

Практическое задание по теме «Байесовские рассуждения»

Солоткий Михаил, 417 группа ВМК МГУ

25 сентября 2018 г.

1 Вывод формул в модели посещаемости

1.1 Модель (1)

Рассмотрим вероятностную модель

$$P(a, b, c, d) = P(d|c) P(c|a, b) P(a) P(b),$$

$$d|c \sim c + \text{Bin}(c, p_3),$$

$$c|a, b \sim \text{Bin}(a, p_1) + \text{Bin}(b, p_2),$$

$$a \sim \text{Unif}[a_{\min}, a_{\max}],$$

$$b \sim \text{Unif}[b_{\min}, b_{\max}].$$

В варианте №2 задания предлагается вычислить следующие распределения: $P(a)$, $P(b)$, $P(c)$, $P(d)$, $P(b|a)$, $P(c|a)$, $P(c|b)$, $P(b|d)$, $P(b|a, d)$.

- $P(a) = \frac{1}{a_{\max} - a_{\min} + 1}$
- $P(b) = \frac{1}{b_{\max} - b_{\min} + 1}$
- $P(c) = \sum_{a,b} P(c|a, b) \cdot P(a, b) = \sum_{a,b} P(c|a, b) \cdot P(a) \cdot P(b)$
 $P(c = k|a, b) = \sum_{i=0}^k P(\text{Bin}(a, p_1) = i) \cdot P(\text{Bin}(b, p_2) = k - i)$
 $P(c = k) = \sum_{a=a_{\min}}^{a_{\max}} \sum_{b=b_{\min}}^{b_{\max}} \sum_{i=0}^k P(\text{Bin}(a, p_1) = i) \cdot P(\text{Bin}(b, p_2) = k - i) \cdot P(a) \cdot P(b)$
 $P(c = k) = \sum_{a=a_{\min}}^{a_{\max}} \sum_{b=b_{\min}}^{b_{\max}} \left[\sum_{i=0}^k \binom{a}{i} \cdot p_1^i \cdot (1 - p_1)^{a-i} \cdot \binom{b}{k-i} \cdot p_2^{k-i} \cdot (1 - p_2)^{b-k+i} \right] P(b) P(a)$
- $P(d) = \sum_c P(d|c) \cdot P(c)$
 $P(d = k|c) = P(\text{Bin}(c, p_3) = k - c)$
 $P(d = k) = \sum_{c=0}^{a_{\max} + b_{\max}} P(\text{Bin}(c, p_3) = k - c) \cdot P(c)$
- $P(b|a) = P(b)$
- $P(c|a) = \sum_b P(c|a, b) \cdot P(b|a) = \sum_b P(c|a, b) \cdot P(b)$
 $P(c = k|a) = \sum_{b=b_{\min}}^{b_{\max}} \left[\sum_{i=0}^k P(\text{Bin}(a, p_1) = i) \cdot P(\text{Bin}(b, p_2) = k - i) \right] \cdot P(b)$
- $P(c|b) = \sum_a P(c|a, b) \cdot P(a|b) = \sum_a P(c|a, b) \cdot P(a)$
 $P(c = k|b) = \sum_{a=a_{\min}}^{a_{\max}} \left[\sum_{i=0}^k P(\text{Bin}(a, p_1) = i) \cdot P(\text{Bin}(b, p_2) = k - i) \right] \cdot P(a)$
- $P(b|d) = \frac{P(d|b) \cdot P(b)}{P(d)} = \frac{\left[\sum_c P(d|b, c) \cdot P(c|b) \right] \cdot P(b)}{P(d)} = \frac{\left[\sum_c P(d|c) \cdot P(c|b) \right] \cdot P(b)}{P(d)}$

- $$P(b|a, d) = \frac{\sum_c P(d|c) \cdot P(c|a, b) \cdot P(a) \cdot P(b)}{\sum_{b,c} P(d|c) \cdot P(c|a, b) \cdot P(a) \cdot P(b)}$$

$$\sum_{b,c} P(d|c) \cdot P(c|a, b) \cdot P(a) \cdot P(b) = \sum_c P(d|c) \left[\sum_b P(c|a, b) \cdot P(a) \cdot P(b) \right] = \sum_c P(d|c) \left[\sum_b P(a, b, c) \right] = \sum_c P(d|c) P(a, c) =$$

$$= \sum_c P(d|c) \cdot P(c|a) \cdot P(a)$$

$$P(b|a, d) = \frac{\sum_c P(d|c) \cdot P(c|a, b) \cdot P(a) \cdot P(b)}{\sum_c P(d|c) \cdot P(c|a) \cdot P(a)} = \frac{\sum_c P(d|c) \cdot P(c|a, b) \cdot P(b)}{\sum_c P(d|c) \cdot P(c|a)}$$

1.2 Модель (2)

Теперь изменим немного модель: заменим распределение $c|a, b$ на $\text{Poiss}(ap_1 + bp_2)$. Тогда поменяются формулы для $P(c)$, $P(c|a)$, $P(c|b)$.

- $$P(c = k|a, b) = P(\text{Poiss}(ap_1 + bp_2) = k) = \frac{e^{-ap_1 - bp_2} \cdot (ap_1 + bp_2)^k}{k!}$$

$$P(c = k) = \sum_{a=a_{\min}}^{a_{\max}} \sum_{b=b_{\min}}^{b_{\max}} \frac{e^{-ap_1 - bp_2} \cdot (ap_1 + bp_2)^k}{k!} \cdot P(b) \cdot P(a)$$
- $$P(c|a) = \sum_b P(c|a, b) \cdot P(b|a) = \sum_b P(c|a, b) \cdot P(b)$$

$$P(c = k|a) = \sum_{b=b_{\min}}^{b_{\max}} P(\text{Poiss}(ap_1 + bp_2) = k) \cdot P(b)$$
- $$P(c|b) = \sum_a P(c|a, b) \cdot P(a|b) = \sum_a P(c|a, b) \cdot P(a)$$

$$P(c = k|b) = \sum_{a=a_{\min}}^{a_{\max}} P(\text{Poiss}(ap_1 + bp_2) = k) \cdot P(a)$$

2 Вывод математических ожиданий и дисперсий случайных величин

- $$E(a) = \sum_a [a \cdot P(a)] = \sum_{a=a_{\min}}^{a_{\max}} a \cdot \frac{1}{a_{\max} - a_{\min} + 1} = \frac{1}{a_{\max} - a_{\min} + 1} \cdot \sum_{a=a_{\min}}^{a_{\max}} a =$$

$$= \frac{1}{a_{\max} - a_{\min} + 1} \cdot \left(\sum_{a=1}^{a_{\max}} a - \sum_{a=1}^{a_{\min}-1} a \right) = \frac{1}{a_{\max} - a_{\min} + 1} \cdot \left(\frac{a_{\max} \cdot (a_{\max} + 1)}{2} - \frac{(a_{\min} - 1) \cdot a_{\min}}{2} \right) =$$

$$= \frac{1}{a_{\max} - a_{\min} + 1} \cdot \left(\frac{a_{\max}^2}{2} + \frac{a_{\max}}{2} - \frac{a_{\min}^2}{2} + \frac{a_{\min}}{2} \right) = \frac{1}{a_{\max} - a_{\min} + 1} \cdot$$

$$\cdot \left(\frac{(a_{\max} - a_{\min}) \cdot (a_{\max} + a_{\min})}{2} + \frac{a_{\max} + a_{\min}}{2} \right) = \frac{a_{\max} + a_{\min}}{2 \cdot (a_{\max} - a_{\min} + 1)} \cdot (a_{\max} - a_{\min} + 1) =$$

$$\frac{a_{\max} + a_{\min}}{2}$$

$$D(a) = E(a^2) - E(a)^2 = \frac{1}{a_{\max} - a_{\min} + 1} \cdot \sum_{a=a_{\min}}^{a_{\max}} a^2 - \frac{1}{(a_{\max} - a_{\min} + 1)^2} \cdot \left(\sum_{a=a_{\min}}^{a_{\max}} a \right)^2 =$$

$$= \frac{(a_{\max} - a_{\min} + 1)^2 - 1}{12}$$
- $$E(b) = \frac{b_{\max} + b_{\min}}{2}$$

$$D(b) = \frac{(b_{\max} - b_{\min} + 1)^2 - 1}{12}$$
- $$c_{\max} = a_{\max} + b_{\max}$$

$$E(c) = \sum_{c=0}^{c_{\max}} c \cdot P(c)$$

$$D(c) = E(c^2) - E(c)^2 = \sum_{c=0}^{c_{\max}} [c^2 \cdot P(c)] - \left(\sum_{c=0}^{c_{\max}} c \cdot P(c) \right)^2$$

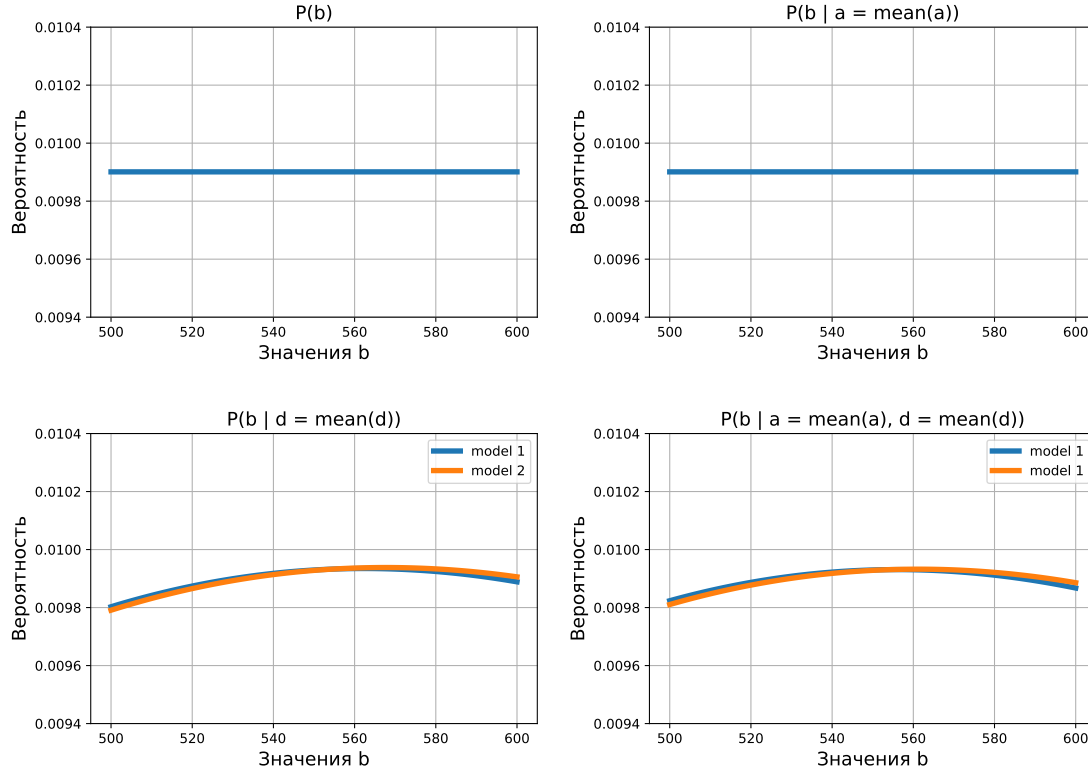
- $d_{max} = 2 \cdot c_{max}$

$$E(d) = \sum_{d=0}^{d_{max}} d \cdot P(d) = \sum_{d=0}^{d_{max}} \sum_{c=0}^{c_{max}} d \cdot P(\text{Bin}(c, p_3) = d - c) \cdot P(c)$$

$$D(d) = E(d^2) - E(d)^2 = \sum_{d=0}^{d_{max}} [d^2 \cdot P(d)] - \left(\sum_{d=0}^{d_{max}} d \cdot P(d) \right)^2$$

3 Влияние косвенной информации на прогноз величины b

Ниже приведены графики распределений величины b .



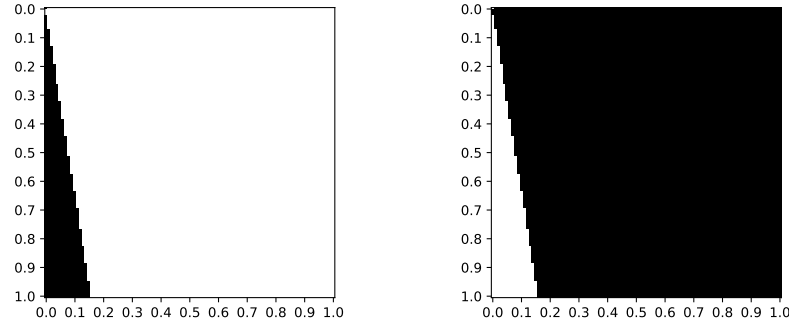
Распределения $P(b)$ и $P(b|a)$ не отличаются, так как величины a и b независимы. Распределения $P(b | d = E(d))$ и $P(b | a = E(a), d = E(d))$ отличаются не сильно, как можно видеть на графиках. Различия в вероятностях при одинаковых значениях b между ними порядка 10^{-5} . Как видно, само по себе наличие косвенной информации (зависимой) чуть-чуть меняет вид распределения, то есть более ярко выражаются наиболее правдоподобные значения b . Уточнение происходит относительно небольшое, такой же вывод можно сделать по значениям дисперсии – она не сильно уменьшилась. Ещё можно заметить, что для каждой отдельной модели добавление величины a в качестве условия уменьшает дисперсию и меняет мат. ожидание.

Величина	Модель	Мат. ожидание	Дисперсия
b	1, 2	550	850
$b d = E(d)$	1	550.07	848.037
$b d = E(d)$	2	550.09	848.128
$b a = E(a), d = E(d)$	1	550.03	848.031
$b a = E(a), d = E(d)$	2	550.06	848.123

4 Влияние параметров вероятностей на относительную важность параметров a и b для оценки величины c

Проведём вычисление значений $D(c | b = E(b))$ и $D(c | a = E(a))$ при различных значениях параметров p_1, p_2 вероятностей биномиальных распределений в модели (1). На равномерной сетке квадрата $[0, 1] \times [0, 1]$ (по 100 значений для каждой размерности). Ниже приведены графики соотношений дисперсий, белым обозначена область, где условие выполняется. По графикам кажется, что множества линейно разделимы.

$$\text{var}(c | b = \text{mean}(b)) < \text{var}(c | a = \text{mean}(a)) \quad \text{var}(c | b = \text{mean}(b)) \geq \text{var}(c | a = \text{mean}(a))$$



5 Доказательство линейной разделимости

Формула полной дисперсии:

$$D(Y) = E(D(Y|X)) + D(E(Y|X))$$

Возьмём $Y = c|b$, $X = a|b$

$$D(Y|X) = D(c|a, b) = D(\text{Bin}(a, p_1) + \text{Bin}(b, p_2)) = a \cdot p_1 \cdot (1 - p_1) + b \cdot p_2 \cdot (1 - p_2)$$

$$E(Y|X) = E(c|a, b) = E(\text{Bin}(a, p_1) + \text{Bin}(b, p_2)) = a \cdot p_1 + b \cdot p_2$$

$$E(D(Y|X)) = E(a) \cdot p_1 \cdot (1 - p_1) + E(b) \cdot p_2 \cdot (1 - p_2)$$

$$D(E(Y|X)) = D(a) \cdot p_1$$

$$D(c|b = E(b)) = E(a) \cdot p_1 \cdot (1 - p_1) + E(b) \cdot p_2 \cdot (1 - p_2) + D(a) \cdot p_1$$

В силу симметрии:

$$D(c|a = E(a)) = E(a) \cdot p_1 \cdot (1 - p_1) + E(b) \cdot p_2 \cdot (1 - p_2) + D(b) \cdot p_2$$

Пусть $c_1 = D(a)$, $c_2 = D(b)$, заметим, что $c_1 > 0$, $c_2 > 0$.

$$f(p_1, p_2) = D(c|b = E(b)) - D(c|a = E(a)) = c_1 \cdot p_1 - c_2 \cdot p_2$$

Если приравнять к 0, получим уравнение прямой, которая разбивает плоскость на 2 полуплоскости, причём в одной $f > 0$, а в другой $f < 0$. Заметим, что на квадрате $[0, 1] \times [0, 1]$ тоже есть точки, где функция f разного знака.

6 Замеры времени

Ниже приведены замеры времени вычисления распределений. Для величин, стоящих в условиях вычислялись вероятности для всех значений от минимального до максимального, то есть, к примеру $0 \leq d \leq 2 \cdot (a_{\max} + b_{\max})$. Замеры проводились с помощью утилиты **timeit** ipython notebook, которая запускает код 10 раз и оценивает по ним стандартное отклонение и среднее затраченного времени.

Распределение	c	$c a$	$c b$	$b a$	$b d$	$b a, d$	d
Время вычисления 1 модели: t_1	99.2 ms	111 ms	108 ms	2.48 μs	233 ms	2.42 s	165 ms
Стандартное отклонение t_1	1.46 ms	2.74 ms	1.35 ms	114 ns	13.9 ms	198 ms	2.18 ms
Время вычисления 2 модели: t_2	59.7 ms	58.6 ms	63.8 ms	2.63 μ s	172 ms	2.19 s	115 ms
Стандартное отклонение t_2	1.55 ms	345 μ s	2.31 ms	72.2 ns	4.29 ms	41 ms	2.17 ms

Можно заметить, что распределение $b | a, d$ вычисляется сильно дольше, чем остальные, но это и понятно, ведь оно единственное выводилось по определению условной вероятности и разложению совместного распределения. В процессе вычисления появлялись трёхмерные тензоры, в которых надо поменять некоторые размерности перед умножением. Быстрее всего вычислялось $b | a$, так как величины a и b независимы, а b распределено равномерно, то есть достаточно было заполнить матрицу одинаковым значением. Почти все распределения второй модели вычисляются быстрее соответствующих распределений первой модели. В одном случае получилось наоборот, но это не стат. значимый результат.

7 Сравнение моделей

Как видно из графиков распределений величины $b | a$ и величины $b | a, d$ для разных моделей разница в вероятностях незначительная, это значит, что можно не считать биномиальные распределения

и приближать их распределениями Пуассона (в случае конкретных использовавшихся в экспериментах параметрах вероятностей $p_1 = 0.1$ и $p_2 = 0.01$). Исходя из проведённых экспериментов нельзя сказать, что какая-то из моделей существенно лучше предсказывает какие-либо величины. Исходя из на практике использовать модель (2) предпочтительней (не факт, что тоже самое можно сказать в случае других значений параметров вероятностей посещения).