

ТАНКИ, УТКИ И ШАРЫ

Сергей Николенко

СПбГУ – Санкт-Петербург

13 ноября 2025 г.

Random facts:

- 13 ноября 1002 г., в день Святого Брайса, по приказу короля Англии Этельреда II Неразумного были убиты практически все даны, жившие в Англии; эффект оказался логичным, но вряд ли ожидаемым: Свен Вилобородый приплыл в Англию с огромным флотом и жаждой мести и к 1013 году полностью господствовал в Англии
- 13 ноября 1841 г. Джеймс Брейд впервые посетил демонстрацию «животного магнетизма» Чарльза Лафонтена; Брейд отнёсся скептически, но впоследствии нашёл рациональное зерно и стал первым современным исследователем гипноза
- 13 ноября 1862 г. Чарльз Доджсон начал выполнять обещание, данное Алисе Лидделл; из дневника: «Начал писать сказку об Алисе, надеюсь закончить её к Рождеству»
- 13 ноября 1887 г. — «кровавое воскресенье», правда, не в России; 10000 человек вышли в Лондоне на марш против происходящего в Ирландии и были разогнаны, 75 получили тяжёлые травмы; пехота была на месте, но приказа стрелять так и не получила
- 13 ноября 1947 г. официально завершилась разработка автомата Калашникова (АК-47); к этому типу принадлежит около 1/5 всего автоматического стрелкового оружия в мире



ОБУЧЕНИЕ РАВНОМЕРНОГО РАСПРЕДЕЛЕНИЯ

КАК ПОДСЧИТАТЬ НЕМЕЦКИЕ ТАНКИ

- Как во время второй мировой понять, сколько танков произвела Германия в этом году?
- Есть традиционные средства разведки: оценить число немецких танков там, где мы смогли их увидеть, и как-то экстраполировать на всё остальное
- Но нашёлся и другой неожиданный путь: предположим, что вы захватили немецкий танк (например, вывезли подбитый с поля боя); оказалось, что на разных деталях танка дотошные немецкие производители... наносят серийные номера!
- Мы можем посмотреть на, условно говоря, коленвал очередной «Пантеры» и узнать, что это коленвал «Пантеры» номер 273
- Что это нам скажет?

КАК ПОДСЧИТАТЬ НЕМЕЦКИЕ ТАНКИ

- Задача кажется несложной: даны несколько чисел, взятых из равномерного распределения на интервале $1, \dots, N$, и требуется оценить N
- Начнём с самого простого случая, когда число только одно; мы захватили танк с серийным номером m и хотим оценить, сколько всего их могло бы быть
- Функция правдоподобия:

$$p(m|N) = \begin{cases} \frac{1}{N}, & \text{если } N \geq m, \\ 0, & \text{если } N < m \end{cases}$$

- Гипотеза максимального правдоподобия — это $N = m$, тогда правдоподобие исхода m (и всех остальных) равно $\frac{1}{m}$
- Странный вывод, да? Это смещённая оценка, а нам бы несмешённую...

- Ответ из статистики: найти распределение максимума M из выборки, то есть оценить $\mathbb{E}[M | N]$, а затем подставить в качестве этого ожидания эмпирический максимум $m = \max_{i=1}^k x_i$ и решить уравнение относительно N
- Для $k = 1$ (одно число в выборке) это совсем просто: ожидание максимума совпадает с ожиданием одного x_i и равно $\frac{N+1}{2}$, а значит,

$$\hat{N} = 2m - 1$$

- Это звучит вполне интуитивно: мы ожидаем, что равномерно взятое случайное число в среднем попадёт в середину интервала, значит, мы ожидаем, что сам интервал будет примерно вдвое шире этого числа

КАК ПОДСЧИТАТЬ НЕМЕЦКИЕ ТАНКИ

- Для $k > 1$ придётся заняться комбинаторикой;
распределение $p(M|N)$ равно

$$p(M = m|N) = \frac{\text{число выборок из } k \text{ чисел с максимумом } m}{\text{общее число выборок из } k \text{ чисел}}$$

- В знаменателе число сочетаний из N по k , $\binom{N}{k}$, а в числителе, соответственно, число сочетаний из $m - 1$ по $k - 1$: мы уже зафиксировали одно из чисел m , и это максимум, то есть остальные $k - 1$ число должны быть взяты от 1 до $m - 1$
- Итого получается, что

$$p(M = m|N) = \frac{\binom{m-1}{k-1}}{\binom{N}{k}}.$$

КАК ПОДСЧИТАТЬ НЕМЕЦКИЕ ТАНКИ

- А ожидание максимума, значит, равно

$$\begin{aligned}\mathbb{E}[M \mid N] &= \sum_{m=k}^N mp(M=m|N) = \sum_{m=k}^N m \frac{\binom{m-1}{k-1}}{\binom{N}{k}} = \frac{1}{\binom{N}{k}} \sum_{m=k}^N m \binom{m-1}{k-1} = \\ &= \frac{1}{\binom{N}{k}} \sum_{m=k}^N m \frac{(m-1)!}{(k-1)!(m-k)!} = \frac{k}{\binom{N}{k}} \sum_{m=k}^N \frac{m!}{k!(m-k)!} = \frac{k}{\binom{N}{k}} \sum_{m=k}^N \binom{m}{k};\end{aligned}$$

в предпоследнем равенстве мы умножили и разделили на k

- Воспользуемся известным тождеством о биномиальных коэффициентах: $\sum_{a=c}^b \binom{a}{c} = \binom{b+1}{c+1}$; его легко доказать индукцией по b : для $b = c$ слева и справа единица, а если мы уже доказали тождество для некоторого b , то индукционный переход выглядит так:

$$\sum_{a=c}^{b+1} \binom{a}{c} = \sum_{a=c}^b \binom{a}{c} + \binom{b+1}{c} = \binom{b+1}{c+1} + \binom{b+1}{c} = \binom{b+2}{c+1};$$

последнее равенство — свойство треугольника Паскаля

- Применяя доказанное тождество, получим, что

$$\begin{aligned}\mathbb{E}[M \mid N] &= \frac{k}{\binom{N}{k}} \binom{N+1}{k+1} = \\ &= k \frac{k!(N-k)!}{N!} \frac{(N+1)!}{(k+1)!(N-k)!} = (N+1) \frac{k}{k+1}.\end{aligned}$$

- Получилась очень простая формула, которую легко теперь и обратить:

$$\hat{N} = m \left(1 + \frac{1}{k} \right) - 1 = m + \left(\frac{m}{k} - 1 \right).$$

- Вполне интуитивный результат: к максимуму нужно добавить «бонус», который уменьшается с ростом k

КАК ПОДСЧИТАТЬ НЕМЕЦКИЕ ТАНКИ

- А как по-байесовски? Мы хотим найти

$$p(N|M, k) = \frac{p(M|N, k) p(N|k)}{p(M|k)}.$$

- Мы уже подсчитали, что $p(M = m|N) = \binom{m-1}{k-1} / \binom{N}{k}$;
знаменатель будем считать как нормировочную константу

$$p(M|k) = \sum_{N=M}^{\infty} p(M|N, k) p(N|k).$$

- Теперь нужно выбрать априорное распределение $p(N|k)$:
какое взять? Мы бы хотели отразить в нём отсутствие
информации о значении N , но равномерное распределение,
как для параметра монетки, выбрать не получится, ведь нам
нужно распределение на всех натуральных числах

КАК ПОДСЧИТАТЬ НЕМЕЦКИЕ ТАНКИ

- Давайте возьмём равномерное априорное распределение $p(N|k)$ на большом объемлющем интервале k, \dots, Ω (меньше k номеров при k наблюдениях получиться не может):

$$p(N|k) = \frac{1}{\Omega - k + 1} \quad \text{для всех } k \leq N, N \leq \Omega.$$

- Соответственно, получается, что для $k \leq N$ и $N < \Omega$

$$\begin{aligned} p(N|M, k) &= \frac{p(M|N, k) p(N|k)}{\sum_{n=M}^{\Omega} p(M|n, k) p(n|k)} = \frac{\frac{1}{\Omega-k+1} p(M|N, k)}{\sum_{n=M}^{\Omega} \frac{1}{\Omega-k+1} p(M|n, k)} = \\ &= \frac{\binom{M-1}{k-1} / \binom{N}{k}}{\sum_{n=M}^{\Omega} \binom{M-1}{k-1} / \binom{n}{k}} = \frac{\binom{N}{k}^{-1}}{\sum_{n=M}^{\Omega} \binom{n}{k}^{-1}}, \end{aligned}$$

а для других значений N эта вероятность равна нулю.

- Как выбрать верхний предел Ω ? Можно просто заведомо большой, но технически препятствий к тому, чтобы положить $\Omega = \infty$, в формулах нет, ряд $\sum_{n=M}^{\infty} \binom{n}{k}^{-1}$ сходится для $k \geq 2$
- Это называется *некорректным априорным распределением* (improper prior): распределение «расходится», но после умножения на любое правдоподобие начинает сходиться, и байесовский вывод можно вести как обычно
- Например, бета-распределение для монетки без влияния на правдоподобие, то есть $B(0, 0)$:

$$p(\theta|0, 0) \propto \frac{1}{\theta} \frac{1}{1-\theta} \quad \text{для } \theta \in (0, 1).$$

КАК ПОДСЧИТАТЬ НЕМЕЦКИЕ ТАНКИ

- А для танков, применив некорректное априорное распределение с $\Omega = \infty$, приходим к апостериорному распределению

$$p(N|M, k) = \frac{\binom{N}{k}^{-1}}{\sum_{n=M}^{\infty} \binom{n}{k}^{-1}}.$$

- Ряд просуммировать трудно, давайте просто поверим, что

$$\sum_{n=M}^{\infty} \binom{n}{k}^{-1} = \binom{M}{k}^{-1} \frac{M}{k-1}, \quad \text{а значит,}$$

$$p(N|M, k) = \frac{k-1}{M} \binom{M}{k} \binom{N}{k}^{-1} \quad \text{для } N = M, M+1, \dots.$$

КАК ПОДСЧИТАТЬ НЕМЕЦКИЕ ТАНКИ

- Ответ (этот результат имеет смысл только для $k > 1$):

$$p(N|M, k) = \frac{k-1}{M} \binom{M}{k} \binom{N}{k}^{-1} \quad \text{для } N = M, M+1, \dots$$

- Получается *факториальное распределение* на $N - k$, его среднее и дисперсия известны:

$$\mathbb{E}[N | M, k] = \frac{k-1}{k-2} (M-1) \quad \text{для } k > 2,$$

$$\text{Var}[N | M, k] = \frac{(k-1)(M-1)(M-k+1)}{(k-2)^2(k-3)} \quad \text{для } k > 3.$$

КАК СЧИТАТЬ УТОК

- Сколько уток водится в России? Как мы можем это оценить?

КАК СЧИТАТЬ УТОК

- Сколько уток водится в России? Как мы можем это оценить?
- Метода Линкольна – Петерсена (Lincoln–Petersen method, 1896 и 1930):

$$\frac{\text{Число снова пойманных}}{\text{Размер второй выборки}} = \frac{\text{Число отмеченных}}{\text{Размер выборки}},$$

то есть

$$\text{Размер выборки} = \frac{\text{Размер второй выборки} \times \text{Число отмеченных}}{\text{Число снова пойманных}}.$$

- Кажется, тут трудно что-то добавить? Но какие здесь предположения?

- Популяция предполагается замкнутой, т.е. птицы не рождаются и не умирают во время эксперимента
- Предполагается, что обе выборки делаются *равномерно* (на практике это противоречит первому)
- Равномерность выборки означает ещё одно предположение: даже в уже «хорошо перемешанной» совокупности ещё надо предположить, что каждая особь попадётся в наши сети с одной и той же вероятностью
- Отметки предполагаются постоянными, то есть птица не может сбросить или уничтожить кольцо
- Немало предположений, да? Но есть и математическая проблема: метод Линкольна – Петерсена несколько раз использует оценку вероятности через отношение выборок $p \approx \frac{n}{N}$; это работает, если выборки большие, но для малых n и N работать перестаёт

- Давайте по-байесовски: правдоподобие в точности попадает в определение гипергеометрического распределения
- Нужно получить k успехов в выборке из n элементов без замещения из конечной популяции размера N с K успешными вариантами:

$$p(k|N, K, n) = \frac{\binom{K}{k} \binom{N-K}{n-k}}{\binom{N}{n}}.$$

- Какое выбрать априорное распределение?

- Давайте попробуем снова воспользоваться неинформативным априорным распределением и предположить, что апостериорное распределение $p(N|k)$ будет пропорционально правдоподобию:

$$p(N|k, K, n) = \frac{p(k|N)}{\sum_{M=k+n-k}^{\infty} p(k|M)} = \frac{\binom{K}{k} / \binom{N-K}{n-k} \binom{N}{n}}{\sum_{M=k+n-k}^{\infty} \binom{K}{k} \binom{M-K}{n-k} / \binom{M}{n}}$$

- Ряд в знаменателе сходится для $k \geq 2$, и моменты апостериорного распределения можно подсчитать через гипергеометрические функции (мы не будем):

$$\mathbb{E}[N | k, K, n] = \frac{(K-1)(n-1)}{(k-2)} \quad \text{для } k > 2,$$

$$\text{Var}[N|k, K, n] = \frac{(K-1)(n-1)(K-k+1)(n-k+1)}{(k-2)^2(k-3)} \quad \text{для } k > 3.$$

Спасибо!

Спасибо за внимание!

