

# Баннерная крутилка с весами

---

## Постановка задачи

---

Необходимо из  $n$  различных взвешенных баннеров выбрать  $k$  различных баннеров в соответствии с весами.

**На вход:**

набор взвешенных баннеров

**На выход:**

случайные разные  $k$  баннеров

---

## Подход

---

Интерпретируем сумму весов как отрезок. На отрезке стоят метки: первая метка – вес первого баннера, вторая метка – вес первого + вес второго и так далее. Равновероятно выбираем некоторую точку на отрезке. Тот интервал, в который попала брошенная точка, и будет говорить нам о том, какой баннер показывать.

Прямолинейная реализация – при выборе случайной точки на отрезке каждый раз суммировать веса начиная с первого, пока не превысим случайно выбранное число.

Сложность этой реализации  $O(n)$

Был выбран алгоритм, построенный на основе дерева отрезков. Ниже он будет подробно разобран на примере. А затем приведен полный алгоритм отбора  $k$  разных баннеров. В самом конце приведено замечание про сдвиг вероятностей показа каждого баннера при увеличении числа  $k$ .

---

## Сложность алгоритма:

---

Каждый раз выбирая баннер мы ищем правильный интервал за  $O(\log n)$ .

Также стоит учитывать процесс обмена двух элементов, который приводит к обновлению дерева. Обновление дерева производится за  $O(\log n)$ .

И тогда суммарная сложность алгоритма:  $kO(\log n)$

---

## Пример работы алгоритма.

---

На вход:

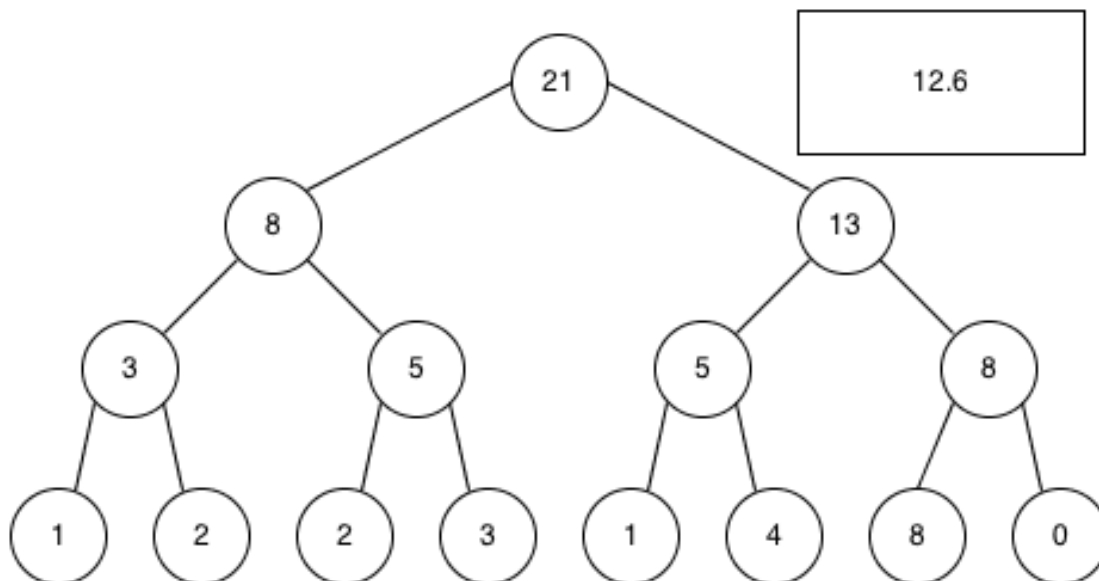
$[('a', 1), ('b', 2), ('c', 2), ('d', 3), ('e', 1), ('f', 4), ('g', 8)]$

Строим дерево отрезков:

$[0, 21, 8, 13, 3, 5, 5, 8, 1, 2, 2, 3, 1, 4, 8, 0]$

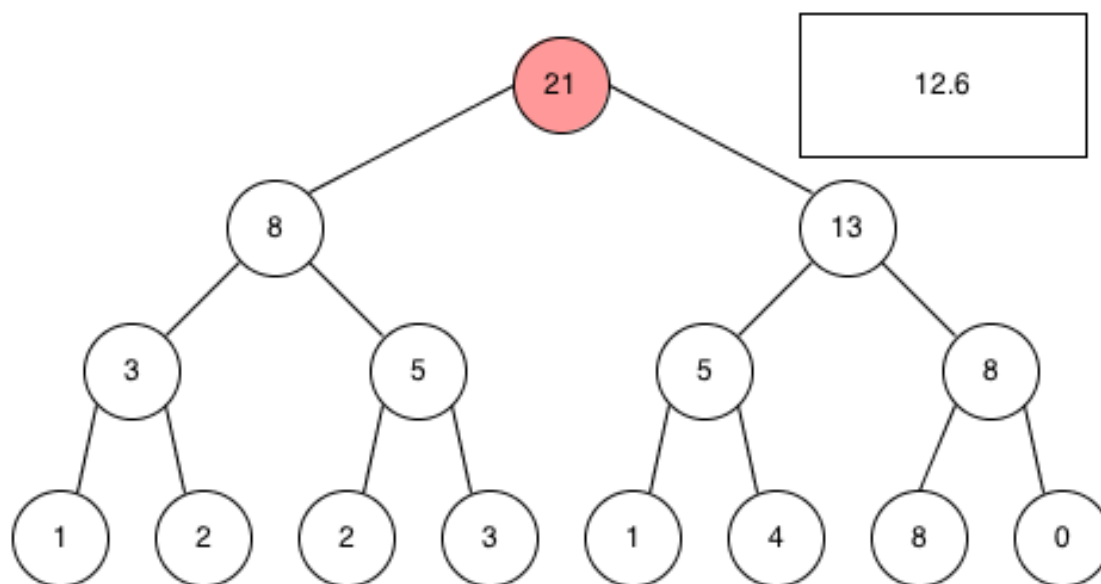
Выбираем случайно число от 0 до 1. Пускай это будет 0.6. Суммарное количество весов: 21.

Значит, нашим генерированным значением считаем  $\text{generated} = 21 * 0.6 = 12.6$



В дереве отрезков нам необходимо найти наименьший отрезок, который начинается с первого элемента, и у которого сумма больше, чем  $\text{generated}$ .

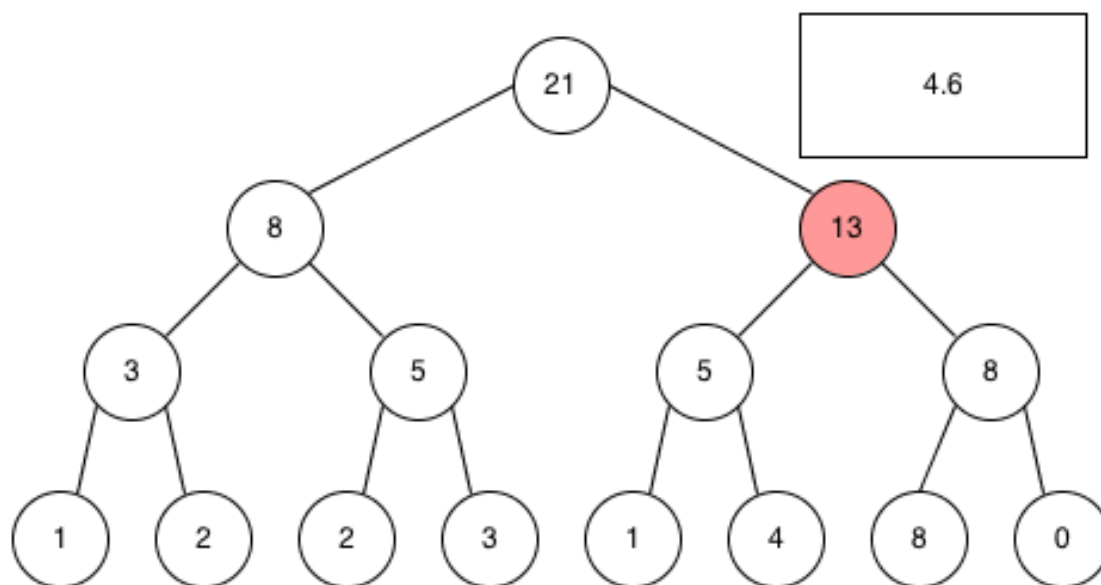
Помечаем корень базовой вершиной.



Смотрим в его левого сына.

$8 < 12$ .

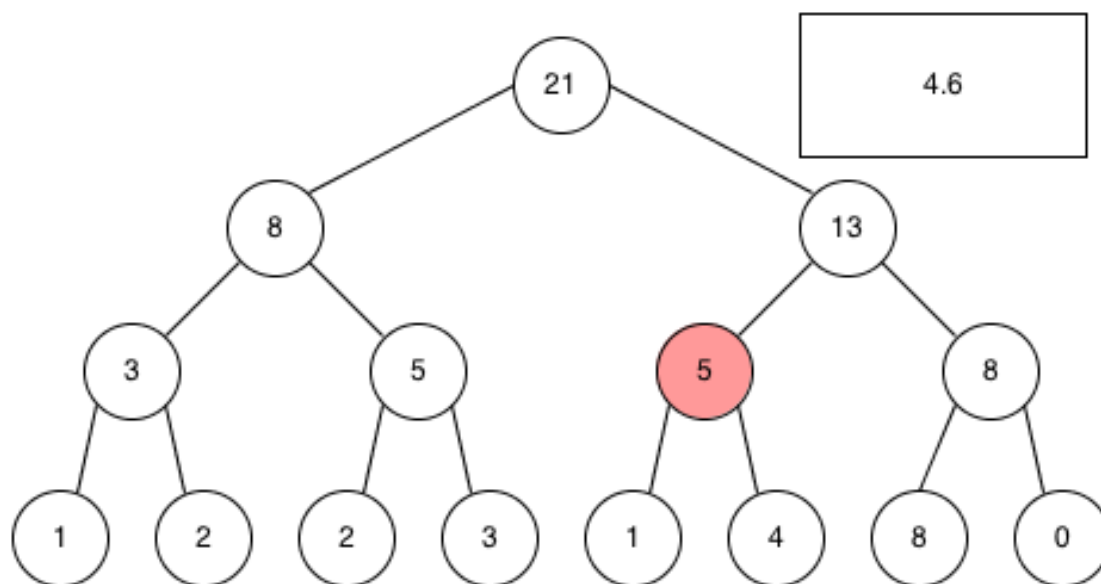
Вычитаем из generated значение в левом сыне и помечаем правого сына базовой вершиной.



Смотрим в левого сына.

$5 > 4$ .

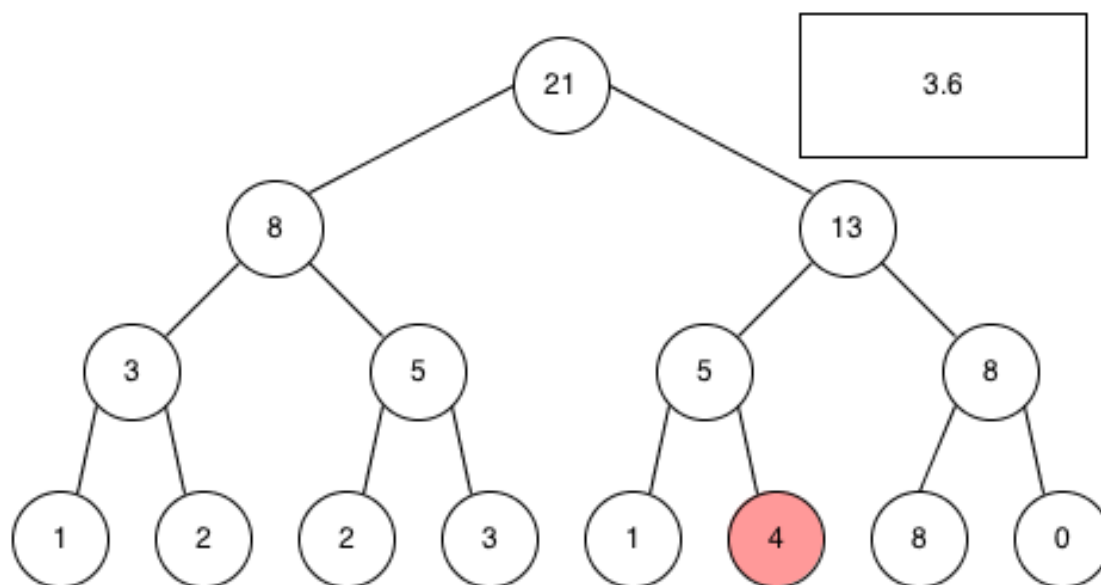
Поэтому помечаем левого сына базовой вершиной



Смотрим на левого сына.

$5 > 1$ .

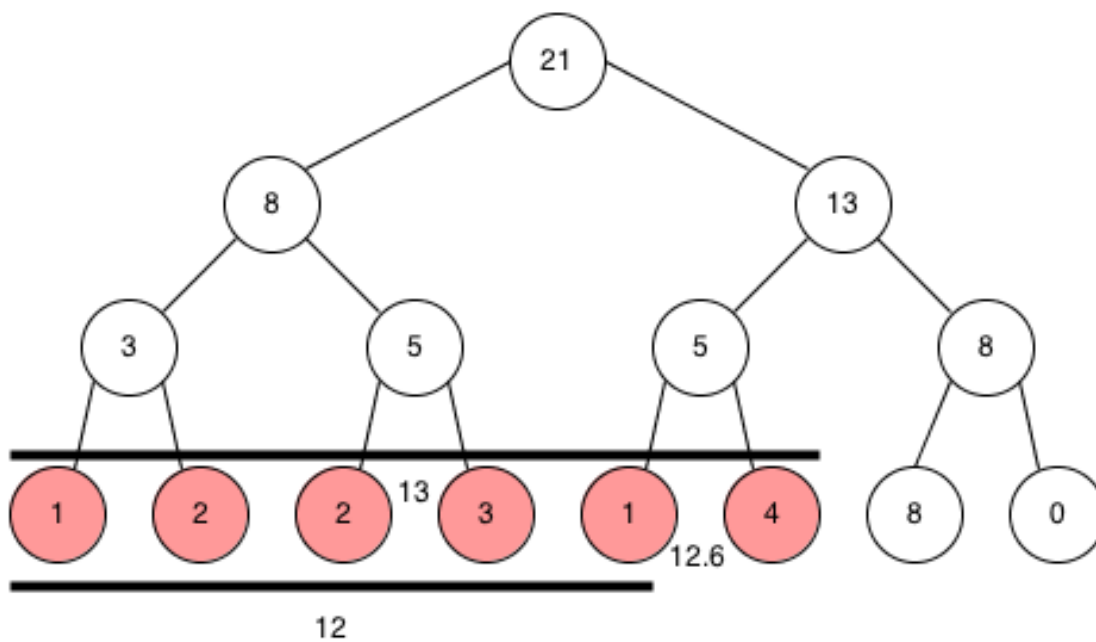
Поэтому вычитаем из generated 1 и помечаем базовой вершиной правого сына.



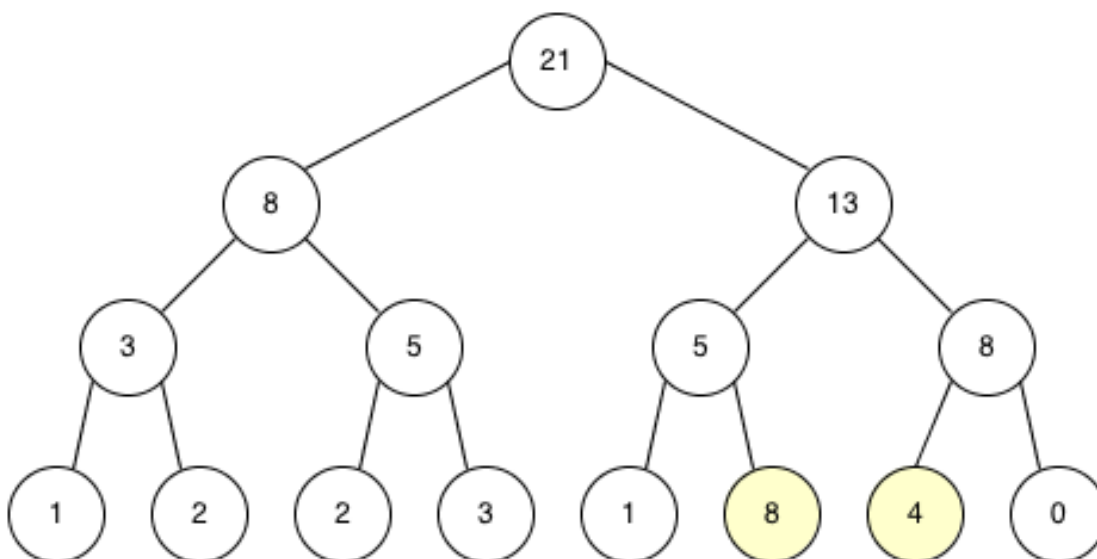
Наша базовая вершина стала листом.

Процесс завершается.

Индекс этой вершины говорит нам о том, что нам нужно показывать 6й баннер

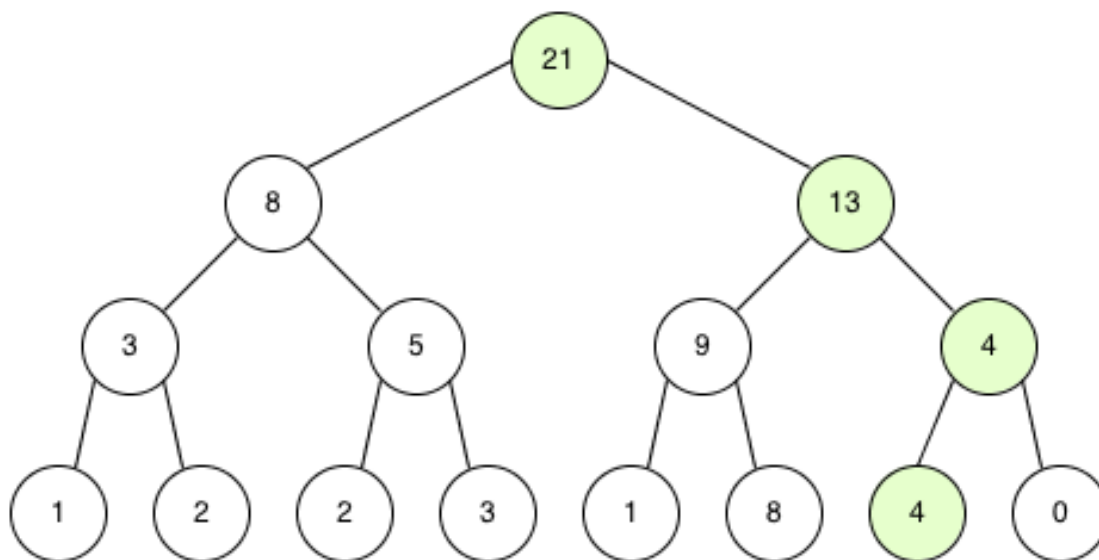
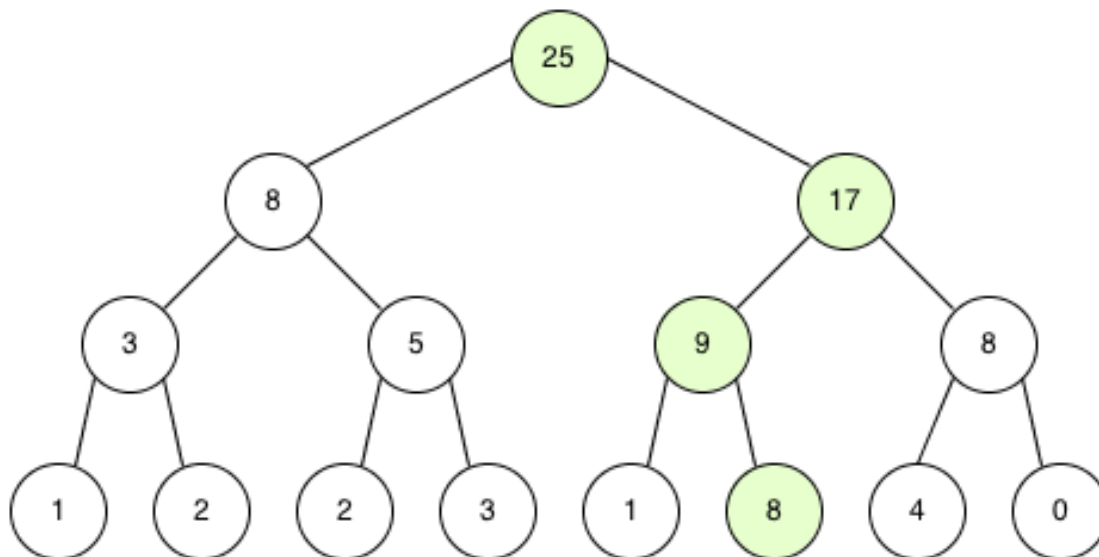


Теперь делаем swap найденного баннера с последним из имеющихся, чтобы при выборе следующего баннера его можно было не рассматривать.



Теперь нужно обновить дерево.

Делаем это для каждой из этих двух вершин.



Теперь мы получили обновленное дерево и можем продолжать выбирать следующий баннер для показа из набора еще не рассмотренных. Из суммарного количества весов для всех баннеров вычитается количество рассмотренных весов:  $21 - 4 = 17$ . Генерируется снова число от 0 до 1, пускай 0.5, умножается на новое суммарное количество весов:  $17 * 0.5 = 8.5$ , и в снова производится поиск по дереву. Процедура повторяется, пока не покажем все k баннеров.

---

## Алгоритм:

---

1. Строим дерево отрезков для подсчета сумм.
  - а. Добиваем до степени 2ки наш массив с весами.

- b. Создаем массив по размеру в два раз больший. Во вторую половину записываем наш массив
  - c. В элементы с индексом  $k$  записываем сумму элементов по индексам  $2k$  и  $2k+1$
- 2. Генерируем случайное число  $x$  из интервала  $[0,1)$
- 3. Умножаем его на сумму всех рассматриваемых весов
- 4. Начинаем искать интервал в дереве отрезков для сгенерированного значения
  - a. Помечаем корень базовой вершиной
  - b. Пока базовая вершина не является листом,
    - i. Если значение левого сына базовой вершины меньше, чем сгенерированное значение
      - 1. Помечаем левого сына базовой вершиной.
      - 2. Повторяем пункт b
    - ii. Иначе
      - 1. Из сгенерированного значения вычитаем значение левого сына
      - 2. Помечаем правого сына базовой вершиной
      - 3. Повторяем с пункта b
  - c. Определяем индекс базовой вершины, которая в данный момент уже находится в листе. Это и будет наш интервал.
- 5. Отображаем баннер, который находится по полученному индексу
- 6. Меняем баннеры по полученному индексу с последним
- 7. Обновляем в дереве суммы
- 8. Повторяем, начиная со второго пункта для нового набора баннеров без последнего элемента. И повторяем, пока не получим  $k$  баннеров

---

## Смещение вероятностей

---

Но хочется при этом отметить, что вероятность появления баннера на экране не будет соответствовать пропорциям.

Пусть у нас есть 4 баннера  $\{A:n_1, B:n_2, C:n_3, D:n_4\}$ , из них нам нужно получить 2 разных баннера.

$$\begin{aligned}
P(\text{показа } A) &= P(A \text{ будет показан в первом блоке}) \\
&+ \sum_{i \in \{B, C, D\}} P(i \text{ будет показан в первом блоке}) \\
&* P(A \text{ будет показан в 1 блоке выбирая из баннеров } \{A, B, C, D\} - \{i\})
\end{aligned}$$

Сравним соотношения вероятностей А и В.

Изначально, когда мы выбираем 1 баннер, соотношение было  $P(B) = cP(A)$ , то есть  $n_2 = cn_1$

Посчитаем, будет ли теперь такое же соотношение, когда мы выбираем два баннера

$$\begin{aligned}
P(A) &= \frac{n_1}{n} + \frac{n_2}{n} \left( \frac{n_1}{n - n_2} \right) + \frac{n_3}{n} \left( \frac{n_1}{n - n_3} \right) + \frac{n_4}{n} \left( \frac{n_1}{n - n_4} \right) \\
&= \frac{n_1}{n} + c \frac{n_1}{n} \left( \frac{n_1}{n - cn_1} \right) + \frac{n_3}{n} \left( \frac{n_1}{n - n_3} \right) + \frac{n_4}{n} \left( \frac{n_1}{n - n_4} \right) \\
P(B) &= \frac{n_2}{n} + \frac{n_1}{n} \left( \frac{n_2}{n - n_1} \right) + \frac{n_3}{n} \left( \frac{n_2}{n - n_3} \right) + \frac{n_4}{n} \left( \frac{n_2}{n - n_4} \right) \\
&= c \frac{n_1}{n} + c \frac{n_1}{n} \left( \frac{n_1}{n - n_1} \right) + c \frac{n_3}{n} \left( \frac{n_1}{n - n_3} \right) + c \frac{n_4}{n} \left( \frac{n_1}{n - n_4} \right)
\end{aligned}$$

Все дроби в  $P(A)$ , кроме  $c \frac{n_1}{n} \left( \frac{n_1}{n - n_1} \right)$ , очевидно больше в  $c$  раз чем соответствующие им в  $P(B)$ .

$$\frac{c \frac{n_1}{n} \left( \frac{n_1}{n - n_1} \right)}{c \frac{n_1}{n} \left( \frac{n_1}{n - cn_1} \right)} = \frac{n - cn_1}{n - n_1} = c - \frac{n(c - 1)}{n - n_1}$$

Значит при выборе нескольких баннеров соотношение уже выполняться не будет.

Для более общего случая рассмотрено доказательство не было. Но практическим экспериментом было подтверждено, что соотношение вероятностей быть показанными сдвигаются.