

# Глава 7

## Регулярные модели

### 7.1 Базовая часть

#### 7.1.1 Информация Фишера

##### Условия регулярности

Пусть имеется выборка размера 1 из бернуллиевского распределения. Насколько хорошо я могу оценить вероятность успеха? Очевидно, что не могу сказать ничего определенного, одно наблюдение может помочь мне составить только самое общее представление о параметре. При этом одно наблюдение из  $R[0, \theta]$ , скажем, будет несколько более информативным — я заведомо смогу вычеркнуть из потенциальных значений  $\theta$  отрезок  $[0, X_1]$ .

В рамках сегодняшнего материала нам хотелось бы отказаться от второго типа оценок, которые позволяют отдельным наблюдениям сужать область рассматриваемого параметра. Для этого, нам хотелось бы некоторую структуру пространства  $\Theta$  и плотностей распределения  $\theta$ . Наложим следующие условия регулярности:

R1) Пусть  $\Theta$  — открытый интервал прямой (возможно бесконечный).

R2) Распределения  $F_\theta$  имеют плотности  $f_\theta(x)$  и носители этих плотностей (распределений) не зависят от  $\theta$ .

R3) Плотность  $f_\theta(x)$  имеет конечную производную по  $\theta$  в каждой точке  $\Theta$  при каждом  $x$ , для которого  $f_\theta(x) > 0$ .

R4) Для величины  $U_1(X_1) = \frac{\partial}{\partial \theta} \ln f_\theta(X_1)$   $\mathbf{E}U_1(X_1) = 0$ ,  $0 < I_1(\theta) = \mathbf{D}U_1(X_1) < \infty$ .

Условие R1) необязательно, мы можем просто исключить из рассмотрения  $\theta$ , лежащие на границе. Условие R2) запрещает нам полностью отвергать значения  $\theta$  на основе наблюдений. Условия R3) и R4) говорят о том, что зависимость плотности от  $\theta$  достаточно гладка.

Аналогичные условия наложим в дискретном случае, рассматривая при этом вместо  $f_\theta(x)$  функцию  $\mathbf{P}_\theta(X = x)$ .

**Определение 1.** Величина

$$U(x_1, \dots, x_n; \theta) = \sum_{i=1}^n \frac{\partial}{\partial \theta} \ln f_\theta(x_i) = U_1 + \dots + U_n$$

называется функцией вклада выборки.

Если модель регулярна, то

$$\mathbf{E}U(X_1, \dots, X_n; \theta) = 0$$

в силу условия R4). Отметим, что это не очень сильное условие,

$$\mathbf{E}_\theta U(X_1; \theta) = \int_{\mathbb{R}} \frac{\partial}{\partial \theta} \ln f_\theta(x) f_\theta dx = \int_{\mathbb{R}} \frac{\partial}{\partial \theta} f_\theta dx = \frac{\partial}{\partial \theta} \int_{\mathbb{R}} f_\theta dx = 0.$$

По сути требуется лишь уметь переставлять дифференцирование с интегрированием.

## 7.1.2 Информация Фишера

Нам потребуется количественная характеристика информации о параметре, содержащейся в выборке.

**Определение 2.** Назовем информацией Фишера одного элемента  $X$  следующую величину:

$$I(\theta) = I_1(\theta) = \mathbf{E}_\theta \left( \frac{\partial}{\partial \theta} \ln f_\theta(X) \right)^2 = \mathbf{D}_\theta \frac{\partial}{\partial \theta} \ln f_\theta(X).$$

Иначе говоря, функция плотности одной случайной величины  $f_\theta(x)$  логарифмируется и затем дифференцируется по  $\theta$ , в полученную функцию вместо аргумента  $x$  подставляется случайная величина  $X$  и у этой величины подсчитывается дисперсия.

**Определение 3.** Информацией Фишера выборки  $X_1, \dots, X_n$  называют функцию

$$I_n(\theta) = \mathbf{E}_\theta \left( \frac{\partial}{\partial \theta} \ln f_\theta(X_1, \dots, X_n) \right)^2.$$

Иначе говоря,

$$I_n(\theta) = \mathbf{D}_\theta U(X_1, \dots, X_n; \theta)$$

В нашем случае н.о.р. величин  $I_n(\theta) = nI(\theta)$ .

Отметим, что иногда удобнее использовать формулу

$$I_n(\theta) = -\mathbf{E} \left( \frac{\partial^2}{\partial \theta^2} \ln f_\theta(X_1, \dots, X_n) \right),$$

справедливую в силу соображений

$$\begin{aligned} \mathbf{E} \left( \frac{\partial^2}{\partial \theta^2} \ln f_\theta(X_1, \dots, X_n) \right) &= \int_{\mathbb{R}} \frac{\partial^2}{\partial \theta^2} f_\theta(x_1, \dots, x_n) dx_1 \dots dx_n - \int_{\mathbb{R}} \left( \frac{\frac{\partial}{\partial \theta} f_\theta(x_1, \dots, x_n)}{f_\theta(x_1, \dots, x_n)} \right)^2 f_\theta(x_1, \dots, x_n) dx_1 \dots dx_n = \\ &= \frac{\partial^2}{\partial \theta^2} \int_{\mathbb{R}} f_\theta(x_1, \dots, x_n) dx_1 \dots dx_n - I_n(\theta) = -I_n(\theta). \end{aligned}$$

Здесь мы воспользовались тем, что

$$(\ln f(x))'' = \frac{f''(x)}{f(x)} - \left( \frac{f'(x)}{f(x)} \right)^2.$$

**Пример 1.** Подсчитаем информацию Фишера для бернуллиевских величин с параметром  $\theta \in (0, 1)$ . Запишем  $\ln f_\theta$  и продифференцируем его:

$$\ln f_\theta(x) = \ln(\theta^x (1-\theta)^{1-x}) = (1-x) \ln(1-\theta) + x \ln \theta, \quad \frac{\partial}{\partial \theta} \ln f_\theta(x) = \frac{x}{\theta} - \frac{1-x}{1-\theta} = \frac{x-\theta}{\theta(1-\theta)},$$

откуда

$$I(\theta) = \mathbf{D}_\theta \frac{\partial}{\partial \theta} \ln f_\theta(X) = \frac{\mathbf{D}_\theta X}{(\theta(1-\theta))^2} = \frac{1}{\theta(1-\theta)}, \quad I_n(\theta) = \frac{n}{\theta(1-\theta)}.$$

Таким образом, информация о выборке минимальна при  $\theta$  близким к  $1/2$  (в этом случае при небольшом изменении параметра доля 0 и 1 в выборке меняется незначительно), а максимальна (бесконечна) при  $\theta$ , приближающихся к 0 или 1 (в этом случае даже небольшие изменения параметра могут оказаться значимыми — в выборке значительно изменится соотношение 0 и 1).

**Пример 2.** Подсчитаем информацию Фишера для нормальных величин  $\mathcal{N}(\theta, 1)$ .

$$\ln f_{\theta}(x) = -\ln \sqrt{2\pi} - \frac{(x - \theta)^2}{2}, \quad \frac{\partial}{\partial \theta} \ln f_{\theta}(x) = x - \theta.$$

Таким образом,

$$I(\theta) = \mathbf{D}_{\theta}(X - \theta) = 1, \quad I_n(\theta) = n.$$

Еще проще было бы искать информацию Фишера с помощью второй производной  $\ln f_{\theta}(x)$  по  $\theta$ , которая равна  $-1$ . Итак, информация Фишера для  $\mathcal{N}(\theta, 1)$  постоянна и равна 1, информация Фишера выборки равна  $n$ .

## Информация Фишера статистики

**Определение 4.** Информацией о параметре, содержащейся в статистике  $T$  назовем

$$I_T(\theta) = \mathbf{D}_{\theta} \frac{\partial}{\partial \theta} (\ln f_{T(X_1, \dots, X_n), \theta}(T(X_1, \dots, X_n))).$$

Можно доказать, что  $I_{T(X_1, \dots, X_n)} \leq I_{X_1, \dots, X_n}$ , то есть информация, содержащаяся в выборке не меньше, чем информация, содержащаяся в функции от нее, что вполне естественно. Кроме того равенство  $I_{T(X_1, \dots, X_n)}(\theta) = I_{X_1, \dots, X_n}(\theta)$  выполняется тогда и только тогда, когда  $T$  достаточна.

**Пример 3.** Найдем в модели предыдущего примера информацию, содержащуюся в  $\bar{X}$ . Распределение этой величины  $\mathcal{N}(\theta, 1/n)$ , откуда

$$\ln f_{\theta}(x) = -\ln \sqrt{2\pi/n} - \frac{n(x - \theta)^2}{2}, \quad \frac{\partial}{\partial \theta} \ln f_{\theta}(x) = n(x - \theta),$$

и

$$I_T(\theta) = n^2 \mathbf{D} \bar{X} = n = I_n(\theta).$$

## 7.1.3 Неравенство информации

### Одномерное неравенство информации

Перейдем к базовому результату, связанному с понятием информации Фишера: информационному неравенству:

**Теорема 1.** (Неравенство Рао-Крамера). Пусть выполнены условия регулярности и  $\hat{\theta}$  — оценка с математическим ожиданием  $g(\theta)$ . Тогда

$$\mathbf{D}_{\theta} \hat{\theta}(X_1, \dots, X_n) \geq \frac{(g'(\theta))^2}{I_n(\theta)}.$$

В частности, для несмещенных оценок имеем  $\mathbf{D}_{\theta} \hat{\theta} \geq 1/(nI(\theta))$ .

**Пример 4.** Рассмотрим оценку  $\bar{X}$  в схеме Бернулли. Ее дисперсия равна  $\theta(1 - \theta)/n$ . При этом в силу неравенства Рао-Крамера несмещенные оценки в этой модели не могут иметь дисперсию лучше чем  $\theta(1 - \theta)/n$  (см. пример 3). Значит  $\bar{X}$  — оптимальная.

Оптимальные оценки, чья дисперсия совпадает с нижней границей из неравенства Рао-Крамера называются эффективными. Эффективные оценки бывают только в определенных моделях и у определенных функций, но об этом чуть позже.

**Пример 5.** Оценка  $\hat{\theta} = e^{\bar{X}-1/(2n)}$  в модели  $\mathcal{N}(\theta, 1)$  имеет математическое ожидание

$$e^{\theta-1/(2n)} \mathbf{E} e^{\bar{X}-\theta} = e^{\theta-1/(2n)} \mathbf{E} e^{Z/\sqrt{n}} = \frac{e^{\theta-1/(2n)}}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{x/\sqrt{n}} e^{-x^2/2} dx = \frac{e^{\theta}}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-(x-1/\sqrt{n})^2/2} dx = e^{\theta},$$

где  $Z = \sqrt{n}(\bar{X} - \theta) \sim \mathcal{N}(0, 1)$ . Следовательно,  $\hat{\theta}$  является оптимальной для функции  $e^{\theta}$  как функция от полной достаточной статистики. Покажем, что она не является эффективной оценкой той же функции. Аналогично предыдущим оценкам  $\mathbf{E} e^{2\bar{X}} = e^{2\theta+2/n}$ , откуда

$$\mathbf{D}_{\theta} \hat{\theta} = \mathbf{E} e^{2\bar{X}-1/n} - e^{2\theta} = e^{2\theta+1/n} - e^{2\theta} = e^{2\theta}(e^{1/n} - 1)$$

При этом  $I(\theta) = 1$ ,  $g(\theta) = e^{\theta}$ , откуда нижняя граница в неравенстве Рао-Крамера есть  $e^{2\theta}/n$ . Асимптотически эти границы эквивалентны, но при фиксированном  $n$   $e^{1/n} - 1 > 1/n$ , откуда нижняя граница неравенства Рао-Крамера оказывается недостижимой.

Неравенство Рао-Крамера показывает наилучший порядок приближения любой оценкой параметра в регулярной модели. Действительно

$$\mathbf{E}_{\theta}(\hat{\theta} - \theta)^2 = \mathbf{D}_{\theta} \hat{\theta} + (\mathbf{E}_{\theta} \hat{\theta} - \theta)^2.$$

Эта величина по порядку не меньше  $1/n$ , откуда  $\hat{\theta} - \theta$  имеет порядок не меньше, чем  $n^{-1/2}$ . Вот почему, рассматривая асимптотически нормальные оценки нам было наиболее интересно смотреть случай  $\sigma_n(\theta) = \sigma(\theta)n^{-1/2}$ . При этом мы видели, что в нерегулярных моделях такого может и не быть, скажем  $(n+1)X_{(n)}/n$  в модели  $R[0, \theta]$  имеет дисперсию порядка  $1/n^2$ .

Из доказательства неравенства Рао-Крамера нетрудно вывести, что равенство там возможно только в так называемых экспоненциальных моделях, в которых

$$f_{\theta}(x) = \exp(-A(\theta)B(x) + C(\theta))D(x).$$

В них существует эффективная оценка для функции  $g(\theta) = C'(\theta)/A'(\theta)$  существует и имеет вид

$$T(x_1, \dots, x_n) = n^{-1} \sum_{i=1}^n B(x_i).$$

Нетрудно понять, что эффективные оценки будут также для линейных функций от таких статистик. Других возможных эффективно оцениваемых функций не будет. Таким образом, эффективные оценки существуют лишь для небольшого набора функций лишь в некоторых моделях. В остальных случаях неравенство Крамера не дает наилучших оценок.

Этот результат следует принять во внимание, но пользоваться им как доказанным в рамках занятий нельзя

**Пример 6.** Для нормального распределения  $\mathcal{N}(\theta, 1)$  плотность имеет вид  $f_{\theta}(x) = (\sqrt{2\pi})^{-1} e^{-(x-\theta)^2/2} = (\sqrt{2\pi})^{-1} e^{-x^2} e^{\theta x - \theta^2/2}$ , откуда  $A(\theta) = \theta$ ,  $B(x) = x$ ,  $C(\theta) = \theta^2/2$ ,  $D(x) = e^{-x^2}$ . Соответственно, оценка  $\bar{X}$  является эффективной оценкой  $\theta$ .

## 7.1.4 ОМП в регулярном случае

При этом асимптотически эффективные оценки (то есть асимптотически нормальные оценки с асимптотической дисперсией  $I(\theta)$ ) существуют в значительно большем количестве моделей. Оказывается, что

при определенных условиях такой оценкой является ОМП

**Теорема 2.** Пусть модель сильно регулярна, то есть выполнены условия R1)-R4) + условия

R5) Плотности  $f_\theta$  трижды дифференцируемы по  $\theta$ ;

R6) При каждом  $\theta_0$  существует функция  $h(x)$ , такая что при некотором  $\delta$  и всех  $\theta \in [\theta_0 - \delta, \theta_0 + \delta]$

$$\left| \frac{\partial^3}{\partial \theta^3} f_\theta(x) \right| \leq h(x),$$

причем  $\int_{\mathbb{R}} h(x) f_\theta(x) dx < \infty$ .

Тогда существует ОМП  $\hat{\theta}(X_1, \dots, X_n)$ , а точнее сказать, решение уравнения максимального правдоподобия

$$\sum_{i=1}^n \frac{\partial}{\partial \theta} \ln f_\theta(x_i) = 0,$$

такое что

$$\sqrt{n}(\hat{\theta}(X_1, \dots, X_n) - \theta) \xrightarrow{d} Z \sim \mathcal{N}\left(0, \frac{1}{I(\theta)}\right).$$

Таким образом, ОМП в сильно регулярных моделях асимптотически эффективна. Возникает вопрос — действительно ли это наименьшая возможная асимптотическая дисперсия? Вообще говоря, нет, ведь из сходимости распределений не следует сходимость их дисперсий. Примером такого случая служит задача 2.1.3, в которой для модели  $X_i \sim \mathcal{N}(\theta, 1)$  строится асимптотически нормальная оценка с асимптотической дисперсией 1 при  $\theta \neq 0$  и  $b < 1$  при  $\theta = 0$ , тогда как информация Фишера модели 1.

Однако, при дополнительном условии непрерывности асимптотической дисперсии, оказывается, что асимптотическая дисперсия не может быть меньше  $I(\theta)$ . Отсюда следует еще одно свойство ОМП в сильно регулярных моделях — их асимптотическая дисперсия наилучшая возможная среди непрерывных. Этим же свойством обладает и ОМС, которая в сильно регулярных моделях имеет ту же асимптотическую дисперсию.

### 7.1.5 Многомерный случай

Рассмотрим случай многомерного параметра  $\theta$ .

**Определение 5.** Матрица

$$I_{i,j}(\theta) = \mathbf{E}_\theta \left( \frac{\partial}{\partial \theta_i} \ln f_\theta(X_1) \frac{\partial}{\partial \theta_j} \ln f_\theta(X_1) \right),$$

называется информационной.

Аналогичным образом матрица  $I^{(n)}(\theta) = nI(\theta)$  является ковариационной матрицей функции вклада выборки

$$U(X_1, \dots, X_n; \theta) = (U(X_1, \dots, X_n; \theta)_j) = \left( \sum_{i=1}^n \frac{\partial}{\partial \theta_j} \ln L(x; \theta_1, \dots, \theta_k) \right).$$

Неравенство Рао-Крамера для несмещенных оценок  $\hat{\theta}$  при этом принимает такой вид:

$$\Sigma^2(\theta) \geq (I^{(n)})^{-1}(\theta),$$

где  $\Sigma^2$  — ковариационная матрица вектора  $\hat{\theta}$ , а под  $A \geq B$  для матриц мы понимаем следующее:  $A - B$  положительно определена. Иначе говоря, для любого вектора  $s = (s_1, \dots, s_m)$

$$\mathbf{D}_\theta \left( \sum_{i=1}^m s_i \hat{\theta}_i \right) \geq s^t (I^{(n)})^{-1}(\theta) s.$$

Таким образом, в многомерном случае мы получаем оценку снизу для дисперсий всех линейных комбинаций наших оценок.

Как и в одномерном случае мы можем искать  $I$  в форме

$$I_{i,j}(\theta) = -\mathbf{E}_\theta \frac{\partial^2}{\partial \theta_i \partial \theta_j} \ln f_\theta(X_1).$$

Аналогичным образом, информационной матрицей статистики  $T$  будем называть

$$I_{i,j}(\theta) = \mathbf{E}_\theta \frac{\partial}{\partial \theta_i} \ln f_{T,\theta}(T) \frac{\partial}{\partial \theta_j} \ln f_{T,\theta}(T),$$

где  $f_\theta(t)$  — плотность статистики  $T$ . Как и прежде достаточность равносильна тому, что информационная матрица статистики равна  $nI(\theta)$ .

## 7.2 Факультатив

### 7.2.1 Расстояние Кульбака-Лейблера

Давайте представим себе, что у нас есть два возможных распределения с плотностями  $f$  и  $g$  и есть некоторые априорные вероятности  $p, q$  того, что распределение имеет плотность  $f$  или  $g$  соответственно. Тогда после наблюдения реализации выборки  $x_1, \dots, x_n$  у нас появляются апостериорные вероятности

$$\tilde{p} = \frac{pf(x_1)\dots f(x_n)}{pf(x_1)\dots f(x_n) + qg(x_1)\dots g(x_n)}, \quad \tilde{q} = \frac{qg(x_1)\dots g(x_n)}{pf(x_1)\dots f(x_n) + qg(x_1)\dots g(x_n)}.$$

Априорно мы считали, что первая плотность в  $d = p/q$  раз вероятнее второй, апостериорно стали считать то же самое с отношением  $\tilde{d} = \tilde{p}/\tilde{q}$ . Соответственно  $\tilde{d}/d$  это некоторый показатель, отражающий количество информации, помогающей нам выбрать из  $f$  и  $g$ , содержащейся в выборке. Глядя на наши формулы, мы можем увидеть, что это отношение есть  $\frac{f(x_1)\dots f(x_n)}{g(x_1)\dots g(x_n)}$ . Давайте в качестве меры различимости мер  $f$  и  $g$  возьмем

$$\ln \tilde{d}/d = \sum_{i=1}^n \ln \frac{f(x_i)}{g(x_i)}.$$

Назовем *расстоянием Кульбака-Лейблера* между вероятностными распределениями  $F, G$  с плотностями  $f, g$  следующую величину:

$$I(F : G) = \int_{\mathbb{R}} \ln \left( \frac{f(x)}{g(x)} \right) f(x) dx,$$

в дискретном случае, когда  $F, G$  имеют вероятности  $f(x), g(x)$  для отдельных значений  $x$ ,

$$I(F : G) = \sum_x \ln \left( \frac{f(x)}{g(x)} \right) f(x).$$

Рассматриваемая величина есть  $\mathbf{E}_F \ln(f(X)/g(X))$ , где индекс  $F$  означает, что  $X$  имеет функцию распределения  $F$ . Это являет собой среднее значение нашего  $\ln \tilde{d}/d$  при условии, что распределение  $F$ .

Строго говоря, рассматриваемая величина не является метрикой, поскольку несимметрична, но неотрицательна и обращается в ноль лишь при  $F = G$ :

$$I(F : G) = -\mathbf{E}_F \ln \left( \frac{g(X)}{f(X)} \right) \geq -\ln \left( \mathbf{E}_F \frac{g(X)}{f(X)} \right) = -\ln \left( \int_{\mathbb{R}} g(x) dx \right) = 0.$$

Введенное расстояние по сути измеряет, насколько вероятны большие различия между  $F$  и  $G$ , то есть насколько легко различить какое было распределение, глядя на выборку.

**Пример 7.** Найдем расстояние между  $Bernoulli(p)$  и  $Bernoulli(1-p)$ . Тогда

$$I(p : 1-p) = \ln(p/(1-p))p + \ln((1-p)/p)(1-p) = (2p-1) \ln(p/(1-p)).$$

При отдалении  $p$  от  $1/2$  оно возрастает до бесконечности.

С нашей, сугубо статистической позиции, вполне разумно рассматривать расстояния между распределениями, соответствующими двум различным значениям параметра  $\theta$  на основе выборки размера  $n$ :

$$I(\theta_1 : \theta_2, n) = \int_{\mathbb{R}^n} \ln \left( \frac{f_{\theta_1}(x_1, \dots, x_n)}{f_{\theta_2}(x_1, \dots, x_n)} \right) f_{\theta_1}(x_1, \dots, x_n) dx_1 \dots dx_n$$

В привычном для нас случае н.о.р. наблюдений имеем

$$I(\theta_1 : \theta_2, n) = \mathbf{E}_{\theta_1} \left( \ln \left( \frac{f_{\theta_1}(X_1) \dots f_{\theta_1}(X_n)}{f_{\theta_2}(X_1) \dots f_{\theta_2}(X_n)} \right) \right) = nI(\theta_1 : \theta_2, 1),$$

т.е. расстояние между параметрами линейно растет с увеличением выборки.

Предположим, что параметр меняется непрерывно и  $f_\theta$  зависит от него достаточно гладко. Тогда посмотрим насколько хорошо параметр  $\theta$  отделяется от окрестных значений:

$$\begin{aligned} I(\theta : \theta + \Delta\theta, 1) &= -\mathbf{E}_\theta \ln \left( \frac{f_{\theta+\Delta\theta}(X)}{f_\theta(X)} \right) = -\mathbf{E}_\theta \ln \left( 1 + \frac{\Delta\theta f'_\theta(X) + \Delta\theta^2 f''_\theta(X)/2 + o(\Delta\theta^2)}{f_\theta(X)} \right) = \\ &= -\Delta\theta \mathbf{E}_\theta \frac{f'_\theta(X)}{f_\theta(X)} - \Delta\theta^2 \mathbf{E}_\theta \frac{f''_\theta(X)}{2f_\theta(X)} + \frac{1}{2} \Delta\theta^2 \mathbf{E}_\theta \left( \frac{f'_\theta(X)}{f_\theta(X)} \right)^2 + o(\Delta\theta^2). \end{aligned}$$

Если семейство  $f_\theta$  устроено достаточно хорошо, чтобы можно было менять порядок интегрирования и дифференцирования, то первые два слагаемых последней суммы нулевые, поскольку

$$\int_{\mathbb{R}} f'_\theta(x) dx = \left( \int_{\mathbb{R}} f_\theta(x) dx \right)' = 0, \quad \int_{\mathbb{R}} f''_\theta(x) dx = \left( \int_{\mathbb{R}} f_\theta(x) dx \right)'' = 0,$$

откуда

$$\frac{2I(\theta : \theta + \Delta\theta, 1)}{\Delta\theta^2} \rightarrow \mathbf{E}_\theta \left( \frac{\partial}{\partial\theta} \ln f_\theta(X) \right)^2, \quad \Delta\theta \rightarrow 0.$$

Правая часть характеризует, насколько хорошо параметр  $\theta$  хорошо отделяется на основе выборки из одного элемента от окрестных значений параметра. Будем называть ее *информацией Фишера*  $I(\theta) = I_1(\theta)$ . Аналогично

$$I_n(\theta) = \mathbf{E}_\theta \left( \frac{\partial}{\partial\theta} \ln f_\theta(X_1, \dots, X_n) \right)^2.$$

## 7.2.2 ОМП и расстояние Кульбака

Отметим, что ОМП минимизирует расстояние Кульбака в следующем смысле. Расстояние

$$I(F; F_\theta) = \mathbf{E}_F \ln f(X) - \mathbf{E}_F \ln f_\theta(X)$$

мы хотим минимизировать, подобрав наилучшее  $\theta$ . При этом мы максимизируем второе слагаемое, которое можно аппроксимировать обычным путем

$$\mathbf{E}_\theta \ln f_\theta(X) \approx \frac{1}{n} \sum_{i=1}^n \ln f_\theta(X_i).$$

Максимизация левой части подменяется максимизацией правой, которая и дает ОМП.

Если наша параметрическая модель окажется неверной (то есть  $X_i \sim F$ , где  $F$  не лежит в парамет-

рическом семействе), то ОМП будет подбирать в классе  $F_\theta$  ближайшее к  $F$  распределение в смысле расстояния Кульбака и в условиях регулярности с некоторыми дополнениями (конечность  $\mathbf{E}_F \ln f(X)$ , где  $f$  – настоящая плотность, равномерная мажорируемость  $\ln f_\theta(x)$  функцией  $x$  с конечным математическим ожиданием и единственность указанного ближайшего распределения).

**Пример 8.** Пусть выборка  $X_i$  взята из биномиального распределения  $F$  с параметрами  $n, p$ , а мы считаем ее пуассоновской с ф.р.  $F_\theta$ . Тогда

$$I(F : F_\theta) = \sum_{k=0}^n \left( \ln (C_n^k p^k (1-p)^{n-k}) - \ln \left( \frac{\lambda^k e^{-\lambda}}{k!} \right) \right) C_n^k p^k (1-p)^{n-k} = \\ \mathbf{E}_F \ln (C_n^X p^X (1-p)^{n-X}) - \mathbf{E}_F \ln X! - \ln(\lambda) \mathbf{E}_F X + \lambda \mathbf{E}_F 1.$$

Мы хотим найти самое близкое к биномиальному пуассоновское распределение, а значит мы должны максимизировать величину

$$\ln(\lambda) \mathbf{E}_F X - \lambda \mathbf{E}_F 1 = \ln \lambda np - \lambda,$$

где мы воспользовались тем, что математическое ожидание биномиальной величины равно  $np$ . Максимум этой величины по  $\lambda$  будет при

$$\lambda = np.$$

Следовательно, ближайшее к биномиальному  $n, p$  пуассоновское распределение – это  $Poiss(np)$ .

### 7.2.3 Критерий Бхаттачария

Неравенство Рао-Крамера дает возможность проверять оптимальность оценок функции  $g(\theta)$ . Оценка является эффективной (а значит оптимальной) в том и только том случае, если оценка представляет собой величину

$$T(\vec{x}) = g(\theta) + \left\langle a(\theta), \frac{\partial}{\partial \theta} U(\vec{x}; \theta) \right\rangle,$$

где  $\langle, \rangle$  – скалярное произведение,  $a(\theta)$  – некоторая векторная функция параметра  $\theta$ . Это критерий эффективности и достаточное условие оптимальности, однако, довольно редкое. Оказывается что можно обобщить этот результат:

**Теорема 3.** *Предположим, что статистика  $T(\vec{x})$  допускает представление*

$$T(\vec{x}) = g(\theta) + \frac{1}{L(\vec{x}; \theta)} \left( \sum_i a_i(\theta) \frac{\partial}{\partial \theta_i} L(\vec{x}; \theta) + \sum_{i,j} a_{i,j}(\theta) \frac{\partial^2}{\partial \theta_i \partial \theta_j} L(\vec{x}; \theta) + \dots + \sum_{i_1, \dots, i_s} a_{i_1, \dots, i_s}(\theta) \frac{\partial^s}{\partial \theta_{i_1} \dots \partial \theta_{i_s}} L(\vec{x}; \theta) \right)$$

*при некоторых функциях  $a_i(\theta)$ . Тогда  $T$  – оптимальная оценка.*

**Пример 9.** Пусть  $X_i$  имеют бернуллиевское распределение с параметром  $\theta$ . Будем искать оценку по критерию Бхаттачария при  $s = 2$ . Найдем

$$\frac{1}{L} \frac{\partial}{\partial \theta} L = \frac{\partial}{\partial \theta} \ln L = \frac{x_1 + \dots + x_n}{\theta} - \frac{n - x_1 - \dots - x_n}{1 - \theta} = \frac{x_1 + \dots + x_n - n\theta}{\theta(1 - \theta)}; \\ \frac{1}{L} \frac{\partial^2}{\partial \theta^2} L = \frac{\partial^2}{\partial \theta^2} \ln L + \left( \frac{\partial}{\partial \theta} \ln L \right)^2 = -\frac{x_1 + \dots + x_n}{\theta^2} - \frac{n - x_1 - \dots - x_n}{(1 - \theta)^2} + \\ \left( \frac{x_1 + \dots + x_n - n\theta}{\theta(1 - \theta)} \right)^2 = \frac{n(n+1)}{(1 - \theta)^2} - \frac{(\theta^2 + (1 - \theta)^2 + 2n\theta)(x_1 + \dots + x_n)}{\theta^2(1 - \theta)^2} + \frac{(x_1 + \dots + x_n)^2}{\theta^2(1 - \theta)^2}.$$



Следовательно,

$$\begin{aligned} \widehat{\theta}(x_1, \dots, x_n) = & \theta^2 + a_1(\theta) \cdot \frac{1}{L} \frac{\partial}{\partial \theta} L + a_{1,1}(\theta) \cdot \frac{1}{L} \frac{\partial^2}{\partial \theta^2} L = \frac{a_{1,1}(\theta)(x_1 + \dots + x_n)^2}{\theta^2(1 - \theta)^2} + \\ & \frac{(x_1 + \dots + x_n)(-a_{1,1}(\theta)(-\theta^2 + (1 - \theta)^2 + 2n\theta) + a_1(\theta)\theta(1 - \theta))}{\theta^2(1 - \theta)^2} + \theta^2 + \frac{a_{1,1}(\theta)n(n - 1) - n(1 - \theta)a_1(\theta)}{(1 - \theta)^2}. \end{aligned}$$

Чтобы получить оценку в левой части мы должны добиться того, чтобы все три коэффициента не зависели от  $\theta$ .

Из первого уравнения  $a_{1,1}(\theta) = c_1\theta^2(1 - \theta)^2$  для некоторой константы  $c_1$ . Из третьего уравнения для некоторого  $c_2$  выполнено

$$\theta^2(c_1n(n - 1) + 1) - n\frac{a_1(\theta)}{1 - \theta} = c_3.$$

Из второго уравнения, таким образом,

$$-c_1(-\theta^2 + (1 - \theta)^2 + 2n\theta) + \frac{\theta^2(c_1n(n - 1) + 1) - c_3}{n\theta} = c_2,$$

где  $c_2$  — некоторая константа. Приравнявая все коэффициенты при  $\theta$  в правой и левой части, получаем  $c_3 = 0$ ,

$$((n - 1)_1 + 1/(n) - 1(2n - 2) = 0,$$

откуда

$$c_1 = \frac{1}{n(n - 1)}, \quad c_2 = -\frac{1}{n(n - 1)}.$$

Следовательно,

$$\widehat{\theta}(X_1, \dots, X_n) = \frac{1}{n(n - 1)} \sum_{i=1}^n X_i \left( \sum_{i=1}^n X_i - 1 \right).$$

является оптимальной.