。実際、下から上までの異なる経路は全部で12本あります。では、登る人がある特定の経路をたどる確率はどれくらいでしょうか。

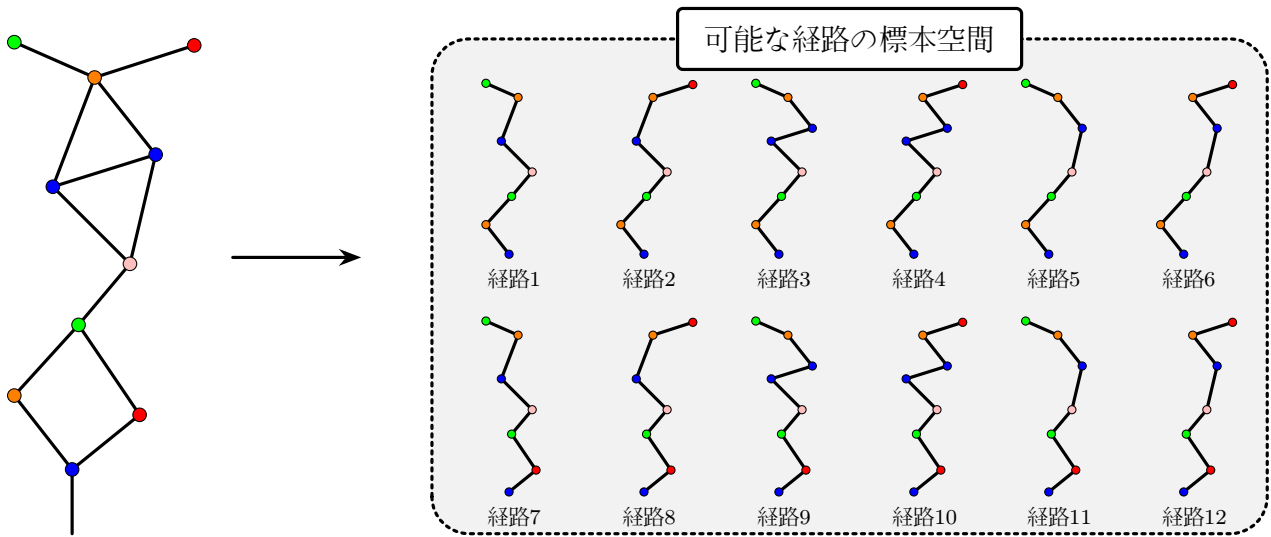
1.1 素朴な数え上げで解こうとしてみる

前回の講義では、古典的確率の設定で成り立つ基本公式を見ました。

$$P(\text{事象}) = \frac{\text{その事象に有利な結果の数}}{\text{結果の総数}}$$

要するに、一般の事象の確率は、その事象に対応する結果の数を数え、それを可能性の総数で割ることで計算できることが多いのです。

今回の岩登りの例でも、同じことを試してみたいくなります。まず標本空間を記述する必要がありますが、この場合それは下から上までのすべての経路の集合です。少し考えると、この岩場には異なる経路が12本あることが分かります。それらを下に図示します。



したがって、この設定における標本空間は

$$S = \{\text{経路1}, \text{経路2}, \text{経路3}, \dots, \text{経路12}\}$$

という集合です。したがって、上の公式を使ってある特定の経路の確率を計算したくなるのもっともです。たとえば経路1について、

$$P(\text{経路1}) \stackrel{?}{=} \frac{\text{その事象に有利な結果の数}}{\text{結果の総数}} = \frac{1}{12}$$

と推測したくなります。

しかし実際には、これは正しくありません。あとでこれを正しく計算しますが、今の段階では単に次のことを観察すれば十分です。経路全体の空間を見ると、ある経路は他の経路より複雑です。これらの異なる経路は同様には起こりません。したがって、上の公式はここでは正しいアプローチではないのです。

2 決定木

2.1 決定木とは何か

前回の講義では、条件付き確率という考え方を簡単に導入しました。これは、複合事象「A and B」の確率を計算する公式の中に現れていました。依存する事象A とB について、事象「A and B」の確率は

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

で与えられました。同値な形として、

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

とも書けます。ここで $P(B | A)$ という項が条件付き確率 であり、大まかには「A が起こると仮定したもとのB の確率」を意味します。一般に条件付き確率とは、仮定によって情報を更新した結果として生じる確率の変化のことです。

前回、条件付き確率は

$$P(B | A) = \frac{P(A \text{ and } B)}{P(A)} \quad (\text{ただし } P(A) > 0)$$

によって定義できることを見ました。古典的確率では各結果が等確率なので、上の式は数え上げによって

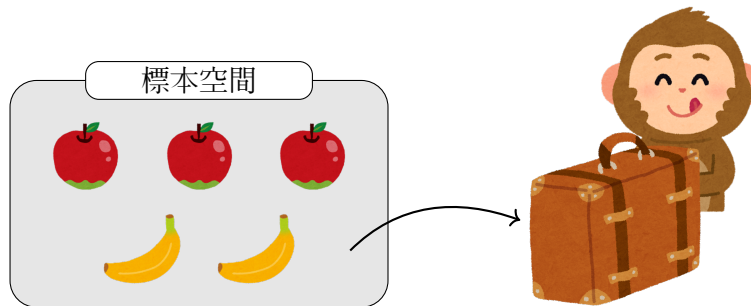
$$P(B | A) = \frac{\text{size of } A \cap B}{\text{size of } A}$$

と計算できます。これは集合論的な共通部分 $A \cap B$ の大きさを数える方法です。すなわち、まず事象 A の中にある結果だけに注目し、そのうちどれだけがさらに B にも属しているかを数えるのです。

条件付き確率は、逐次確率、つまり複数段階からなる実験を考えるときにとくに有用になります。各段階ごとに異なる可能な結果があるような実験です。決定木とは、このような多段階実験の可能な結果を漏れなく並べた図のことで、結果の集まりは枝分かれしていき、木の枝として表されます。¹ 要するに、決定木は可能な結果の空間を整理して、各個別の結果の確率を計算しやすくするための道具なのです。

2.2 例：ジャングルと泥棒

あなたがジャングルにいて、かばんをそのまま放置したとしましょう。すると近くのサルが木から降りてきて、あなたのかばんの中を無作為にあさり始めます。たまたま、昼食用に入っていた果物が5個あり、りんごが3個、バナナが2本入っていました。サルは果物を無作為に盗み、見つけたものを食べてからまたかばんの中を見ます。確率の用語では、これを復元しない抽出といいます。



演習1

- (a) サルが最初にりんごを盗む確率は何か。
- (b) サルが最初にバナナを盗む確率は何か。

¹少し正確に言うと、枝とは木の始点から終点までを結ぶ線のことで、

解答

果物は全部で5個なので、単純に数え上げればよいです。

$$(a) P(\text{最初がりんご}) = \frac{3}{5}.$$

$$(b) P(\text{最初がバナナ}) = \frac{2}{5}.$$

ここで、次の逐次的な問いに注目しましょう。

サルがまずバナナを盗み、そのあとでりんごを盗む確率は何か。

2回目の結果は1回目の結果によって明らかに影響を受けるので、ここでは条件付き確率を使うべきです。

順序をはっきりさせるために、 B_1 を「1回目にバナナ」、 A_2 を「2回目にりんご」とします。すると「バナナ、そのあとりんご」は事象「 B_1 and A_2 」です。

これは条件付き確率の公式を用いて計算できます。ここでは、事象 A を「サルがりんごを取る」、事象 B を「サルがバナナを取る」と対応させます。すると複合事象「 B and A 」は「サルがバナナを取り、そのあとでりんごを取る」と解釈できます。したがって

$$P(B_1 \cap A_2) = P(B_1) \cdot P(A_2 | B_1)$$

を用いて計算できます。

この公式を順に見ていきましょう。

Step 1: すでに $P(B_1) = \frac{2}{5}$ であることを見ました。

Step 2: 次に $P(A_2 | B_1)$ を計算します。バナナが1本取り除かれて（そしてサルに食べられて）いると、2回目に果物を取ろうとするとき、サルは異なる標本空間を持つこととなります。最初の時点では果物は5個ありました。しかしバナナが1本すでに取り除かれているので、今では残りは4個、すなわちりんご3個とバナナ1本だけです。この新しい標本空間の大きさは4であり、サルが選べるりんごは3個あります。したがって条件付き確率は

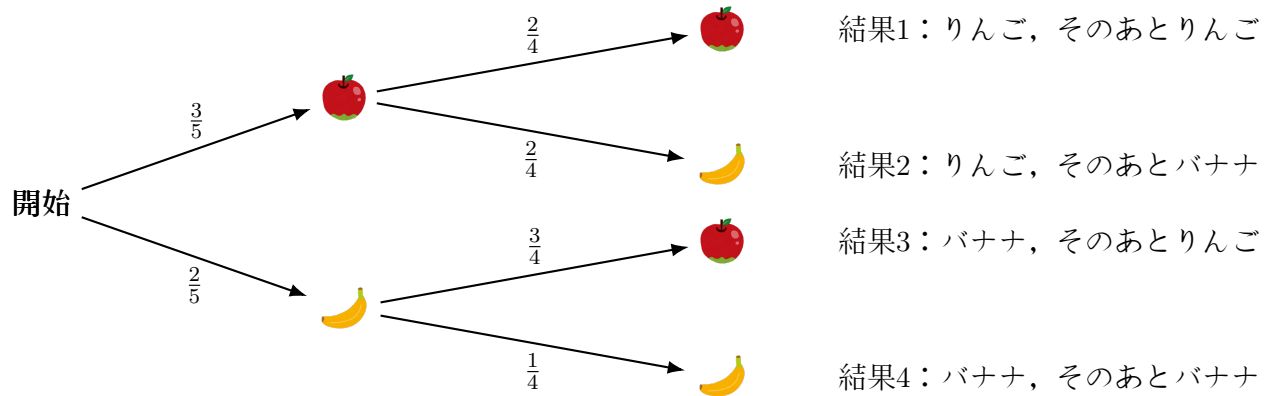
$$P(A_2 | B_1) = \frac{3}{4}$$

です。

Step 3: 最後に、この2つの数を掛けて、「バナナを取ってそのあとりんごを取る」確率を得ます。

$$P(B_1 \text{ and } A_2) = P(B_1) P(A_2 | B_1) = \frac{2}{5} \cdot \frac{3}{4} = \frac{3}{10}.$$

この一連の過程は、下のような木に整理できます。



ご覧のように、先ほど行った計算

$$P(\text{バナナ, そのあとりんご}) = P(B_1) \cdot P(A_2 | B_1) = \frac{2}{5} \cdot \frac{3}{4} = \frac{3}{10}$$

は、上の木で結果3に対応する枝を見れば視覚的にも確認できます。ここでは単に3本目の枝に書かれた数を掛けているだけです。実は、ここで一般的なコツの最初の例を見たこととなります。つまり、多段階実験の結果を木として整理できるなら、どの結果の確率も、その枝に現れるすべての確率を掛けることで求められるのです。

この観察を用いれば、上の木にある他の3つの結果もすぐに計算できます。

- 結果1: $P(A \text{ then } A) = \frac{3}{5} \cdot \frac{2}{4} = \frac{3}{10}$,
- 結果2: $P(A \text{ then } B) = \frac{3}{5} \cdot \frac{2}{4} = \frac{3}{10}$,
- 結果4: $P(B \text{ then } B) = \frac{2}{5} \cdot \frac{1}{4} = \frac{1}{10}$.

これら4つの確率の和は

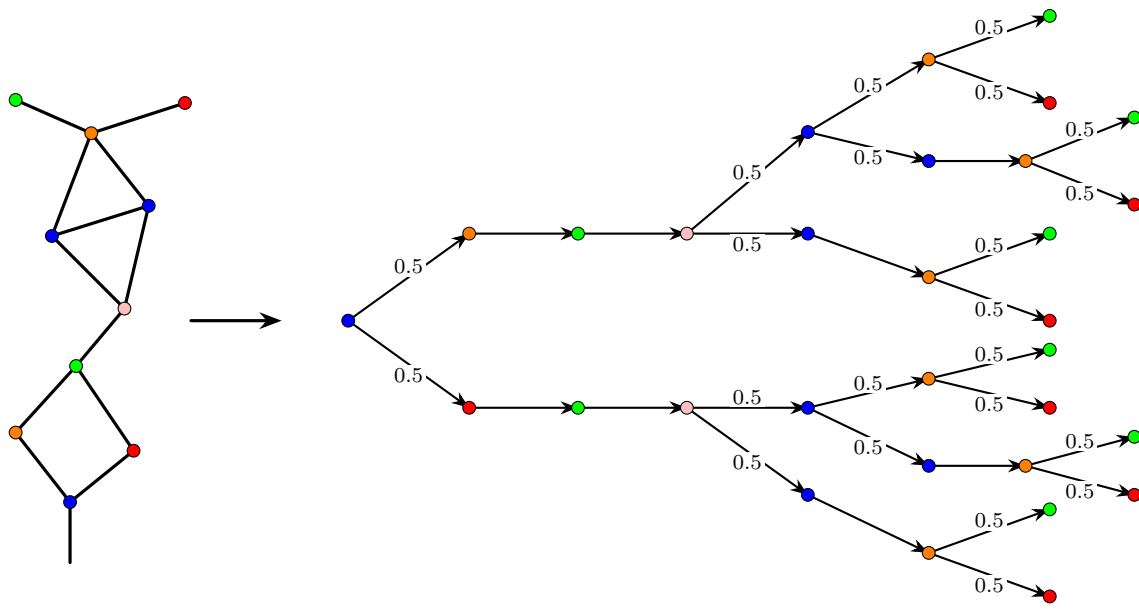
$$\frac{3}{10} + \frac{3}{10} + \frac{3}{10} + \frac{1}{10} = 1$$

になります。これは、4つの結果のうちどれか1つは必ず起こる、ということと同じです。

2.3 岩登りの問題に戻る

今見たように、決定木を使うと、多段階実験のすべての可能な結果を整理できます。ある特定の経路の確率を求めるには、その経路に沿って書かれた確率をすべて掛ければよいのです。この方法を使って、岩場を登る人が壁を進む各経路の確率を求めることができます。

まずこの系に対する木を書き、各「選択点」に仮定された50-50の確率を書き込みます。完全な木は下の通りです。



この木は、各経路が順番に並ぶように整理されています。最初の枝が経路1, 2 本目の枝が経路2, という具合です。ここでも、ある経路は他の経路より複雑であることに気づくはずですが。たとえば経路12 は、経路9 や10 より単純に見えます。これは、壁を登るルートの中には3 回の選択だけで済むものもあれば、4 回の選択を必要とするものもあるからです。実際、ピンクの節点で左へ曲がる選択を含む経路は、右へ曲がる選択を含む経路より複雑になります。

いずれにせよ、決定木の理論を使えば、岩場を登る任意の経路の確率を計算できます。さきほどのサルの場合と同様に、ここではその枝に書かれている値を掛け合わせれば全体の確率が得られます。まとめると以下の通りです。

各経路の確率

各経路で必要な選択回数を区別します。

- 経路1,2,3,4,7,8,9,10 はすべて4 回の選択を必要とします。決定木では、これらの経路に対応する枝に0.5 が4 回現れるので、全体の確率は

$$\frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{16}$$

となります。

- 経路5,6,11,12 は全部で3 回の選択しか必要としません。したがって、これらの経路をたどる確率は他より大きく、

$$\frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{8}$$

となります。

計算が正しいかどうかは、すべての答えの和が1 になることを確かめればいつでも確認できます。実際、

$$8 \left(\frac{1}{16} \right) + 4 \left(\frac{1}{8} \right) = \frac{1}{2} + \frac{1}{2} = 1$$

で、正しいことが分かります。ここで私たちは偶然もうひとつの事実も見つけています。登る人が単純に見える経路のどれかをたどる確率は50%であり、そうならない確率も50%です。

3 結果の数え上げ

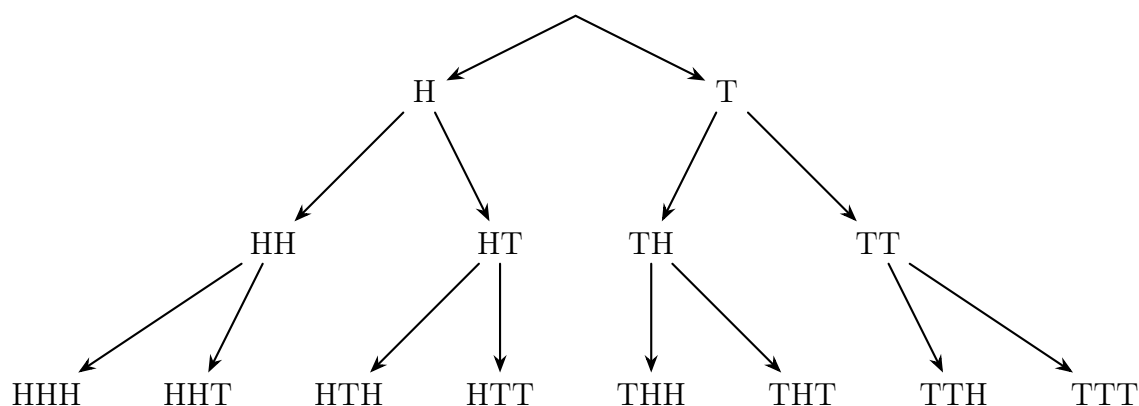
3.1 コインを3回投げる

コインを3回投げるとしましょう。

演習2

- (a) ちょうど2回裏が出る確率は何か。
- (b) ちょうど1回裏が出る確率は何か。

この標本空間では各結果が等確率なので、ここでは単純にそれぞれの場合について有利な結果の数を数えるのが最善の方法です。もちろん、必要ならこれも木として整理できます。あとで便利なように、この木は実験の始まりが上に来るように描いておきます。



演習2の解答

上の木から分かるように、3回のコイン投げの標本空間は大きさ8の集合です。ここではすべての結果が等確率なので、直接数え上げれば確率を計算できます。

- (a) ちょうど2回裏：この事象に有利な結果は3つ、すなわち{HTT, THT, TTH}なので、

$$P(\text{ちょうど2回裏}) = \frac{3}{8}.$$

- (b) ちょうど1回裏：この事象に有利な結果は3つ、すなわち{THH, HTH, HHT}なので、

$$P(\text{ちょうど1回裏}) = \frac{3}{8}.$$

3.2 n 回のコイン投げで裏を数える

前の例を見ると、コイン投げに関する確率の計算はとても簡単だと思えるかもしれません。し

しかし、回数が増えるにつれて状況はかなり複雑になります。まず、上の確率木の各段の大きさが時間とともに増えていくことに注意しましょう。1 段目には2つの結果、2 段目には4つの結果、3 段目には8つの結果があります。もう1回コインを投げると、すでにある各HとTの列に対してさらに2つの結果が生じます。したがって可能な結果の数は再び2倍になり、4 段目には16個の結果が現れます。つまり、4回のコイン投げのあり方は16通りあります。一般に、コインを n 回投げると、全結果数は 2^n になります。

では、 n 回投げたあとに特定の枚数の裏、たとえばちょうど r 回の裏を数えたいとしましょう。原理的には、前と同じように、ちょうど r 回裏がある結果の数を数え、それを 2^n で割ればよいのです。しかし、この作業はすぐに大変になります。たとえば、6回のコイン投げの結果のうち、ちょうど2回裏であるものはいくつあるのでしょうか。これは、HとTの64通りの結果の中から、Tが2個、Hが4個あるものを探し出す作業になります。²

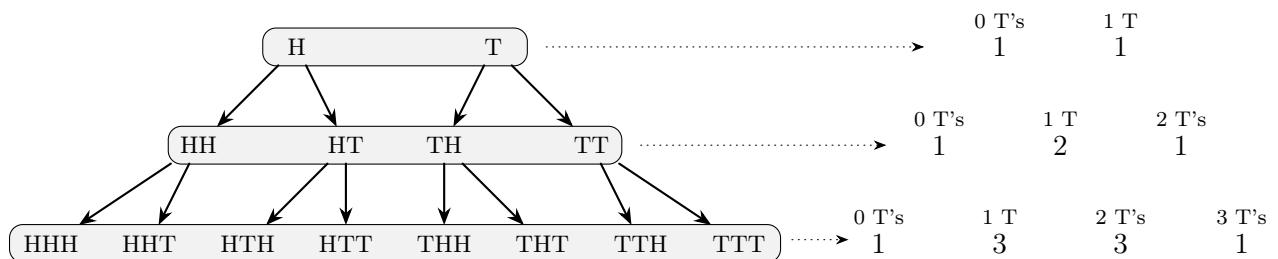
実際には、HとTを大量に書き出してそれを読み取らなくても、ちょうど r 回裏となる結果の数を数えられる巧妙な方法があると便利です。実はそのような方法はあります。しかも、それを私たちはすでに見ています。

3.2.1 n 回の投げでちょうど r 回裏を数える

この方法を説明するために、もう一度上の木を眺め、そこに現れる結果を数えてみましょう。

- 1 回投げ：コインを1回投げると、結果集合は $\{H, T\}$ です。Tを0個含む結果が1つ、Tを1個含む結果が1つあります。
- 2 回投げ：2回投げると結果は $\{HH, HT, TH, TT\}$ です。Tを0個含む結果は1つ (HH)、Tを1個含む結果は2つ (HT と TH)、Tを2個含む結果は1つ (TT) あります。
- 3 回投げ：このとき結果集合の大きさは8で、上の木の最下段に当たります。この数え上げはすでに行いました。Tを0個含む結果は1つ (HHH)、Tを1個含む結果は3つ (HHT, THH, HTH)、Tを2個含む結果は3つ (HHT, THT, TTH)、Tを3個含む結果は1つ (TTT) です。

これをもっと示唆的にするために、木の横に「Tの個数カウンター」を書いてみます。



ご覧のように、図の右側にある三角形はパスカルの三角形の最初の数行にとってもよく似ています。実際、これはまさにパスカルの三角形です。なぜでしょうか。コイン投げでちょうど r 回裏が出る結果を数えることは、「 n 個の位置があって、そのうち r 個を裏として選ぶ方法をすべて見たい」と言っているのと同じだからです。

²これが楽だと感じるなら、 $n = 10$, $r = 4$ を考えてみてください。

パスカルの三角形を、「何個の裏があり得るか」という注記つきで描き直せば、面倒な数え上げをしなくても数を読み取ることができます。

${}^0 T\text{'s}$ 1	$2^0 = 1$ 通り
${}^0 T\text{'s}$ 1 1 ${}^1 T$	$2^1 = 2$ 通り
${}^0 T\text{'s}$ 1 ${}^1 T$ 2 ${}^2 T\text{'s}$ 1 2 1	$2^2 = 4$ 通り
${}^0 T\text{'s}$ 1 ${}^1 T$ 3 ${}^2 T\text{'s}$ 3 ${}^3 T\text{'s}$ 1 3 3 1	$2^3 = 8$ 通り
${}^0 T\text{'s}$ 1 ${}^1 T$ 4 ${}^2 T\text{'s}$ 6 ${}^3 T\text{'s}$ 4 ${}^4 T\text{'s}$ 1 4 6 4 1	$2^4 = 16$ 通り
${}^0 T\text{'s}$ 1 ${}^1 T$ 5 ${}^2 T\text{'s}$ 10 ${}^3 T\text{'s}$ 10 ${}^4 T\text{'s}$ 5 ${}^5 T\text{'s}$ 1 5 10 10 5 1	$2^5 = 32$ 通り
${}^0 T\text{'s}$ 1 ${}^1 T$ 6 ${}^2 T\text{'s}$ 15 ${}^3 T\text{'s}$ 20 ${}^4 T\text{'s}$ 15 ${}^5 T\text{'s}$ 6 ${}^6 T\text{'s}$ 1 6 15 20 15 6 1	$2^6 = 64$ 通り

パスカルの三角形を使うほうが、H と T のすべての組合せを書き出すよりずっと簡単です。たとえば6回のコイン投げでちょうど2回裏が出る確率を求めるには、三角形の $n = 6$ の行で $r = 2$ の項を見れば、ちょうど2回裏になる結果が15通りあることが分かります。そして全結果数は $2^6 = 64$ です。したがって、これらをすぐに数え上げの公式に代入して

$$P(6 \text{ 回投げてちょうど} 2 \text{ 回裏}) = \frac{6 \text{ 回の投げでちょうど} 2 \text{ 回裏になる方法の数}}{6 \text{ 回投げたときの結果の総数}} = \frac{15}{64} \approx 23.4\%$$

と結論できます。

パスカルの三角形の各項は $C_{n,r}$ です。一般に、この「組合せを数える」方法を使うと、 n 回のコイン投げでちょうど r 回裏が出る確率の一般公式

$$P(n \text{ 回投げてちょうど} r \text{ 回裏}) = \frac{n \text{ 回の投げでちょうど} r \text{ 回裏になる方法の数}}{n \text{ 回投げたときの結果の総数}} = \frac{C_{n,r}}{2^n}$$

を得ます。

もちろん、表が出る確率と裏が出る確率は同じなので、上の公式は対称的であり、 n 回の投げでちょうど r 回表が出る確率にもそのまま使えます。

4 おまけ：誕生日のパラドックス

A を「部屋の中に誕生日が同じ2人かいる」という事象とします。このとき、確率 $P(A)$ は部屋にいる人数 n によって大きく変わります。たとえば $n = 2$ なら、この2人の誕生日が同じである確率は極めて小さいはずですが、しかし $n > 365$ なら、誕生日は365通りしかないので、部屋の中に誕生日が同じ2人かいることは必ず起こります。このように、 n がとても小さいときは $P(A)$ も小さく、 n が大きくなると $P(A)$ はついには1になります。ここで次の問いを考えます。

どの n において、 $P(A)$ は50%を超えるのか。

言い換えると、部屋の中に何人いれば、「2人が同じ誕生日である」ことのほうが、そうでないことより起こりやすくなるのでしょうか。実は答えは驚くべきもので、必要なのはたった23人です。

この驚くべき事実は「誕生日のパラドックス」と呼ばれます。これは、この講義の前のほうで扱った床屋のパラドックスのような意味でのパラドックスではありません。あちらでは、真と偽が無限ループして論理の働きそのものが破綻しました。それに対して誕生日のパラドックスは、「真なる逆説 (veridical paradox)」の一例です。つまり、数学的には正しい結果なのに、あまりに意外なので間違っているように感じられるものです。しかしこれから示すように、誕生日のパラドックスは確率の基本的な応用からきちんと導かれます。

4.1 誕生日のパラドックスを解く

誕生日のパラドックスを解くために、問題の見方をひっくり返して考えます。 $P(A)$ を直接計算する代わりに、 A の補事象の確率 $P(A^c)$ を計算します。補事象は確率論における「not」に対応するものでした。ここでは補事象 A^c は「部屋の中の全員の誕生日がすべて異なる」という意味です。

補事象の確率は単に $1 - P(A)$ であることを思い出しましょう。したがって、もし $P(A^c) < 50\%$ なら、 $P(A) > 50\%$ です。なぜなら

$$P(A) + P(A^c) = 1$$

だからです。各誕生日は等確率で選ばれると仮定し、どの誕生日も確率 $1/365$ を持つとしましょう。まずは小さな n について $P(A^c)$ を計算して、どのような過程になるのかを見ます。

$n = 2$ の場合、求めたいのは2人の誕生日が異なる確率です。明らかに1人目の誕生日は自由に選べます。しかし2人の誕生日が異なるためには、2人目の誕生日は1人目の誕生日を避けなければなりません。1人目の誕生日と異なる日は364日あるので、(これをランダムな事象として扱うなら) 2人目の誕生日が異なる確率は

$$P(A^c) = \frac{1 \text{ 人目の誕生日と異なる日の数}}{1 \text{ 年の日数の総数}} = \frac{364}{365} \approx 0.997 = 99.7\%$$

となります。

次に3人目を部屋に加える、すなわち $n = 3$ の場合を考えましょう。2人の誕生日が異なる確率はすでに $\frac{364}{365}$ であることが分かっています。しかしさらに、3人目の誕生日が1人目と2人目の両方の誕生日と異ならなければなりません。この次の人の誕生日をランダムに選ぶことは条件付きの事象です。なぜなら、その確率は前の選択の結果に依存するからです。このような「and」の確率を計算するには、対応する確率を掛け合わせます。

$$\begin{aligned} P(A^c) &= \frac{364}{365} \times \frac{1 \text{ 人目と2人目の誕生日の両方と異なる日の数}}{1 \text{ 年の日数の総数}} \\ &= \frac{364}{365} \cdot \frac{363}{365} \approx 0.992 = 99.2\%. \end{aligned}$$

このパターンはさらに大きな n に対しても繰り返せます。 n を1増やすたびに、 $(n+1)$ 人目が選べる誕生日の候補を少しずつ「使い切って」いくことになります。したがって、各段階でだんだん小さい数を掛けていかなければなりません。たとえば

$$n = 4 \text{ のとき: } P(A^c) = \left(\frac{364}{365}\right) \left(\frac{363}{365}\right) \left(\frac{362}{365}\right) \approx 0.983 = 98.3\%$$

$$n = 5 \text{ のとき: } P(A^c) = \left(\frac{364}{365}\right) \left(\frac{363}{365}\right) \left(\frac{362}{365}\right) \left(\frac{361}{365}\right) \approx 0.973 = 97.3\%$$

$$n = 6 \text{ のとき: } P(A^c) = \left(\frac{364}{365}\right) \left(\frac{363}{365}\right) \left(\frac{362}{365}\right) \left(\frac{361}{365}\right) \left(\frac{360}{365}\right) \approx 0.960 = 96.0\%$$

となります。

一般の n に対しては、

$$P(A^c) = \left(\frac{365}{365}\right) \left(\frac{364}{365}\right) \left(\frac{363}{365}\right) \cdots \left(\frac{365 - (n - 1)}{365}\right)$$

という公式を得ます。

このパラドックスは、確率 $P(A^c)$ が50%を下回る最小の n を見つければ解決されます。なぜなら、それは $P(A)$ が50%を上回ることを意味するからです。公式を使って n を大きくしていくと、 $n = 23$ のとき

$$P(A^c) = \left(\frac{365}{365}\right) \left(\frac{364}{365}\right) \left(\frac{363}{365}\right) \cdots \left(\frac{342}{365}\right) \approx 0.490 < 0.5$$

となることが分かります。したがって、 $n = 23$ のとき $P(A) > 50\%$ です。