

ЗАДАЧИ С УЛИЦЫ

Охота на принцесс

Каждый родитель сталкивается с этой задачей, когда его ребенок получает первую игрушку из киндер-сюрприза. Трудно назвать шоколадом массу, из которой сделано яйцо, но это неважно — вся радость внутри. В желтом контейнере прячется игрушка, которую нужно собрать, например, приделать клюшку пингвину-хоккеисту или нацепить плащ на бэтмена. Но главное — игрушки образуют коллекции. Если ребенок получил одну фигурку, то не успокоится, пока не добудет все. Это замечательный маркетинговый ход, который придуман отнюдь не компанией «Kinder». В начале XX века американская компания «Dixie», выпускавшая солдатские котелки¹ и прохладительные напитки, использовала этот трюк, чтобы привлечь покупателей. Под крышечкой бутылочки покупатель находил купон с изображением персонажа цирка «Dixie». Счастливый обладатель полной коллекции купонов получал одну бутылку с газировкой бесплатно. Потом в ход пошли купоны с игроками бейсбольных команд и другими популярными личностями.

Коллекция из десяти киндер-бегемотов уже стала классикой (рис. 1).



Рис. 1

Еще были пингвины, акулы, крокодилы, феечки с крылышками. Мы будем собирать коллекцию из восьми принцесс (рис. 2).



Рис. 2

Проблема в том, что, в отличие от коллекционирования марок, киндер-коллекционирование происходит вслепую: вы не видите, какую игрушку покупаете. Принцессы начинают повторяться, выкинуть дубликаты жалко, и потому они скапливаются в ящиках и на полках, умножая количество домашнего хлама.

Чтобы упростить словоупотребление, будем называть *экспонатами* тех принцесс, которые появились впервые и тем самымполнили коллекцию.

Сам по себе процесс такого слепого коллекционирования — это невероятно увлекательный математический сюжет, если его немного формализовать. Например, *будем предполагать, что в каждом следующем купленном яйце может с равными шансами оказаться любая из фигурок коллекции.*

¹ В английском языке слово dixie стало нарицательным — солдатский походный котелок.

Опять принцеску притащил! Не надоело?!
Мама, такой у меня еще не было



Классическая задача коллекционера

Задача 1. Сколько потребуется купить киндер-сюрпризов, чтобы собрать полную коллекцию? Разумеется, речь о математическом ожидании.

По причинам, обсуждавшимся выше, эта классическая задача коллекционера известна еще под названиями «The Coupon Collector's Problem» и «Dixie Cup Problem». Первое решение, вероятно, дал Абрахам Муавр в 1712 году, но опубликовано оно было лишь спустя сто лет Лапласом (Laplace, Pierre-Simon (1812), «Théorie analytique des probabilités»).

Предварительные сведения. Рассмотрим случайный опыт, в котором производятся одинаковые независимые испытания до достижения первого успеха. Например, монету бросают, пока не выпадет орел. Или игральную кость бросают, пока не выпадет единица. Или мобильный телефон передает СМС из леса, совершая попытки, пока не достигнет успеха в условиях почти полного отсутствия связи. Если вероятность успеха p при каждой отдельной попытке хоть чуть-чуть больше нуля, то рано или поздно успех наступит.

Пусть X — число испытаний в такой серии, то есть число испытаний, проведенных до момента, когда наступил первый успех (включая последнее успешное испытание). Такая величина имеет *геометрическое распределение*². Нам потребуются математическое ожидание и дисперсия геометрического распределения:

$$EX = \frac{1}{p}. \quad (1)$$

$$DX = \frac{1-p}{p^2} \text{ и } \sqrt{DX} = \frac{\sqrt{1-p}}{p}. \quad (2)$$

Утверждение (1) практически очевидно, во всяком случае, интуитивно ясное. Например, при бросании игрального кубика вероятность единицы p равна $\frac{1}{6}$, поэтому, чтобы получить единицу, в среднем нужно $\frac{1}{p} = 6$ бросков.

Второе утверждение вовсе не очевидно, и равенство (2) имеет весьма «подлую» сущность. В самом деле, если вероятность p велика, скажем $\frac{2}{3}$,

то для достижения успеха потребуется в среднем 1,5 попытки, а реальное число попыток редко когда превысит 4, поскольку стандартное отклонение мало:

$$\sqrt{DX} = \frac{\sqrt{\frac{1}{3}}}{\frac{2}{3}} = \frac{\sqrt{3}}{2} = 0,866\dots,$$

и потому

$$EX + 3\sqrt{DX} \approx 3,264\dots$$

А если p мало, то числитель $\sqrt{1-p}$ почти не отличается от единицы и можно считать, что $\sqrt{DX} \approx \frac{1}{p}$. Поэтому реальное число попыток

может сильно превышать среднее. Например, если $p = 0,1$, то

$$EX = 10 \text{ и } \sqrt{DX} = \frac{\sqrt{0,9}}{0,1} \approx 9,486\dots$$

В таком случае вполне может случиться и 20, и 30, и даже больше неудачных попыток, поскольку

$$EX + 3\sqrt{DX} \approx 38,46.$$

Наш мозг, обычно хорошо оценивающий средние значения, в этой ситуации отказывается столь же хорошо оценивать возможные отклонения. Если мы ждем успеха при 10-й попытке, то 35 неудач подряд покажутся нам чем-то из ряда вон выходящим, хотя на самом деле это в порядке вещей.

Покончив со вступлением, перейдем к самой задаче коллекционера.

Решение задачи коллекционера

Разумеется, будем считать, что принцесс в коллекции не 8, а n : решать задачу в общем случае проще.

Пусть X_k — число испытаний (покупок) нужных, чтобы получить k -й экспонат коллекции после того, как предыдущий уже получен. В частности, $X_1 = 1$, поскольку для получения первой принцессы, кто бы она ни была, требуется ровно один киндер-сюрприз. Значит, $EX_1 = 1$, математическое ожидание константы равно ей самой.

С X_2 уже не все так очевидно, поскольку при покупке второго яйца может повториться первая принцесса, тогда нового экспоната не будет и коллекция не пополнится. Однако это маловероятно, поскольку вероятность p_2 того, что в каждом следующем яйце окажется отсутствующий экспонат, высока:

$$p_2 = \frac{n-1}{n}.$$

Попытки продолжаются, пока не случится успех — второй экспонат, а потому этот опыт является серией испытаний до первого успеха. Значит,

$$EX_2 = \frac{1}{p_2} = \frac{n}{n-1}.$$

Как только появился второй экспонат, начинается охота за третьим. Она потребует X_3 киндер-сюрпризов и тоже является серией испытаний до первого успеха. Вероятность того, что

² Вероятности того, что $X = k$, образуют геометрическую прогрессию $(1-p)^{k-1}p$, где p — вероятность успеха в каждом отдельном испытании. Это распределение обычно обозначают $G(p)$. Можно записать $X \sim G(p)$.

каждая следующая покупка принесет какую-нибудь ныне отсутствующую принцессу, равна $\frac{n-2}{n}$, поэтому

$$EX_3 = \frac{n}{n-2}.$$

И так далее. Последний экспонат — последняя нужная принцесса — будет приобретен в результате X_n попыток, последовавших за появлением предпоследней принцессы. Ясно, что

$$EX_n = \frac{1}{p_n} = n.$$

Общее число испытаний складывается из всех этих чисел:

$$X = X_1 + X_2 + \dots + X_n.$$

Следовательно,

$$EX = EX_1 + EX_2 + \dots + EX_n = 1 + \frac{n}{n-1} + \frac{n}{n-2} + \dots + \frac{n}{2} + \frac{n}{1}.$$

Единица в этой сумме выглядит чужеродно. Однако чужеродность легко устранить:

$$EX = \frac{n}{n} + \frac{n}{n-1} + \frac{n}{n-2} + \dots + \frac{n}{2} + \frac{n}{1} = n \left(1 + \frac{1}{2} + \dots + \frac{1}{n-1} + \frac{1}{n} \right).$$

Выражение в скобках $1 + \frac{1}{2} + \dots + \frac{1}{n}$ называется n -м гармоническим числом. Обычное обозначение H_n . Поэтому результат можно записать коротко:

$$EX = nH_n. \quad (3)$$

Например, поскольку принцесс в полной коллекции восемь, в среднем киндер-сюрпризов требуется

$$EX = 8H_8 = 8 \cdot \left(1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{8} \right) = 21,742\dots,$$

чтобы собрать всю коллекцию. Но это только в среднем. Насколько реальное число может отличаться от этого среднего?

Возьмем в качестве решающего правила три стандартных отклонения, как мы это уже делали, и посмотрим, какова верхняя граница разумного.

К счастью, величины X_k попарно независимы, поскольку число попыток, нужных, чтобы получить очередную принцессу, не зависит от того, сколько времени мы охотились на предыдущую. Поэтому дисперсию тоже можно вычислять по слагаемым³. Мы помним, что $X_1 = 1$, поэтому $DX_1 = 0$. Для вычисления остальных дисперсий нужно воспользоваться формулой (3):

$$DX_2 = \frac{1 - \frac{n-1}{n}}{\left(\frac{n-1}{n}\right)^2} = \frac{n}{(n-1)^2}.$$

Аналогично

$$DX_3 = \frac{2n}{(n-2)^2},$$

и так далее до

$$DX_n = \frac{(n-1)n}{1^2}.$$

³ Дисперсия суммы независимых случайных величин равна сумме их дисперсий (если они существуют).

Тогда

$$\begin{aligned} DX &= DX_1 + DX_2 + \dots + DX_n = \\ &= 0 + \frac{n}{(n-1)^2} + \frac{2n}{(n-2)^2} + \dots + \frac{(n-1)n}{1^2} = \\ &= n \left(\frac{n-1}{1^2} + \frac{n-2}{2^2} + \dots + \frac{1}{(n-1)^2} + \frac{0}{n^2} \right) = \\ &= n^2 \left(1 + \frac{1}{4} + \dots + \frac{1}{(n-1)^2} + \frac{1}{n^2} \right) - n \left(1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n-1} + \frac{1}{n} \right). \end{aligned}$$

Число в первых скобках принято обозначать $H_{n,2}$ и называть гармоническим числом 2-го порядка. Пользуясь этими обозначениями, можно записать:

$$DX = n^2 H_{n,2} - nH_n. \quad (4)$$

В частности, при $n = 8$ получаем:

$$DX = 64H_{8,2} - 8H_8 = 76,012\dots \text{ и } \sqrt{DX} = 8,718\dots$$

Значит,

$$EX + 3\sqrt{DX} = 8H_8 + 3\sqrt{64H_{8,2} - 8H_8} = 47,898\dots$$

Если кому-то ради коллекции из восьми принцесс придется купить 47 киндер-сюрпризов, то он не вправе винить судьбу. Больше число следует считать либо невезением, либо признаком того, что, вопреки нашему предположению, принцессы распределены по киндер-сюрпризам *неравномерно*: то есть одни встречаются чаще других.

Решение с помощью MS Excel и моделирование

Несложно составить электронную таблицу для вычисления среднего и дисперсии числа попыток в классической задаче коллекционера (рис. 3).

		Допуст.отклонение (ст.откл.)		3
Задача коллекционера				
		Объем коллекции n=		200
		EX=	1175,606	Ст.откл.= 253,815
		DX=	64422,256	Макс.= 1937,052
n=	n^2=	EX=	DX=	Ст.откл.
1	1	1,000	0,000	0,000
2	4	=C11/CPГАРМ(\$B\$10:B11)		
3	9	5,500	6,750	2,598
4	16	8,333	14,444	3,801
5	25	11,417	25,174	5,017
6	36	14,700	38,990	6,244
7	49	18,150	55,928	7,479
8	64	21,743	76,012	8,718
9	81	25,461	99,260	9,963
10	100	29,290	125,687	11,211

Рис. 3

Для вычисления гармонических чисел можно воспользоваться сложением в последовательных ячейках. Однако MS Excel имеет стандартную функцию СРГАРМ(), которая вычисляет среднее гармоническое аргументов:

$$\text{СРГАРМ}(x_1; \dots; x_n) = \frac{n}{\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_n}}$$

Таким образом, чтобы найти, например, H_8 , достаточно написать формулу

$$= \frac{8}{\text{СРГАРМ}(\text{диапазон})},$$

где *диапазон* — диапазон ячеек, содержащих числа от 1 до 8.

Таким же способом удобно вычислять суммы обратных квадратов, то есть гармонические числа второго порядка.

Математическое решение задачи сводится к поиску гармонических чисел. Но вот компьютерное моделирование выглядит не так просто: нужно получать случайные натуральные числа из отрезка от 1 до n до тех пор, пока каждое число не появится хотя бы один раз.

В MS Excel это сделать можно, но лишь для фиксированного n , если не пользоваться макросами. Первый способ слишком бедный, а второй нарушает наш принцип минимальности средств. Вместо MS Excel для моделирования можно воспользоваться любым языком программирования. Листинг программы на языке Pascal показан ниже:

{Язык: Pascal}

Программа моделирует задачу коллекционера. Запрашивает объем коллекции n и возвращает последовательность случайных натуральных чисел от 1 до n до тех пор, пока каждое число не появится хотя бы раз, а также число потребовавшихся испытаний}

Program Collector_Problem;

```
var
Number: integer;           {объем коллекции}
Exhibits: array[1..200] of 0..1; {какие экспонаты уже есть}
Count: integer;           {счетчик экспонатов}
Trials: integer;          {счетчик попыток}
Princess: integer;        {текущий случайный элемент}

i: integer;
```

```
begin
randomize;
write ('Объем коллекции: '); {запрос объема коллекции}

readln (Number);
Count:=0; Trials:=0;
```

```
for i:= 1 to Number do Exhibits[i]:= 0;
{инициализация массива}
while Count<Number do begin {пока коллекция не собрана}
Princess:= Trunc(random(Number))+1;
{выбирай случайное число}

Trials:=Trials+1;
write(Princess, ' '); {печать очередного элемента}
if Exhibits[Princess]=0 then begin {если новый экспонат}
Exhibits[Princess]:=1; Count:=Count+1;
{то добавь его в коллекцию}

end;
end;
writeln (''); writeln ('Испытаний=', Trials);
{печать числа испытаний}
end.
```

На странице цикла «Задачи с улицы» сайта лаборатории теории вероятностей (https://ptlab.mccme.ru/Street_problems), где публикуются все файлы MS Excel к задачам цикла, размещен работающий код этой же программы на языке Javascript.

Гармонические числа и приближенное решение задачи коллекционера

Существует приближенная формула для решения задачи коллекционера, не требующая прямого вычисления гармонических чисел. Она основана на том, что гармонические числа H_n мало отличаются от натурального логарифма $\ln n$. А именно

$$\lim_{n \rightarrow \infty} (H_n - \ln n) = \gamma = 0,57721566\dots$$

Число γ , к которому приближается разность между гармоническим числом и логарифмом, называют *константой Эйлера–Маскерони*. Точное вычисление γ требует математического анализа, далеко выходящего за рамки школы, но показать, что

$$\ln n < H_n < \ln n + 1,$$

несложно. Сделаем это на примере $n = 6$. На рисунке 5,а видно, что число

$$1 + \frac{1}{2} + \dots + \frac{1}{6} = H_6$$

больше, чем площадь, заключенная под графиком функции $y = \frac{1}{x}$ на отрезке $[1; 6]$, то есть

$$H_6 = \int_1^6 \frac{dx}{x} = \ln x \Big|_1^6 = \ln 6.$$

На рисунке 5,б видно, что число

$$\frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{6} = H_6 - 1$$

меньше, чем та же самая площадь:

$$H_6 - 1 < \ln 6.$$

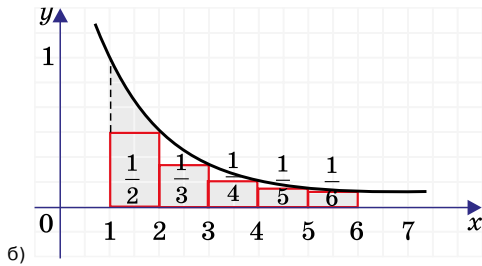
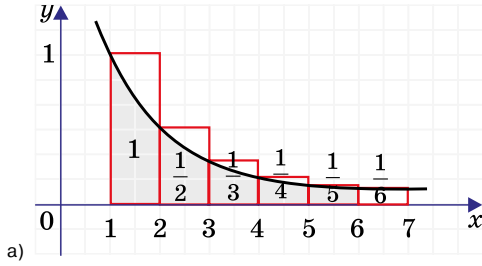


Рис. 5

Таким образом,

$$\ln 6 < H_6 < \ln 6 + 1.$$

Это двойное неравенство верно для любого n , а не только для $n = 6$. Точное равенство

$$H_n = \ln n + \gamma + \frac{1}{2n} + o\left(\frac{1}{n}\right),$$

где символ $o\left(\frac{1}{n}\right)$ обозначает, как обычно, функцию, бесконечно малую по отношению к $\frac{1}{n}$.

А потому можно смело считать, что даже для не очень больших n

$$EX = nH_n \approx n \ln n + 0,58n + 0,5. \quad (5)$$

Для сравнения можно добавить соответствующий столбец в таблицу MS Excel с тем, чтобы посмотреть, велико ли отличие между точной формулой и приближенной (рис. 6, столбец G).

А	В	С	Д	Е	Г	Н
Задача коллекционера						
Объем коллекции n=			200	Допуст.отклон (ст.откл.)		
EX=		1175,606	Ст.откл.=		253,815	
DX=		64422,256	Макс.=		1937,052	
n=	n^2=	EX=	DX=	Ст.откл.	n ln n + n + 0,5	
1	1	1,000	0,000	0,000	1,08	
2	4	3,000	2,000	1,414	3,046294	
3	9	5,500	6,750	2,598	5,535837	
4	16	8,333	14,444	3,801	8,365177	
5	25	11,417	25,174	5,017	11,44719	
6	36	14,700	38,990	6,244	=B15*LN(B15)+0,58*B15+0,5	
7	49	18,150	55,928	7,479	18,18137	
8	64	21,743	76,012	8,718	21,77553	

Рис. 6

Видно, что уже для $n = 8$ отличие около 0,03, а дальше оно только уменьшается.

Еще интереснее обстоят дела с приближенной формулой для дисперсии. Точно так же, как мы показали, что $\ln n < H_n < \ln n + 1$, можно показать, что

$$1 - \frac{1}{n} < H_{n,2} = 1 + \frac{1}{4} + \frac{1}{9} + \dots + \frac{1}{n^2} < 2 - \frac{1}{n},$$

и потому $H_{n,2} < 2$. Но эта оценка грубая. Известно, что последовательности гармонических чисел 2-го порядка сходится:

$$\lim_{n \rightarrow \infty} H_{n,2} < \frac{\pi^2}{6} = 1,6449\dots$$

и для любого n верно, что

$$H_{n,2} < \frac{\pi^2}{6} = 1,6449\dots$$

(доказательство этого удивительного факта выходит далеко за рамки школьной математики). Поэтому универсальной оценкой сверху для стандартного отклонения является неравенство

$$\sqrt{DX} < \sqrt{\frac{\pi^2}{6} n^2 - n \ln n - 0,58n - 0,5}. \quad (6)$$

Дефицит коллекции

В процессе собирательства принцесс возникают разные любопытные вероятностные процессы. Например, можно поставить вопрос о том, каково математическое ожидание недостающих принцесс в момент, когда куплено очередное яйцо. Число недостающих экспонатов будем называть *дефицитом*.

Задача 2. Сколько принцесс не хватает в коллекции после приобретения 10-го киндер-сюрприза? Найти математическое ожидание этой случайной величины. Найти приближенное значение при больших n .

Решение задачи 2

Разумеется, будем решать задачу в общем виде: объем коллекции n , а число накопленных элементов последовательности (купленных киндер-сюрпризов) пусть будет k . Пронумеруем принцесс каким-нибудь образом (пусть, скажем, Аврора будет № 1, Бель — № 2 и т.п.). Для каждого из n экспонатов введем индикатор I_j события «этот экспонат отсутствует» по следующему правилу:

$$I_j = \begin{cases} 1, & \text{если экспонат } j \text{ отсутствует,} \\ 0, & \text{если иначе.} \end{cases}$$

Дефицит легко выражается через индикаторы:

$$\delta = I_1 + I_2 + \dots + I_n.$$

Вероятность того, что $I_j = 1$, равна вероятности того, что ни один из k полученных элементов не является экспонатом j . В силу независимости испытаний эта вероятность равна

$$P(I_j = 1) = \left(\frac{n-1}{n}\right)^k.$$

Следовательно,

$$EI_j = \left(\frac{n-1}{n}\right)^k.$$

Это выражение не зависит от j и поэтому

$$E\delta = n \cdot \left(\frac{n-1}{n}\right)^k.$$

Это равенство дает точное решение задачи.

Чтобы посмотреть, что произойдет при больших n , преобразуем полученное выражение:

$$E\delta = n \cdot \left(1 - \frac{1}{n}\right)^n.$$

Последовательность $\left(1 - \frac{1}{n}\right)^n$ стремится к e^{-1} при $n \rightarrow \infty$. При этом сходимость довольно быстрая. Поэтому при больших n можно полагать, что

$$E\delta \approx ne^{-\frac{k}{n}}. \quad (7)$$

Обобщение задачи на две и более коллекции

Если детей в семье двое и нужно собрать две коллекции на двоих, то задача становится до обидного сложной. Она известна как Double Dixie Cup Problem. Решение для m коллекций было получено в 1960 году Дональдом Ньюманом и Лоуренсом Шеппом. Точная формула для математического ожидания числа нужных покупок выглядит не очень привлекательно, но есть и приближенная, аналогичная формуле (5):

$$E\zeta_m = n(\ln n + (m-1)\ln \ln n + C_m) + o(1).$$

Ньюман и Шепп пишут: «Если первая коллекция «стоит» $n \ln n$, то вторая и каждая последующая стоят $n \ln \ln n$ ». Под стоимостью следует понимать число попыток, то есть затраты коллекционера в подходящих денежных единицах. Природа чисел C_m исследована П. Эрдешем и А. Реньи в 1961 году.

Более подробно о случае двух коллекций можно прочесть в журнале «Математическое просвещение» (сер. 3, вып. 26, 2020, с. 198–220).

Задачи для самостоятельного решения

1. С помощью MS Excel найдите точные значения математического ожидания и стандартного отклонения случайной величины «число попыток, которые потребуются, чтобы собрать коллекцию из 20 экспонатов». Сравните эти значения с числами, которые получаются по приближенным формулам (5) и (6).

2. Ребенок собирает киндер-коллекцию из 14 экспонатов. Сколько в среднем различных экспонатов окажется у него после покупки 20-го киндер-сюрприза?

3. Сестры Аня и Таня собирают две коллекции, каждая из восьми киндер-принцесс. Договорились так: старшая сестра Аня покупает по одному киндер-сюрпризу в день и отдает лишних принцесс в коллекцию Тане. В какой-то момент Аня собрала полную коллекцию.

а) Какова вероятность того, что в этот момент у Тани в коллекции нет еще ни одной принцессы?

б) Что более вероятно: Тане в этот момент не хватает ровно одной принцессы или ровно двух принцесс?

4* Покажите, что математическое ожидание дефицита коллекции является наименьшим, если все принцессы распределены по киндер-сюрпризам с равными вероятностями.

5* В условиях задачи 3 найдите математическое ожидание числа принцесс, которых не хватает в коллекции у Тани в момент, когда Аня собрала всю коллекцию.

Ответы к задачам из статьи

«Задача о такси в Майкопе»

1. а) Не изменится; б) увеличится до 33. **2.** 19. Лишняя информация — номера поездок, когда машина повторялась. Важно только, что их было две. **3.** Функцию наибольшего правдоподобия можно составить как произведение функций, полученных из наблюдений обоих пассажиров. Другой способ: объединить обе серии наблюдений в одну, проверив, нет ли повторяющихся машин в этих сериях. Оба варианта подразумевают, что не было *одновременных* поездок. Если же они были, то при вычислении вероятностей придется учитывать и это. **4.** Пусть k_1 и k_2 — номера первой поездки с повторившимся такси первого и второго пассажиров. Уравнение примет вид $EX = \frac{k_1 + k_2}{2}$. **5.** Например, можно «склеить» последовательности наблюдений обоих пассажиров. Есть и другие способы.

От редакции. В № 8 на с. 51 произошла досадная описка (левая колонка, третий абзац сверху). Следует читать: «Если вы себя относите к математическим пуристам и считаете, что...» Кстати, пурист — приверженец пуризма, человек, выступающий за чистоту нравов, языка и т.п.