

## ЧАСТЬ ВТОРАЯ

# МАТЕМАТИЧЕСКАЯ СТАТИСТИКА

*“ . . . а спорили они о чем угодно, начиная с подлинности Священного писания и кончая вопросом, что на самом деле значит надпись “гарантировано” на банках с джемом.”*

*Мюриэл Спарк, Мисс Джин Броди в расцвете лет  
(Muriel Spark, The prime of Miss Jean Brodie)*

# О Г Л А В Л Е Н И Е

## *часть вторая*

### МАТЕМАТИЧЕСКАЯ СТАТИСТИКА

§1 Проблема статистического вывода	164
§2 Выборочные характеристики. Достаточные статистики	175
§3 Оценка параметров. Метод моментов	189
§4 Оценка параметров. Метод максимального правдоподобия	196
§5 Эффективность оценок	210
§6 Доверительные интервалы	216
§7 Статистическая проверка гипотез (критерии значимости)	228
§8 Равномерно наиболее мощные критерии	242
§9 Проверка модельных предположений. Критерии согласия	249
Литература	259

## §1. Проблема статистического вывода

### Лекция 1

Теория вероятностей создает базу для построения моделей реальных явлений, в основе которых лежат соотношения между частотами появления определенных событий. Располагая вероятностной моделью, мы можем рассчитать вероятности (относительные частоты) этих событий и тем самым оптимизировать свое поведение в условиях неопределенности. Математическая статистика строит модели индуктивного поведения в этих условиях на основе имеющихся вероятностных моделей. Основная проблема состоит в том, чтобы по наблюдениям элементарных исходов (обычно это – значения наблюдаемых случайных величин) дать метод выбора действий, при которых частота ошибок была бы наименьшей. Естественно, эта проблема сопряжена с решением сложных задач на экстремум, но даже в том случае, когда эти задачи не удастся решить, теория вероятностей дает метод для расчета средней величины потерь, которые мы будем нести, используя конкретное, выбранное нами правило индуктивного поведения. Таким образом, *математическая статистика есть теория принятия оптимальных решений, когда последствия от действий, предпринимаемых на основе этих решений, носят случайный характер*. Математическая статистика использует методы теории вероятностей для расчета частоты “неправильных” решений или, более общо, для величины средних потерь, которые неизбежно возникают в условиях случайности, как бы мы ни пытались оптимизировать свое поведение в этих условиях.

Приведем два примера, иллюстрирующих задачи математической статистики и, отчасти, методы их решения, с тем чтобы в последующем формализовать общую проблему статистического вывода.

*Пример 1.1. Определение общего содержания серы в дизельном топливе.* Мы снова обращаемся к примеру 7.2 из курса теории вероятностей, где речь шла о важной в экологическом отношении характеристике дизельного топлива – процентном содержании элементарной серы, которая при сжигании и последующем соединении с водой дает серную кислоту. Необходимость использования методов теории вероятностей при аттестации дизельного топлива по этой характеристике была вызвана значительными расхождениями между результатами  $x_1, \dots, x_n$  параллельных и независимых испытаний  $n$  проб из партии дизельного топлива. Если даже исключить ошибки эксперимента, связанные с неправильным определением

веса пробы и титрованием, то все равно разброс в параллельных испытаниях будет значительным в силу случайного характера процесса сжигания пробы топлива и выпадения части элементарной серы в золу. Но в таком случае возникает естественный вопрос, что же мы измеряем и что же это за характеристика дизельного топлива, которую мы назвали “общим содержанием серы”? В практике лабораторных испытаний обычно говорят о *среднем* значении этой характеристики, и дизельное топливо аттестуется величиной  $\bar{x} = n^{-1} \sum_1^n x_k$  – арифметическим средним результатов параллельных испытаний. Это и есть то “индуктивное поведение” статистика в условиях случайности, о котором мы говорили в начале лекции, и оправдание разумности такого поведения естественно искать в рамках закона больших чисел.

Действительно, в примере 7.2 мы интерпретировали результат  $x$  определения общего содержания серы в одной пробе как результат наблюдения случайной величины  $X$ , распределенной по нормальному закону со средним  $\mu$  и дисперсией  $\sigma^2$ , причем значение (неизвестное экспериментатору) параметра  $\mu$  являлось математическим выражением той, не совсем понятной для нас характеристики испытываемого топлива, которая называлась “общим содержанием серы”. В рамках этой вероятностной модели естественно трактовать результаты  $x_1, \dots, x_n$  параллельных испытаний  $n$  проб дизельного топлива как наблюдения  $n$  независимых копий  $X_1, \dots, X_n$  случайной величины  $X$ . Термин “копия” в данном случае употребляется для обозначения того факта, что каждая из наблюдаемых случайных величин имеет то же распределение, что и  $X$ . Таким образом, постулируется, что  $X_1, \dots, X_n$  независимы и одинаково распределены  $\mathcal{N}(\mu, \sigma^2)$ , так что в силу закона больших чисел при неограниченном возрастании объема испытаний  $n$

$$\bar{X} = \frac{1}{n} \sum_{k=1}^n X_k \xrightarrow{P} \mu.$$

Итак, закон больших чисел гарантирует нам, что при достаточно большом объеме испытаний мы будем близки к истинному значению исследуемой характеристики топлива. Однако на практике в заводских лабораториях обычно сжигаются всего две пробы топлива, и только в исключительных случаях при поверке приборов или тестировании лаборантов делается четыре испытания. Естественно, при  $n = 2$  говорить о законе “больших” чисел просто смешно, – следует искать некоторую количественную характеристику последствий от неточной аттестации партии дизель-

ного топлива. Легко понять, что в основу такой характеристики следует положить ошибку  $|\bar{X} - \mu|$  в оценке параметра  $\mu$ , но, к сожалению, значение  $\mu$  нам неизвестно, а  $\bar{X}$  есть случайная величина, что окончательно делает проблему прогноза ожидаемых ошибок при аттестации конкретной партии топлива неразрешимой. Здесь наблюдается та же ситуация, что и при попытке предсказать сторону монеты, которая выпадет при ее подбрасывании. Точный прогноз невозможен, но методы теории вероятностей позволяют нам рассчитать, как часто мы будем ошибаться в прогнозе при достаточно длительной игре в орлянку. Следовательно, мы должны решить задачу о вычислении вероятности того, что ошибка в оценке  $\mu$  будет слишком большой – превосходить некоторую предписанную величину  $\Delta$ . Эта вероятность  $P(|\bar{X} - \mu| > \Delta)$  обычно называется *риском* оценки  $\bar{X}$ , а вероятность  $P(|\bar{X} - \mu| \leq \Delta)$  противоположного события – *надежностью* этой оценки.

Таким образом, риск оценки указывает частоту тех партий дизельного топлива, в паспорте которых общее содержание серы указано с недопустимо большой ошибкой. Зная риск оценки, мы можем вычислить средние затраты на выплату рекламаций по искам потребителей дизельного топлива. Вывести формулу для вычисления риска не представляет особого труда, если обратиться к теореме сложения для нормального распределения (предложение 12.2 курса ТВ). Выборочное среднее  $\bar{X}$  есть нормированная на  $n$  сумма независимых одинаково распределенных  $\mathcal{N}(\mu, \sigma^2)$  случайных величин. В силу теоремы сложения эта сумма имеет также нормальное распределение, среднее значение которого равно сумме средних  $n\mu$ , а дисперсия равна сумме дисперсий  $n\sigma^2$ . При умножении на  $1/n$  среднее умножается на ту же величину, а дисперсия умножается на ее квадрат. Таким образом,  $\bar{X} \sim \mathcal{N}(\mu, \sigma^2/n)$ , надежность оценки

$$P(-\Delta \leq \bar{X} - \mu \leq \Delta) = \Phi(\Delta\sqrt{n}/\sigma) - \Phi(-\Delta\sqrt{n}/\sigma) = 2\Phi(\Delta\sqrt{n}/\sigma) - 1$$

(напомним,  $\Phi(-x) = 1 - \Phi(x)$ ), а ее риск

$$P(|\bar{X} - \mu| > \Delta) = 2(1 - \Phi(\Delta\sqrt{n}/\sigma)).$$

При вычислении риска оценки необходимо знать величину стандартного отклонения  $\sigma$ . Но значение  $\sigma$ , очевидно, остается постоянным при аттестации различных партий – это параметр, характеризующий точность метода химического анализа топлива, и не имеет отношения к его химическому составу. Естественно, за достаточно короткий срок в лабораториях накапливается большой архивный материал данных испытаний различных партий

топлива, что позволяет оценить значение  $\sigma$  с достаточно высокой точностью. С тем, как это делается, мы познакомимся в одной из ближайших лекций.

Используя формулу риска, мы можем определить минимальный объем испытаний  $n$ , гарантирующий предписанную, достаточно малую величину риска. Действительно, если  $\alpha$  – заданное ограничение на риск оценки, то разрешая неравенство  $2(\Phi(\Delta\sqrt{n}/\sigma) - 1) \leq \alpha$  относительно переменной  $n$ , получаем, что требуемый объем испытаний определяется неравенством

$$n \geq \left( \frac{\Phi^{-1}(1 - \alpha/2)\sigma}{\Delta} \right)^2.$$

**Пример 1.2. Выявление эффекта лечения.** Группа пациентов в количестве 10 человек, обладающих схожими антропометрическими и антропологическими характеристиками, подвергается лечению по некоторой новой методике, подтверждение или опровержение эффективности которой составляет предмет статистического исследования. После лечения дается только качественное заключение о состоянии здоровья каждого пациента, так что результат испытания новой методики можно представить в виде последовательности  $x_1, \dots, x_{10}$ , компоненты которой принимают значения 1 (положительный исход лечения) или 0 (отрицательный исход).

Предлагается следующее статистическое правило: новая методика объявляется эффективной, если  $x_i = 1$  для всех  $i = 1, \dots, 10$ , то есть все пациенты выздоровели. Если же лечение хотя бы одного пациента не привело к положительному исходу, новая методика не рекомендуется к дальнейшему клиническому использованию. Что можно сказать о надежности или, как говорят медики, “достоверности” такого правила индуктивного поведения?

Чтобы ответить на этот вопрос, мы должны построить вероятностную модель проводимых наблюдений. Естественно предполагать, что в силу “однородности” группы пациентов они обладают одинаковой вероятностью  $p$  положительного исхода лечения, и если в процессе лечения они не имели возможности излишне тесного общения, то исходы лечений можно представить в виде реализации десяти независимых бинарных случайных величин  $X_1, \dots, X_{10}$ , каждая из которых принимает значение 1 с вероятностью  $p$  и значение 0 с вероятностью  $1-p$ . Таким образом, мы пришли к модели испытаний в схеме Бернулли с вероятностью  $p$  успешного исхода. Вероятность того, что все 10 исходов были успешными равна  $p^{10}$ , и задавая различ-

ные значения  $p$  мы можем судить о том, как часто возможны различные результаты апробации нового метода лечения.

Предположим сначала, что новая методика неэффективна. При таком предположении значение  $p$  не должно превосходить величины  $1/2$ , и максимальное значение вероятности события  $X_1 = 1, \dots, X_{10} = 1$  равно  $2^{-10} = 1/1024 < 0,001$ . Это очень редкое событие, и поэтому предположение о неэффективности новой методики должно быть отвергнуто. При этом вероятность  $2^{-10}$  можно интерпретировать как риск внедрения в медицинскую практику неэффективного метода лечения: *используя предложенное правило выбора между двумя действиями (внедрение или отклонение методики) при испытаниях последующих методик, мы рискуем в среднем не более чем один раз из тысячи внедрить неэффективный метод лечения.*

Интересно заметить, что в предположении “нейтральности” нового метода ( $p = 1/2$ ) вероятность любого исхода  $X_1 = x_1, \dots, X_{10} = x_{10}$  одинакова и равна  $2^{-10}$ , но исход  $X_1 = 1, \dots, X_{10} = 1$  обладает наибольшей вероятностью принятия действительно эффективной методики, ибо

$$p \sum_1^{10} x_k (1-p)^{n-\sum_1^{10} x_k} \leq p^{10},$$

если  $p > 1/2$ . Столь же просто проверить, что результаты испытаний, в которых лечение только одного пациента окончилось неудачей, имеют вероятность  $p^9(1-p)$ , и такие 10 результатов  $x_1, \dots, x_{10}$  с одним  $x_i = 0$  и другими  $x_j = 1$  обладают большей вероятностью, чем исходы с двумя и более количеством неудач, если в действительности  $p > 1/2$ . Это замечание позволяет нам определить статистическое правило, обладающее наибольшей вероятностью принятия в действительности эффективной методики, но не с таким малым риском, как  $2^{-10}$ .

Дело в том, что в медицинской практике установилась определенная граница риска, равная 0.05, и все события, обладающие меньшей вероятностью, объявляются “редкими” – ими можно пренебречь. В связи с этим позволим себе включить в область принятия новой методики дополнительные исходы с ровно одним неуспехом, и вычислим риск такого статистического правила при  $p = 1/2$ . Используя известные нам формулы биномиальных вероятностей, находим, что

$$P \left( \sum_1^{10} X_k \geq 9 \right) = p^{10} + C_{10}^1 p^9 (1-p),$$

и при  $p = 1/2$  эта вероятность равна  $2^{-10}(1 + 10) = 11/1024 \approx 0,01$ , что по-прежнему достаточно мало по сравнению с 0.05. Следовательно, мы можем включить в область принятия новой методики еще  $C_{10}^2$  результатов испытаний, в которых присутствуют ровно две неудачи. Риск такого статистического правила становится равным

$$P\left(\sum_1^{10} X_k \geq 8\right) = p^{10} + C_{10}^1 p^9(1-p) + C_{10}^2 p^8(1-p)^2,$$

и при  $p = 1/2$  эта вероятность равна  $2^{-10}(1 + 10 + 45) = 56/1024 \approx 0.05$ .

Это как раз соответствует принятой в медицине норме риска статистического правила. Итак, мы рекомендуем новую методику к дальнейшему использованию в клинике, если лечение не более чем двух пациентов из десяти оказалось неудачным, и применение такого правила в испытаниях дальнейших методик может привести к принятию неэффективного метода лечения в среднем в пяти случаях из 100.

Мы рассмотрели две типичных задачи математической статистики – оценка параметров и проверка гипотез. Естественно, круг проблем математической статистики намного шире, но при надлежащей трактовке проблем большинство из них сводится или к задаче оценки параметров, или к задаче выбора одного из нескольких альтернативных высказываний об исследуемом объекте. Опираясь на рассмотренные примеры, мы можем теперь представить достаточно общую схему статистического вывода.



Любое статистическое исследование, проводимое в рамках математической статистики, начинается с описания объекта исследования и формализации *пространства*  $\mathcal{D}$  решений  $d$ , одно из которых статистик принимает на основе наблюдений независимых копий случайной, возможно векторной, величины  $X$ , характеризующей состояние объекта в момент проведения наблюдений. Так, в примере с аттестацией партии дизельного топлива (объект исследования)  $\mathcal{D}$  есть интервал  $(0; 100)$  (напомним, общее содержание серы измеряется в процентах к весу пробы), а в примере с определением эффективности нового метода лечения (объект исследования) пространство  $\mathcal{D}$  состоит из двух точек:  $d_0$  – решение о неэффективности метода (принятие “нулевой” гипотезы) и  $d_1$  – решение о внедрении нового метода в лечебную практику (принятие альтернативной гипотезы).

Наиболее важной и, по-видимому, наиболее сложной частью статистического исследования является этап построения *вероятностной модели*, который состоит в спецификации семейства  $\mathcal{P} = \{P_\theta, \theta \in \Theta\}$  возможных распределений наблюдаемой случайной величины  $X$ . Этот этап связан с достаточно глубоким проникновением в природу исследуемого объекта и метода наблюдений  $X$ , – одной математикой здесь, как правило, не обойтись. Семейство  $\mathcal{P}$  индексируется абстрактным *параметром*  $\theta$ , совокупность значений которого  $\Theta$  называется *параметрическим пространством*.

В первом примере мы выяснили, что семейство возможных распределений  $X$  есть семейство нормальных распределений  $\mathcal{N}(\mu, \sigma^2)$  с двумерным параметром  $\theta = (\mu, \sigma)$  и параметрическим пространством  $\Theta = \mathbb{R} \times \mathbb{R}^+$ . В дальнейшем мы предположили, что значение  $\sigma$  известно, и свели наше параметрическое пространство к евклидовой прямой:  $\Theta = \mathbb{R}$  с  $\theta = \mu$ . Наконец, поскольку общее содержание серы измеряется в процентах, мы должны окончательно положить  $\Theta = (0; 100)$ .

Во втором примере мы имели дело с бинарной случайной величиной  $X$ , принимающей значение 1 с вероятностью  $p$  и значение 0 с вероятностью  $1 - p$ . Таким образом, вероятностная модель представлялась семейством двухточечных распределений  $B(1, p)$  с  $\theta = p$  и параметрическим пространством  $\Theta = (0; 1)$ .

Следующий этап статистического исследования состоит в интерпретации решений  $d$  в терминах высказываний о соответствующих этому решению значениях параметра  $\theta$ . Это необходимо сделать, если мы поставили

себе задачу количественного измерения последствий от принятия неверных решений, – в наших примерах риск используемых правил представлял собой функцию от  $\theta$ . Нетрудно понять, что в первом примере  $\mathcal{D} = \Theta$ , а во втором примеру решению  $d_0$  о неэффективности метода соответствует подмножество параметрического пространства  $(0; 1/2]$ , а альтернативному решению  $d_1$  об использовании новой методики соответствует интервал  $(1/2; 1)$  возможных значений параметра  $\theta = p$ . Именно таким образом мы сводим конкретные задачи по аттестации партии дизельного топлива и выявлению эффективности нового метода лечения к абстрактным задачам математической статистики – оценке параметра (среднего значения)  $\theta$  нормального  $(\theta, \sigma^2)$  распределения и, соответственно, различению двух гипотез  $H_0 : \theta \in (0; 1/2]$  и  $H_1 : \theta \in (1/2; 1)$  о величине вероятности  $\theta$  успешного испытания в схеме Бернулли.

Параметрическая интерпретация решений позволяет статистику задать *потери*  $L(\theta, d)$ , которые он несет от принятия решения  $d$ , когда  $\theta$  представляет истинное значение параметра. Среднее значение этих потерь в длинном ряду однотипных статистических исследований с одним и тем же правилом принятия решения определяет величину риска, связанную с принятием неправильных решений. Так, в наших примерах риск определялся вероятностью принятия решения, отстоящего достаточно далеко от того решения, которое соответствовало истинному значению параметра, и, следовательно, *функция потерь* определялась индикатором некоторого подмножества в  $\Theta \times \mathcal{D}$ . Это так называемые функции потерь типа 0–1. В первом примере  $L(\theta, d) = 1$ , если  $|d - \theta| > \Delta$ , и  $L(\theta, d) = 0$  в противном случае. Во втором примере  $L(\theta, d) = 1$ , если принималось решение  $d_1$ , а  $\theta \in (0; 1/2]$ , или принималось  $d_0$ , а  $\theta \in (1/2; 1)$ , в остальных точках произведения пространств  $\Theta \times \mathcal{D}$  потери  $L(\theta, d)$  полагались равными нулю. Отметим, что в задаче оценки параметров довольно часто используется квадратичная функция потерь  $L(\theta, d) = |d - \theta|^2$ .

Каждое из решений  $d$  статистик принимает на основе результата  $x^{(n)} = x_1, \dots, x_n$  наблюдений над независимыми копиями  $X^{(n)} = (X_1, \dots, X_n)$  случайной величины  $X$ . Строится измеримое отображение  $\delta = \delta(\cdot)$  пространства возможных значений  $X^{(n)}$  в пространство решений  $\mathcal{D}$ , с помощью которого принимается решение  $d = \delta(x^{(n)})$ . Это отображение называется *решающей функцией* или *статистическим правилом*. Так, в первом примере

$\delta(X^{(n)}) = \bar{X}$ , а во втором

$$\delta(X^{(n)}) = \begin{cases} d_0, & \text{если } \sum_1^n X_k < 8, \\ d_1, & \text{если } \sum_1^n X_k \geq 8. \end{cases}$$

Последствия от использования конкретной решающей функции в длинном ряду однотипных статистических исследований определяются величиной средних потерь  $R(\theta; \delta) = \mathbf{E}_\theta L(\theta, \delta(X^{(n)}))$ , которая зависит от  $\theta$ ; функция  $R(\theta, \delta)$ ,  $\theta \in \Theta$ , называется *функцией риска*.

*Основная проблема математической статистики состоит в построении решающих функций  $\delta$ , минимизирующих равномерно по  $\theta \in \Theta$  функцию риска  $R(\theta; \delta)$ .* Мы будем решать эту проблему для задач оценки параметров и проверки гипотез. Естественно, будут также изучаться традиционные, возможно не обладающие оптимальными свойствами, статистические правила, и в этом случае нашей основной задачей будет вычисление их функций риска.

Представленная выше схема статистического вывода весьма далека от общности. Большинство статистических задач имеет дело с наблюдениями одновременно за несколькими объектами, например, новый метод лечения применяется к одной группе пациентов, в то время как другая подвергается лечению традиционным методом, и по данным наблюдений копий двух случайных величин делается вывод о предпочтительности нового метода. Если мы хотим сократить число наблюдений, необходимое для достижения заданной (малой) величины риска, то целесообразно не фиксировать заранее  $n$ , а планировать прекращение испытаний после наблюдения каждой копии в зависимости от полученных ранее результатов. Существует большой класс задач управления наблюдениями – оптимального выбора случайной величины, наблюдаемой на каждом шаге статистического эксперимента, а также правила прекращения наблюдений. Все это далеко выходит за рамки тех “кратких начатков” теории статистических выводов, которые будут представлены в нашем семестровом курсе.

Мы завершим этот параграф набором простейших определений и понятий, которые постоянно используются в математической статистике.

Итак, с исследуемым объектом, относительно которого мы должны принять некоторое решение  $d \in \mathcal{D}$ , соотносится наблюдаемая случайная величина  $X$ , распределение которой  $P_\theta$  известно с точностью до значения параметра  $\theta$ . Семейство распределений  $\mathcal{P} = \{P_\theta, \theta \in \Theta\}$ , как обычно, назы-

вается *вероятностной моделью*. Пусть  $(X, \mathcal{A})$  – измеримое пространство значений  $X$ . В дальнейшем будет всегда предполагаться, что на сигма-алгебре  $\mathcal{A}$  существует такая сигма-конечная мера  $\mu$ , что при любом  $\theta \in \Theta$  распределение  $X$  можно представить в виде интеграла

$$P_\theta(A) = \mathbf{P}(X \in A) = \int_A f(x | \theta) d\mu(x), \quad A \in \mathcal{A},$$

от плотности  $f(x | \theta)$  распределения  $X$  по мере  $\mu$ . В таком случае распределение независимых копий  $X^{(n)} = (X_1, \dots, X_n)$  случайной величины  $X$  на произведении  $(\mathcal{X}^n, \mathcal{A}^n)$  измеримых пространств  $(X, \mathcal{A})$  определяется функцией плотности

$$f_n(x^{(n)} | \theta) = \prod_{k=1}^n f(x_k | \theta)$$

по мере  $\mu_n = \underbrace{\mu \times \dots \times \mu}_n$ , то есть

$$P_{\theta,n}(A_n) = \mathbf{P}(X^{(n)} \in A_n) = \int_{A_n} f_n(x^{(n)} | \theta) d\mu_n(x^{(n)}), \quad A_n \in \mathcal{A}^n.$$

**Определение 1.1.** Вектор  $X^{(n)} = (X_1, \dots, X_n)$  независимых, одинаково распределенных по тому же закону, что и наблюдаемая случайная величина  $X$ , случайных величин называется *случайной выборкой объема  $n$* . Измеримое пространство  $(\mathcal{X}^n, \mathcal{A}^n)$  значений  $X^{(n)}$  называется *выборочным пространством*, а семейство распределений  $\mathcal{P}_n = \{P_{\theta,n}, \theta \in \Theta\}$  на этом пространстве – *статистической структурой* или *статистическим экспериментом*. Вектор  $x^{(n)} = (x_1, \dots, x_n)$  результатов наблюдения случайной выборки  $X^{(n)}$  называется *вектором (или совокупностью) выборочных данных*.

Зная распределение выборки, мы можем вычислять риск любого статистического правила  $\delta$  с помощью  $n$ -кратного интеграла

$$R(\theta; \delta) = \int_{\mathcal{X}} \dots \int_{\mathcal{X}} L(\theta, \delta(x^{(n)})) f_n(x^{(n)} | \theta) d\mu_n(x^{(n)}).$$

Конечно, если удастся найти распределение  $G_\theta$  решающей функции  $\delta$  на измеримом пространстве решений  $(\mathcal{D}, \mathcal{D})$ , то вычисление риска упрощается:

$$R(\theta; \delta) = \int_{\mathcal{D}} L(\theta, a) dG(a).$$

Так, в первом примере с выборкой из нормального  $(\mu, \sigma^2)$  распределения решающей функцией служило выборочное среднее  $\bar{X}$ . Было показано, что  $\bar{X}$  имеет нормальное  $(\mu, \sigma^2)$  распределение, и именно это обстоятельство позволило нам найти простое выражение риска статистического правила через функцию распределения стандартного нормального закона. Точно так же во втором примере с выбором из двухточечного распределения  $B(1, p)$  решающая функция была основана на случайной величине  $\sum_1^n X_k$ , которая имеет распределение Бернулли  $B(n, p)$ . Риск нашего решающего правила по выявлению эффективности метода лечения выражался через функцию распределения  $B(n, p)$ .

Заметим, что функции от выборочного вектора  $X^{(n)}$  играют важную, можно даже сказать самостоятельную, роль в математической статистике.

**Определение 1.2.** Любое измеримое отображение  $T = T(X^{(n)})$  выборочного пространства  $(\mathcal{X}^n, \mathcal{A}^n)$  в некоторое измеримое пространство  $(\mathcal{T}, \mathcal{B})$  называется *статистикой*.

Существует довольно устоявшийся универсальный набор статистик, постоянно используемых в теории и практике статистического вывода; распределения этих статистик интенсивно изучались на протяжении последних двух столетий. В следующем параграфе мы познакомимся с набором статистик, которые являются выборочными аналогами стандартных характеристик распределения наблюдаемой случайной величины, а также рассмотрим статистики, редуцирующие размерность выборочного вектора до размерности параметрического пространства без потери информации.

## §2. Выборочные характеристики. Достаточные статистики

Лекция 3

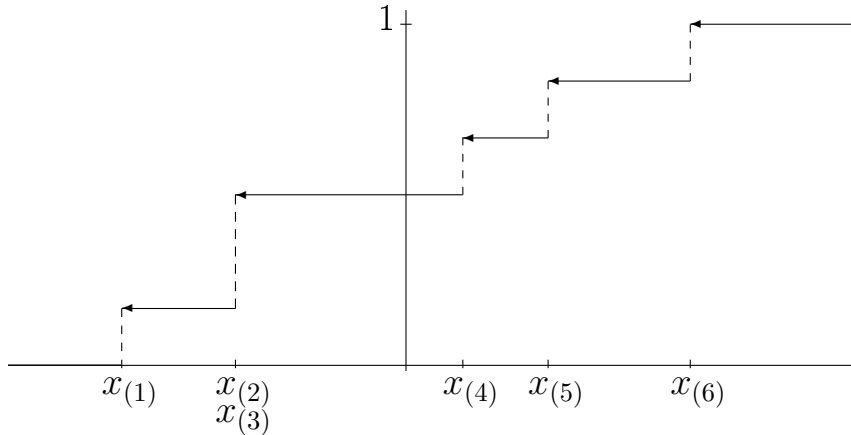
Построение вероятностных моделей в курсе теории вероятностей осуществлялось посредством спецификации функции распределения или функции плотности наблюдаемой случайной величины  $X$ . Любая из этих функций однозначно определяет распределение  $X$  на сигма-алгебре  $\mathcal{A}$  борелевских множеств, порожденной интервалами в пространстве  $\mathcal{X} = \mathbb{R}$  возможных значений  $X$ , и с их помощью вычислялись такие характеристики распределения, как среднее, дисперсия, коэффициенты асимметрии и эксцесса, квантили, мода и пр. В прикладной статистике существует традиция, или, можно сказать, обязательное правило, представлять полученные экспериментальные данные с помощью статистик – выборочных аналогов этих функций и характеристик распределения  $X$ . Выборочные характеристики являются оценками истинных значений своих прообразов и позволяют судить в общих чертах о характере распределения наблюдаемой случайной величины.

Такая “описательная” статистика обычно начинается с построения *вариационного ряда*: выборочные данные  $x_1, \dots, x_n$  упорядочиваются по возрастанию их значений  $x_{(1)} \leq \dots \leq x_{(n)}$ , и полученный таким образом вектор с неубывающими компонентами служит реализацией случайного вектора  $X_{(1)}, \dots, X_{(n)}$ , который, собственно, и следует называть *вариационным рядом*. Компоненты вариационного ряда называются *порядковыми статистиками*, а  $X_{(1)}$  и  $X_{(n)}$  – *крайними членами* вариационного ряда. Мы уже сталкивались с порядковыми статистиками, когда изучали структуру пуассоновского процесса и строили вероятностную модель “слабого звена” (распределение Вейбулла).

Упорядоченные данные наносятся на ось абсцисс, и строится ступенчатая функция, возрастающая скачками величины  $1/n$  в каждой точке  $x_{(1)}, \dots, x_{(n)}$ . Построенная таким образом дискретная функция распределения является реализацией случайной функции

$$F_n(x) = \frac{1}{n} \sum_{k=1}^n \mathbf{I}(X_k < x)$$

( $\mathbf{I}(A)$ , как обычно, индикатор события  $A$ ) и называется *эмпирической функцией распределения*.



Таким образом, дискретное эмпирическое распределение приписывает равные вероятности  $1/n$  каждой из  $n$  компонент выборочного вектора, и при каждом фиксированном  $x \in \mathbb{R}$  случайная величина  $nF_n(x)$  подчиняется биномиальному распределению  $B(n, F(x))$  :

$$P(F_n(x) = k/n) = C_n^k F^k(x)(1 - F(x))^{n-k}, \quad k = 0, 1, \dots, n.$$

В силу закона больших чисел Бернулли  $F_n \xrightarrow{P} F(x)$  при любом  $x \in \mathbb{R}$ . Более того, теорема Гливленко–Кантелли, утверждение которой

$$D_n = \sup_{x \in \mathbb{R}} |F_n(x) - F(x)| \xrightarrow{P} 0$$

мы приводим без доказательства, указывает на равномерность этой сходимости на всей числовой оси  $\mathbb{R}$ .

Мы закончим обсуждение свойств эмпирической функции распределения формулировкой широко известного результата А.Н. Колмогорова:

$$\lim_{n \rightarrow \infty} P(\sqrt{n}D_n < x) = \sum_{k=-\infty}^{+\infty} (-1)^k e^{-k^2 x^2}.$$

Полученная им формула для асимптотического ( $n \rightarrow \infty$ ) распределения статистики  $\sqrt{n}D_n$ , характеризующей величину расхождения между теоретическим  $F$  и эмпирическим  $F_n$  распределениями, используется для построения *критерия согласия* выборочных данных с предположением, что  $F$  является истинной функцией распределения, из которого извлекается выборка (гипотезой о том, что  $F$  есть функция распределения наблюдаемой случайной величины  $X$ ).

Итак, мы установили, что эмпирическое распределение сходится по вероятности к истинному (или, как обычно говорят прикладники, теоретическому) распределению, и теперь можем обратиться к вычислению моментных

и квантильных характеристик распределения  $F_n$ . Его нецентральные

$$a_k = \int_{\mathbb{R}} x^k dF_n(x) = \frac{1}{n} \sum_{i=1}^n X_i^k$$

и центральные

$$m_k = \int_{\mathbb{R}} (x - a_1)^k dF_n(x) = \frac{1}{n} \sum_{i=1}^n (X_i - a_1)^k$$

моменты служат выборочными аналогами соответствующих теоретических моментов  $\alpha_k$ ,  $k = 1, 2, \dots$ , и  $\mu_k$ ,  $k = 2, 3, \dots$ , и называются *выборочными моментами*.

Если теоретические моменты существуют, то в силу закона больших чисел выборочные моменты сходятся по вероятности к своим теоретическим прообразам. Среди выборочных моментов особое место занимают моменты первого и второго порядков. Выборочный момент  $a_1$  называется *выборочным средним* и имеет специальное обозначение  $\bar{X}$ ; *выборочная дисперсия*  $m_2 = a_2 - \bar{X}^2$  обычно обозначается  $S^2$ . Соответствующим образом определяются *выборочный коэффициент асимметрии*  $g_1 = m_3/S^3$  и *выборочный коэффициент эксцесса*  $g_2 = m_4/S^4 - 3$ .

При выборе из  $m$ -мерного,  $m > 1$ , распределения эмпирическое распределение также приписывает массу  $n^{-1}$  каждому выборочному (векторному) значению  $X_i = (X_{1i}, \dots, X_{mi})$ ,  $i = 1, \dots, n$ . В соответствии с этим мы можем определить вектор выборочных средних  $\bar{X} = (\bar{X}_1, \dots, \bar{X}_m)$  с компонентами

$$\bar{X}_k = \frac{1}{n} \sum_{i=1}^n X_{ki}, \quad k = 1, \dots, m,$$

*выборочную ковариационную матрицу*  $S = \|S_{kj}\|$  с элементами

$$S_{kj} = \frac{1}{n} \sum_{i=1}^n (X_{ki} - \bar{X}_k)(X_{ji} - \bar{X}_j) = \frac{1}{n} \sum_{i=1}^n X_{ki}X_{ji} - \bar{X}_k\bar{X}_j, \quad k, j = 1, \dots, m,$$

и матрицу *выборочных коэффициентов корреляции*  $R = \|r_{kj}\|$  с элементами  $r_{kj} = S_{kj}/\sqrt{S_{kk}S_{jj}}$ ,  $k, j = 1, \dots, m$ . Смешанные моменты более высоких порядков в многомерном случае обычно не вычисляются.

Если выбор происходит из распределения, для которого справедлива теорема сложения (предложение 12.2 курса ТВ), то распределение выборочного среднего устанавливается достаточно просто. В общем же случае



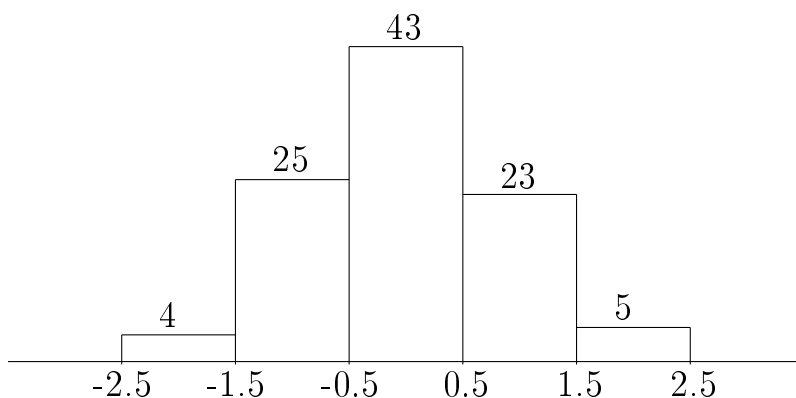
можно только утверждать об асимптотической ( $n \rightarrow \infty$ ) нормальности этой статистики при условии существования второго момента у теоретического распределения. Аналогичное утверждение справедливо и для моментов любого  $k$ -го порядка, если у  $F(x)$  существует момент порядка  $2k$ .

Обратимся теперь к выборочным аналогам квантилей распределения  $F$  наблюдаемой случайной величины  $X$ . Напомним, что для непрерывного распределения квантиль порядка  $p$  определялась как решение  $x_p$  уравнения  $F(x) = p$ , а в случае дискретного распределения – как наибольшее  $x = x_p$  из носителя распределения, при котором  $F(x_p) \leq p$ . Поскольку эмпирическое распределение дискретно, и его функция распределения  $F_n(\cdot)$  возрастает скачками в точках, соответствующих компонентам вариационного ряда, то *выборочная квантиль* порядка  $p$  полагается равной порядковой статистике  $X_{([np])}$ , где  $[x]$ , как обычно, означает целую часть  $x$ . Естественно, для повышения точности оценки истинной квантили  $x_p$  можно проводить интерполяцию между статистиками  $X_{([np])}$  и  $X_{([np]+1)}$ . Так, *выборочная медиана*, будучи квантилью порядка  $p = 0.5$ , обычно определяется как  $(X_{([n/2])} + X_{([n/2]+1)}) / 2$ . Что же касается оценки моды распределения – точки наибольшего сгущения выборочных данных, то здесь нам придется обратиться к выборочным аналогам функции плотности.

При больших объемах наблюдений выборочные данные обычно подвергаются группировке, при этом индивидуальные выборочные значения не приводятся, а указываются лишь количества наблюдений, попавших в интервалы некоторого разбиения множества  $X$  значений наблюдаемой случайной величины. Поясним процедуру группировки на примере выборки из непрерывного одномерного распределения, когда  $X = \mathbb{R}$ .

В декартовой системе координат ось абсцисс разбивается на  $r \geq 2$  интервалов  $(-\infty, a_1]$ ,  $(a_1, a_2]$ ,  $\dots$ ,  $(a_{r-2}, a_{r-1}]$ ,  $(a_{r-1}, +\infty)$ , причем внутренние интервалы выбираются, как правило, одинаковой длины:  $a_i - a_{i-1} = \Delta$ ,  $i = 2, \dots, r - 1$ . Выборочные данные сортируются по интервалам разбиения и подсчитываются частоты  $n_i$ ,  $i = 1, \dots, r$  попадания данных в каждый интервал. Над каждым внутренним интервалом рисуется прямоугольник высоты  $n_i/n\Delta$ , так что площадь  $n_i/n$  каждого прямоугольника с номером  $i = 2, \dots, r - 1$  служит реализацией частотной оценки  $\nu_i/n$  вероятности попадания наблюдаемой случайной величины  $X$  в соответствующий интервал. Здесь  $\nu_i$  – статистика, которую можно записать с помощью индикаторов событий  $A_{ij} = \{X_j \in (a_{i-1}, a_i]\}$ ,  $i = 1, \dots, r$ ,  $a_0 = -\infty$ ,  $a_r =$

$+\infty$ ,  $j = 1, \dots, n$ , а именно  $\nu_i = \sum_{j=1}^n I(A_{ij})$ . Полученная таким образом случайная ступенчатая функция, принимающая нулевые значения на крайних интервалах  $(-\infty, a_1]$ ,  $(a_{r-1}, +\infty)$  и равная  $\nu_i/n\Delta$  на внутренних интервалах с номерами  $i = 2, \dots, r-1$ , называется *гистограммной оценкой*  $f_n$  функции плотности  $f(x)$ ,  $x \in \mathbb{R}$  распределения  $X$ , а ее реализация ( $\nu_i$  заменяются на наблюдаемые частоты  $n_i$ ,  $i = 1, \dots, r$ ) – *гистограммой* выборки  $x^{(n)}$ .



В математической статистике существует ряд теорем, устанавливающих, что при определенных условиях на плотность  $f$ , гистограммная оценка  $f_n(x) \xrightarrow{P} f(x)$  при любом  $x \in \mathbb{R}$ , если  $n \rightarrow \infty$  и одновременно  $r \rightarrow \infty$ , а  $\Delta \rightarrow 0$  со скоростью, зависящей определенным образом от  $n$  и  $r$ .

В случае гистограммной оценки функции плотности естественно считать выборочным аналогом (оценкой) моды распределения  $X$  середину интервала разбиения, в котором гистограмма принимает наибольшее значение.

Заметим также, что вектор частот  $(\nu_1, \dots, \nu_r)$  имеет мультиномиальное распределение  $\mathcal{M}(r, n, \mathbf{p})$  с вероятностями исходов  $p_i = F(a_i) - F(a_{i-1})$ ,  $i = 1, \dots, r$ , что позволяет найти распределение оценки  $f_n(x)$  при любом  $x \in \mathbb{R}$  и построить критерий согласия выборочных данных с гипотезой о виде распределения наблюдаемой случайной величины. Это широко используемый на практике *критерий хи-квадрат*, основанный на статистике (сравните с критерием Колмогорова  $D_n$ )

$$X^2 = \sum_1^r \frac{(\nu_i - np_i)^2}{np_i}.$$

Асимптотическое распределение этой статистики мы изучим в параграфе, посвященном статистической проверке гипотез.

Итак, мы рассмотрели основные выборочные аналоги распределения наблюдаемой случайной величины и его основных характеристик. Мы высказали также ряд утверждений о распределении этих статистик, что позволит нам в последующем вычислять последствия от их использования в качестве решающих функций. Для того, чтобы уяснить, насколько важно знать хотя бы среднее значение статистики, претендующей на роль решающей функции, обратимся снова к примеру 1.1 по аттестации партии дизельного топлива, где обсуждалась сопутствующая проблема оценки дисперсии  $\sigma^2$  наблюдаемой случайной величины  $X \sim \mathcal{N}(\mu, \sigma^2)$ .

Предлагалось оценивать  $\sigma^2$  по накопленному в лаборатории архиву испытаний аттестуемых партий дизельного топлива, то есть по данным большого числа  $N$  выборок  $X_1^{(n)}, \dots, X_N^{(n)}$  малого объема  $n$ . Каждая  $i$ -ая выборка извлекается из нормального  $(\mu_i, \sigma^2)$  распределения, причем средние  $\mu_i$  могут быть различными для разных выборок,  $i = 1, \dots, N$ , но дисперсия  $\sigma^2$  у всех выборок одна и та же. Предлагается следующая оценка  $\sigma^2$ . В каждой выборке вычисляется выборочная дисперсия  $S_i^2$ ,  $i = 1, \dots, N$ , и затем берется их арифметическое среднее:  $\hat{\sigma}_N^2 = (1/N) \sum_1^N S_i^2$ . Распределение каждой  $S_i^2$  не зависит от  $\mu_i$ ,  $i = 1, \dots, N$ , поскольку выборочная дисперсия инвариантна относительно сдвигов  $X_k \rightarrow X_k + a$ . Следовательно, предлагаемая оценка есть нормированная на  $N$  сумма независимых, одинаково распределенных случайных величин – копий статистики  $S^2 = (1/n) \sum_1^n (X_k - \bar{X})^2$ , и в силу закона больших чисел  $\hat{\sigma}_N^2 \xrightarrow{P} \mathbf{E}S^2$  при неограниченном возрастании объема  $N$  архивных данных. Вычислим это математическое ожидание:

$$\mathbf{E}S^2 = \mathbf{E} \left( \frac{1}{n} \sum_1^n X_k^2 - \bar{X}^2 \right) = \mathbf{E}X^2 - \mathbf{E}\bar{X}^2 = \sigma^2 + \mu^2 - \frac{\sigma^2}{n} - \mu^2 = \frac{n-1}{n} \sigma^2,$$

поскольку  $\alpha_2 = \mathbf{E}X^2 = \mathbf{D}X + \mathbf{E}^2X$ , а  $\bar{X} \sim \mathcal{N}(\mu, \sigma^2/n)$ .

Таким образом, предлагаемая оценка обладает значительным смещением при малом объеме  $n$  испытаний каждой партии дизельного топлива. Например, в случае  $n = 2$  мы занижаем дисперсию в два раза, поскольку  $\hat{\sigma}_N^2 \xrightarrow{P} \sigma^2/2$ . Естественно, этот дефект легко устраним – достаточно использовать исправленную на смещение оценку  $\tilde{\sigma}_N^2 = (n/(n-1))\hat{\sigma}_N^2$ .

В завершении этого параграфа мы изучим еще один класс замечательных статистик, используя которые можно редуцировать выборочные данные только к их значениям без потери информации. К сожалению, не все статистические структуры обладают такими статистиками, но, по существу, только в тех структурах, где имеются достаточные статистики, возможно построение оптимального статистического правила, на котором достигается минимум риска.

Идея, состоящая в том, что в определенных случаях для принятия решения без увеличения риска *достаточно* знать только значения некоторых статистик, а не все выборочные данные, не требует введения специальных мер информации, содержащейся в выборочных данных и статистиках, — все становится ясным при рассмотрении следующей простейшей задачи, с которой мы имели дело в самом начале курса теории вероятностей.

Предположим, что мы хотим узнать вероятность наследования доминантного признака в опытах Менделя и располагаем результатами  $x_1, \dots, x_n$  скрещиваний  $n$  пар, где, как обычно, каждое  $x_i$  есть индикатор наследования признака,  $i = 1, \dots, n$ , а совокупность выборочных данных представляет реализацию случайной выборки  $X_1, \dots, X_n$  из двухточечного распределения с функцией плотности  $f(x | \theta) = P_\theta(X = x) = \theta^x(1 - \theta)^{1-x}$ , отличной от нуля только в точках  $x = 0$  и  $1$ . Частотная оценка  $\hat{\theta}_n = T/n$  вероятности  $\theta$  наследования признака определяется статистикой  $T = \sum_1^n X_k$ , выборочное значение  $t = \sum_1^n x_k$  которой соответствует числу потомков в эксперименте, наследовавших доминантный признак. Естественно, возникает вопрос, а нельзя ли извлечь дополнительную информацию о величине параметра  $\theta$  из номеров  $k_1, \dots, k_t$  выборочных данных, принявших значение 1? Нетрудно понять, что это возможно только в том случае, если *распределение выборочного вектора  $X^{(n)}$  при условии, что статистика  $T$  приняла фиксированное значение  $t$ , зависит от параметра  $\theta$* . Действительно, если мы будем наблюдать случайную величину, которая не имеет никакого отношения к интересующему нас параметру, то откуда этой информации взяться? Итак, найдем условное распределение  $X^{(n)}$  относительно  $T$ .

Используя формулу условной вероятности, получаем, что

$$P_\theta \left( X^{(n)} = x^{(n)} \mid T = t \right) = \frac{P_\theta \left( \{X^{(n)} = x^{(n)}\} \wedge \left\{ \sum_1^n X_k = t \right\} \right)}{P_\theta \left( \sum_1^n X_k = t \right)}.$$

Если значения компонент вектора  $x^{(n)}$  таковы, что  $\sum_1^n x_k \neq t$ , то события  $X^{(n)} = x^{(n)}$  и  $\sum_1^n X_k = t$ , очевидно, несовместны, и поэтому в этом случае условная вероятность равна нулю (не зависит от  $\theta$ ). Если же  $\sum_1^n x_k = t$ , то событие  $X^{(n)} = x^{(n)}$  влечет событие  $\sum_1^n X_k = t$ , и формула для вычисления условной вероятности упрощается:

$$P_\theta \left( X^{(n)} = x^{(n)} \mid T = t \right) = \frac{P_\theta \left( X^{(n)} = x^{(n)} \right)}{P_\theta \left( \sum_1^n X_k = t \right)}.$$

Так как

$$P_\theta \left( X^{(n)} = x^{(n)} \right) = f_n(X^{(n)} \mid \theta) = \theta^{\sum_1^n x_k} (1 - \theta)^{n - \sum_1^n x_k},$$

$$P_\theta \left( \sum_1^n X_k = t \right) = C_n^t \theta^t (1 - \theta)^{n-t},$$

то в случае  $\sum_1^n x_k = t$  условное распределение выборочного вектора  $X^{(n)}$  относительно статистики  $T$  имеет вид

$$P_\theta \left( X^{(n)} = x^{(n)} \mid T = t \right) = \frac{1}{C_n^t},$$

и также не зависит от  $\theta$ .

Итак, наши выкладки показывают, что распределение выборочного вектора на “плоскости”  $\sum_1^n X_k = t$  не зависит от  $\theta$ , и поэтому расположение значений  $x_k = 1$  в последовательности  $x_1, \dots, x_n$  при фиксированном количестве таких значений не несет информации о параметре  $\theta$ .

**Определение 2.1.** Статистика  $T = T(X^{(n)})$  называется *достаточной* для статистической структуры  $\mathcal{P}_n = \{P_{\theta,n}, \theta \in \Theta\}$ , если условное распределение выборочного вектора  $X^{(n)}$  относительно статистики  $T$  не зависит от  $\theta$ .

В общей теории статистического вывода в рамках более общего определения статистического правила устанавливается замечательный факт: *если статистическая структура обладает достаточной статистикой  $T$ , то, каково бы ни было статистическое правило  $\delta = \delta(X^{(n)})$ , всегда существует правило  $\delta^* = \delta^*(T)$ , основанное только на  $T$ , риск которого совпадает с риском правила  $\delta$* . Таким образом, построение оптимальных

статистических правил следует начинать с поиска достаточных статистик. Следующая теорема дает критерий существования у статистических структур достаточных статистик и, одновременно, указывает простой способ их нахождения.

**Теорема 2.1.** *Для того, чтобы  $T = T(X^{(n)})$  была достаточной статистикой для статистической структуры, определяемой функцией плотности  $f_n(x^{(n)} | \theta)$ , необходимо и достаточно, чтобы эта функция допускала представление*

$$f_n(x^{(n)} | \theta) = g_\theta \left( T(x^{(n)}) \right) h(x^{(n)}), \quad (1)$$

где функция  $h$  не зависит от параметра  $\theta$ , а функция  $g$  зависит от  $\theta$  и аргумента  $x^{(n)}$  только через значения  $T(x^{(n)})$  статистики  $T = T(X^{(n)})$ ,

Доказательство теоремы мы проведем только для дискретного распределения наблюдаемой случайной величины, когда функция плотности выборки  $f_n(x^{(n)} | \theta) = P_\theta (X^{(n)} = x^{(n)})$ . В случае непрерывного распределения схема доказательства та же, но придется делать замену в  $n$ -кратном интеграле.

*Достаточность.* Пусть выполняется факторизационное представление (1); требуется показать, что условное распределение  $X^{(n)}$  относительно  $T$  не зависит от  $\theta$ . Как и в только что рассмотренном примере с двухточечным распределением, воспользуемся формулой условной вероятности для вычисления условной плотности  $X^{(n)}$  относительно  $T$ :

$$P_\theta \left( X^{(n)} = x^{(n)} | T(X^{(n)}) = t \right) = \frac{P_\theta \left( \{X^{(n)} = x^{(n)}\} \wedge \{T(X^{(n)}) = t\} \right)}{P_\theta(T(X^{(n)}) = t)}.$$

События, стоящие в числителе, будут несовместными, если  $T(x^{(n)}) \neq t$ , и в этом случае условная вероятность равна нулю (не зависит от  $\theta$ ). Если же  $T(x^{(n)}) = t$ , то первое по порядку событие в числителе влечет второе, и поэтому формула для вычисления условной вероятности упрощается:

$$P_\theta \left( X^{(n)} = x^{(n)} | T(X^{(n)}) = t \right) = \frac{P_\theta (X^{(n)} = x^{(n)})}{P_\theta(T(X^{(n)}) = t)}.$$

Так как  $P_\theta(X^{(n)} = x^{(n)}) = f_n(x^{(n)} | \theta)$ , то используя представление (1), получаем, что (напомним,  $T(x^{(n)}) = t$ )

$$P_\theta \left( X^{(n)} = x^{(n)} | T(X^{(n)}) = t \right) =$$

$$\frac{g_{\theta}(T(x^{(n)}))h(x^{(n)})}{\sum_{y^{(n)}:T(y^{(n)})=t} g_{\theta}(T(y^{(n)}))h(y^{(n)})} = \frac{h(x^{(n)})}{\sum_{y^{(n)}:T(y^{(n)})=t} h(y^{(n)})}.$$

Таким образом, условное распределение не зависит от  $\theta$ , и поэтому статистика  $T$  достаточна для  $\mathcal{P}_n$ .

*Необходимость.* Пусть  $T$  – достаточная статистика, так что условное распределение  $P_{\theta}(X^{(n)} = x^{(n)} | T(X^{(n)}) = t) = K(x^{(n)}, t)$ , где функция  $K$  не зависит от  $\theta$ . Требуется показать, что в этом случае для функции плотности выборки справедливо представление (1).

Имеем

$$\begin{aligned} f_n(x^{(n)} | \theta) &= P_{\theta}(X^{(n)} = x^{(n)}) = \\ &P_{\theta}(\{X^{(n)} = x^{(n)}\} \wedge \{T(X^{(n)}) = T(x^{(n)})\}) = \\ &P_{\theta}(T(X^{(n)}) = T(x^{(n)})) \cdot P_{\theta}(X^{(n)} = x^{(n)} | T(X^{(n)}) = T(x^{(n)})). \end{aligned}$$

Мы получили представление (1) с  $g_{\theta}(T(x^{(n)})) = P_{\theta}(T(X^{(n)}) = T(x^{(n)}))$  и  $h(x^{(n)}) = K(x^{(n)}, T(x^{(n)}))$ . Теорема доказана.

Рассмотрим несколько примеров на применения полученного критерия достаточности к статистическим структурам, соответствующим вероятностным моделям из нашего курса теории вероятностей. Начнем с двухточечного распределения (выбор в схеме Бернулли), где мы непосредственными вычислениями условного распределения убедились в достаточности статистики, реализующей число успешных испытаний, – посмотрим, как это делается с помощью представления (1).

1<sup>0</sup>. *Двухточечное распределение*  $B(1, \theta)$  имеет функцию плотности

$$f(x | \theta) = \theta^x (1 - \theta)^{1-x},$$

отличную от нуля только в точках  $x = 0$  и  $1$ . Параметрическое пространство этого распределения  $\Theta = (0; 1)$ , а функция плотности случайной выборки

$$f_n(x^{(n)} | \theta) = \theta^{\sum_1^n x_k} (1 - \theta)^{n - \sum_1^n x_k}.$$

Представление (1) выполняется с  $h(x^{(n)}) \equiv 1$  и  $T(x^{(n)}) = \sum_1^n x_k$ . Следовательно,  $T = \sum_1^n X_k$  – достаточная статистика.

2<sup>0</sup>. Распределение Пуассона  $P(\theta)$ , для которого

$$f(x | \theta) = \frac{\theta^x e^{-\theta}}{x!}, \quad x = 0, 1, \dots, \quad \Theta = \mathbb{R}_+,$$

функция плотности выборки

$$f_n(x^{(n)} | \theta) = \theta \sum_1^n x_k e^{-n\theta} / \prod_1^n x_k!.$$

Следовательно, в представлении (1)

$$h(x^{(n)}) = \left[ \prod_1^n x_k! \right]^{-1},$$

и  $T = \sum_1^n X_k$  – достаточная статистика.

3<sup>0</sup>. Показательное распределение  $E(\theta)$  с

$$f(x | \theta) = \frac{1}{\theta} \exp \left\{ -\frac{x}{\theta} \right\}, \quad x \geq 0, \quad \Theta = \mathbb{R}_+,$$

и

$$f_n(x^{(n)} | \theta) = \frac{1}{\theta^n} \exp \left\{ -\frac{1}{\theta} \sum_1^n x_k \right\}$$

Также обладает достаточной статистикой  $T = \sum_1^n X_k$ .

4<sup>0</sup>. Равномерное распределение  $U(a, b)$ , функция плотности которого

$$f(x | \theta) = \frac{I_{[a; b]}(x)}{b - a}$$

отлична от нуля и постоянна на отрезке  $[a; b]$ , на что указывает стоящая в числителе индикаторная функция отрезка  $[a; b]$ . В этом распределении  $\theta = (a, b)$  – двумерный параметр и параметрическое пространство  $\Theta = \{(a, b) : (a, b) \in \mathbb{R}^2, a < b\}$ . Статистическая структура определяется функцией плотности

$$f_n(x^{(n)} | \theta) = \frac{\prod_1^n I_{[a; b]}(x_k)}{(b - a)^n},$$

и поскольку произведение индикаторов, стоящее в числителе этой функции, принимает значение 1 в случае

$$a \leq \min_{1 \leq k \leq n} x_k \leq \max_{1 \leq k \leq n} x_k \leq b$$



и значение 0 при нарушении этих неравенств, то вектор  $T = (X_{(1)}, X_{(n)})$  крайних членов вариационного ряда является достаточной статистикой.

5<sup>0</sup>. *Нормальное распределение*  $\mathcal{N}(\mu, \sigma^2)$ . Это распределение обладает двумерным параметром  $\theta = (\mu, \sigma)$  с областью значений (параметрическим пространством)  $\Theta = \mathbb{R} \times \mathbb{R}_+$ . Функции плотности наблюдаемой случайной величины  $X$  и случайной выборки  $X^{(n)}$  определяются соответственно как

$$f(x | \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(x - \mu)^2}{2\sigma^2} \right\}$$

и

$$f_n(x^{(n)} | \theta) = \frac{1}{(2\pi)^{n/2}\sigma^n} \exp \left\{ -\frac{1}{2\sigma^2} \sum_1^n (x_k - \mu)^2 \right\} =$$

$$\frac{1}{(2\pi)^{n/2}\sigma^n} \exp \left\{ -\frac{1}{2\sigma^2} \left( \sum_1^n x_k^2 - 2\mu \sum_1^n x_k + n\mu^2 \right) \right\}.$$

Последнее выражение для плотности  $X^{(n)}$  показывает, что двумерная статистика  $T = (T_1, T_2)$  с  $T_1 = \sum_1^n X_k$   $T_2 = \sum_1^n X_k^2$  достаточна для статистической структуры нормального распределения. Кроме того, поскольку  $T_1 = n\bar{X}$  и  $T_2 = n(S^2 + \bar{X}^2)$ , то факторизационное равенство (1) указывает на достаточность статистик  $\bar{X}$  и  $S^2$ , которые имеют конкретную статистическую интерпретацию и поэтому более удобны для практического использования. Понятно, что это замечание носит общий характер: *любые взаимно однозначные преобразования достаточной статистики наследуют свойство достаточности.*

Отметим также, что в случае известного (фиксированного)  $\sigma$  статистическая структура имеет параметрическое пространство, совпадающее с областью значений параметра  $\mu$ , и достаточной статистикой будет выборочное среднее  $\bar{X}$ . Аналогичное утверждение имеет место для статистики  $S^2$  при фиксированном  $\mu$ .

6<sup>0</sup>. *Гамма-распределение*  $G(\lambda, a)$  имеет функцию плотности

$$f(x | \theta) = \frac{1}{a^\lambda \Gamma(\lambda)} x^{\lambda-1} \exp \left\{ -\frac{x}{a} \right\}, \quad x > 0, \quad \theta = (a, \lambda), \quad a > 0, \quad \lambda > 0,$$

так что функция плотности выборочного вектора

$$f(x | \theta) = \frac{1}{a^{n\lambda} \Gamma^n(\lambda)} \left[ \prod_1^n x_k \right]^{\lambda-1} \exp \left\{ -\frac{1}{a} \sum_1^n x_k \right\}.$$

Тождество (1) указывает, что достаточной является двумерная статистика  $(\sum_1^n X_k, \prod_1^n X_k)$  или более удобная в вычислительном отношении статистика  $(\sum_1^n X_k, \sum_1^n \ln X_k)$ . Для этого распределения можно сделать то же замечание, что и для нормального: первая компонента достаточной статистики “отвечает” за масштабный параметр  $a$ , в то время как вторая соответствует параметру формы  $\lambda$ .

7<sup>0</sup>. *Биномиальное распределение*  $B(m, p)$ . Это дискретное распределение, сосредоточенное в точках  $x = 0, 1, \dots, m$ , с функцией плотности

$$f(x | \theta) = C_m^x p^x (1 - p)^{m-x},$$

зависящей от двумерного параметра  $\theta = (m, p)$ , первая компонента  $m$  которого может принимать только значения из множества  $\mathbb{N} = \{1, 2, \dots\}$ , а вторая компонента  $p \in (0; 1)$ . Функция плотности выборочного вектора

$$f_n(x^{(n)} | \theta) = \prod_{k=1}^n C_m^{x_k} \cdot p^{\sum_1^n x_k} (1 - p)^{nm - \sum_1^n x_k}.$$

Применение критерия (1) показывает, что для статистической структуры с параметрическим пространством  $\Theta = \mathbb{N} \times (0; 1)$  достаточной статистикой может быть только весь выборочный вектор  $X^{(n)}$ , но если  $\Theta = (0; 1)$  (значение параметра  $m$  известно), то  $\sum_1^n X_k$  — достаточная статистика.

8<sup>0</sup>. *Распределение Коши*  $C(a, b)$  имеет функцию плотности выборочного вектора

$$f_n(x^{(n)} | \theta) = \pi^{-n} b^{-n} \prod_{k=1}^n \left( 1 + \left( \frac{x_k - a}{b} \right)^2 \right)^{-1},$$

и в силу критерия (1) его статистическая структура обладает только *тривиальной* достаточной статистикой  $T = X^{(n)}$ .

Мы не будем выписывать статистические структуры многомерных распределений в силу их чрезвычайной громоздкости, но нетрудно установить по аналогии с рассмотренными примерами, что у структуры мультиномиального распределения  $\mathcal{M}(m, 1, \mathbf{p})$  с  $m \geq 2$  исходами и вектором  $\mathbf{p} = (p_1, \dots, p_m)$  вероятностей соответствующих исходов достаточным будет вектор, состоящий из частот этих исходов в мультиномиальной схеме испытаний, а у структуры многомерного нормального распределения

$\mathcal{N}_m(\mu, \Lambda)$  достаточную статистику образуют вектор выборочных средних и выборочная ковариационная матрица.

На этом завершается вводная часть нашего курса математической статистики. Мы сделали постановку проблемы статистического вывода, провели классификацию основных статистических структур и теперь мы готовы к решению конкретных статистических проблем по оценке параметров распределения наблюдаемой случайной величины и проверке гипотез, касающихся структуры параметрического пространства этого распределения.

### §3. Оценка параметров. Метод моментов

Лекция 5

Мы приступаем к решению статистической проблемы оценки неизвестного значения параметра  $\theta$ , индексирующего семейство  $\mathcal{P} = \{P_\theta, \theta \in \Theta\}$  возможных распределений наблюдаемой случайной величины  $X$ . Будут рассматриваться только конечномерные параметрические пространства  $\Theta = \mathbf{R}^k$ ,  $k \geq 1$ . Информация о значении  $\theta$  поступает к нам в виде выборочных данных  $x^{(n)} = (x_1, \dots, x_n)$  – результатов наблюдений  $n$  независимых копий  $X^{(n)} = (X_1, \dots, X_n)$  случайной величины  $X$ . Напомним, семейство  $\mathcal{P}$  мы назвали вероятностной моделью, а случайный вектор  $X^{(n)}$  – случайной выборкой объема  $n$ .

В этой проблеме, о которой мы несколько раз упоминали в предыдущем параграфе, пространство решений  $\mathcal{D}$  совпадает с параметрическим пространством  $\Theta$ , решающая функция  $\delta = \delta(X^{(n)})$  – статистика с областью значений  $\mathcal{T} = \Theta$  – называется *оценкой параметра  $\theta$*  и обычно обозначается  $\theta_n$ ,  $\hat{\theta}_n$ ,  $\theta_n^*$  и тому подобное. Функции потерь  $L(\theta, d)$  в проблеме оценивания обычно выбираются в виде неубывающей функции расстояния  $|d - \theta|$  (в евклидовой метрике) между значением оценки  $d = \hat{\theta}_n(x^{(n)})$  и истинным значением  $\theta$  оцениваемого параметра.

*Основная задача статистической теории оценивания состоит в построении оценки  $\theta_n^* = \theta_n^*(X^{(n)})$ , минимизирующей равномерно по  $\theta \in \Theta$  функцию риска*

$$R(\theta; \hat{\theta}_n) = \mathbf{E}_\theta L(\theta, \hat{\theta}_n(X^{(n)})).$$

Таким образом, какова бы ни была статистическая оценка  $\hat{\theta}_n$ , для оценки  $\theta_n^*$  с равномерно минимальным риском при любом  $\theta \in \Theta$  справедливо неравенство  $R(\theta; \theta_n^*) \leq R(\theta; \hat{\theta}_n)$ .

Мы рассмотрим одно из решений этой задачи в случае оценки скалярного параметра ( $\Theta = \mathbf{R}$ ) при квадратичной функции потерь  $L(\theta, d) = (d - \theta)^2$ , но сначала познакомимся с традиционно используемыми в статистической практике методами оценки параметров и изучим распределение этих оценок с целью вычисления их функции риска.

Конечно, далеко не все используемые на практике методы приводят к оптимальным оценкам, иногда бывает трудно найти оценку, обладающую хоть какими-нибудь привлекательными свойствами. Понятно, что считать оценкой любое измеримое отображение выборочного пространства  $\mathcal{X}^n$  в параметрическое пространство  $\Theta$  не совсем разумно, и поэтому мы введем

некоторые условия, которым должна удовлетворять статистика  $\hat{\theta}_n$ , чтобы претендовать на роль *оценки*. Разрабатывая в дальнейшем методы оценивания и предлагая конкретные оценки, мы всегда будем проверять выполнимость этих условий.

**Определение 3.1.** Оценка  $\hat{\theta}_n$  параметра  $\theta$  называется *состоятельной*, если  $\hat{\theta}_n(X^{(n)}) \xrightarrow{P} \theta$  при любом  $\theta \in \Theta$ , когда объем выборки  $n \rightarrow \infty$ . Оценка  $\hat{\theta}_n$  называется *несмещенной в среднем*, если  $\mathbf{E}_\theta \hat{\theta}_n(X^{(n)}) = \theta$ , каково бы ни было значение  $\theta \in \Theta$ .

Напомним, что  $\hat{\theta}_n(X^{(n)}) \xrightarrow{P} \theta$  означает, что для любого  $\varepsilon > 0$

$$\lim_{n \rightarrow \infty} P_\theta \left( \left| \hat{\theta}_n(X^{(n)}) - \theta \right| > \varepsilon \right) = 0,$$

или, что то же,

$$\lim_{n \rightarrow \infty} P_\theta \left( \left| \hat{\theta}_n(X^{(n)}) - \theta \right| \leq \varepsilon \right) = 1. \quad (1)$$

Здесь, как обычно, в случае векторного параметра  $\theta$  запись  $|\theta_1 - \theta_2|$  означает расстояние между точками  $\theta_1$  и  $\theta_2$  эвклидова пространства  $\Theta$ .

В предыдущем параграфе мы показали, что выборочные моменты  $a_i = (1/n) \sum_1^n X_j^i$  являются состоятельными оценками соответствующих “теоретических” моментов  $\alpha_i = \mathbf{E}_\theta X^i$ , которые являются функциями оцениваемого параметра:  $\alpha_i = \alpha_i(\theta)$ ,  $i = 1, 2, \dots$ . Этот результат указывает нам довольно простой метод построения состоятельных оценок в случае существования у распределения  $P_\theta$  наблюдаемой случайной величины  $X$  момента порядка  $k$ , где  $k$  – число компонент  $\theta_1, \dots, \theta_k$  оцениваемого параметрического вектора  $\theta$ .

Приравняем теоретические моменты выборочным и разрешим полученную таким образом систему уравнений  $\alpha_i(\theta_1, \dots, \theta_k) = a_i$ ,  $i = 1, \dots, k$  относительно переменных  $\theta_1, \dots, \theta_k$ . Любое решение  $\hat{\theta}_n(\mathbf{a}) = (\hat{\theta}_{1n}(\mathbf{a}), \dots, \hat{\theta}_{kn}(\mathbf{a}))$ ,  $\mathbf{a} = (a_1, \dots, a_n)$ , этой системы называется *оценкой  $\theta$  по методу моментов*, и прежде, чем исследовать свойства таких оценок, рассмотрим несколько примеров на применения *метода моментов*.

В курсе теории вероятностей, изучая новые вероятностные модели, мы всегда вычисляли их моментные характеристики. Например, мы знаем, что средние значения двухточечного распределения  $B(1, \theta)$ , распределения Пуассона  $P(\theta)$  и показательного распределения  $E(\theta)$  равны  $\theta$ . Следовательно, выборочное среднее  $\bar{X}$  есть оценка по методу моментов параметра  $\theta$

любого из этих распределений. Легко видеть, что эта оценка состоятельна и несмещена. Точно так же у нормального распределения  $\mathcal{N}(\mu, \sigma^2)$  параметр  $\mu$  означает среднее значение, а  $\sigma^2$  – дисперсию этого распределения. Следовательно, выборочное среднее  $\bar{X}$  и выборочная дисперсия  $S^2$  есть состоятельные оценки соответствующих компонент  $\mu$  и  $\sigma^2$  параметрического вектора  $\theta = (\mu, \sigma^2)$ . Исправляя смещение оценки  $S^2$  компоненты  $\sigma^2$ , получаем несмещенную оценку  $\theta$ . Замечательно то, что все оценки являются достаточными статистиками, и это обстоятельство, как будет видно в дальнейшем, предопределяет их оптимальные свойства. Распределение оценки  $\bar{X}$  легко получить, используя теоремы сложения для распределений В, Р, Е и  $\mathcal{N}$ , распределение же  $S^2$  при выборе из нормального распределения мы найдем несколько позже.

Рассмотрим теперь примеры, в которых приходится решать систему уравнений, и найденные оценки по методу моментов не являются функциями достаточных статистик.

*Пример 3.1. Оценка параметров биномиального распределения В( $m, p$ ).* Проблема состоит в оценке обеих компонент  $m$  и  $p$  двумерного параметра  $\theta = (m, p)$ . Из курса теории вероятностей нам известно, что среднее значение биномиального распределения равно  $mp$ , а дисперсия –  $mp(1-p)$ . Приравнявая эти теоретические моменты их выборочным аналогам, получаем систему для определения оценок по методу моментов:  $mp = \bar{X}$ ,  $mp(1-p) = S^2$ . Разделив второе уравнение на первое, находим оценку  $\hat{p}_n = (\bar{X} - S^2)/\bar{X}$  параметра  $p$ , после чего, обращаясь к первому уравнению, находим оценку  $\hat{m}_n = \bar{X}^2/(\bar{X} - S^2)$  параметра  $m$ . Легко показать, что эти оценки обладают свойством состоятельности (общий метод доказательства таких утверждений смотрите в приведенной ниже теореме 3.1), но при малых  $n$  велика вероятность получить отрицательные значения оценок, оценка параметра  $m$ , как правило, не будет целым числом, наконец, можно показать, что оценка  $\hat{p}_n^* = \bar{X}$  параметра  $p$  будет обладать меньшим квадратичным риском, чем оценка  $\hat{p}_n$ . Все это, конечно, печально, однако другие методы, приводящие к более точным оценкам, обладают значительными вычислительными трудностями.

*Пример 3.2. Оценка параметров гамма-распределения G( $\lambda, a$ ).* У этого двухпараметрического распределения среднее равно  $\lambda a$ , а дисперсия –  $\lambda a^2$ . Решение системы уравнений  $\lambda a = \bar{X}$ ,  $\lambda a^2 = S^2$  дает оценки  $\hat{a}_n = S^2/\bar{X}$ ,  $\hat{\lambda}_n = \bar{X}^2/S^2$ , которые, как и в предыдущем примере, не являются

функциями достаточной статистики  $\left(\sum_1^n X_k, \prod_1^n X_k\right)$ , и как показывают не совсем простые вычисления, их риски далеки от возможного минимума. Тем не менее, очевидная вычислительная простота оценок параметров гамма-распределения по методу моментов обеспечивает их популярность в практических применениях.

Изучим теперь асимптотические свойства оценок по методу моментов – установим условия их состоятельности и исследуем поведение их распределений при больших объемах выборок. Для простоты мы ограничимся случаем одномерного параметра  $\theta$ , оценка которого определяется решением уравнения  $\mu(\theta) = \mathbf{E}_\theta X = \bar{X}$ , и предположим, что это уравнение имеет единственное решение  $\hat{\theta}_n = h(\bar{X})$ . Понятно, что  $h(\cdot) = \mu^{-1}(\cdot)$ , так что  $h(\mu(\theta)) \equiv \theta$ . О возможности распространения наших результатов на случай векторного  $\theta$  мы поговорим отдельно.

**Теорема 3.1.** *Если наблюдаемая случайная величина  $X$  имеет конечное среднее значение  $\mu = \mu(\theta)$  и функция  $h(\cdot)$  непрерывна в области значений выборочного среднего  $\bar{X}$ , то  $\hat{\theta}_n = h(\bar{X})$  является состоятельной оценкой параметра  $\theta$  по методу моментов.*

*Доказательство.* Мы воспользуемся формулой (1) в определении состоятельности оценки и покажем, что для любых  $\varepsilon > 0$  и  $\alpha > 0$  существует такое  $N(\varepsilon, \alpha)$ , что для всех  $n > N(\varepsilon, \alpha)$  вероятность

$$P_\theta ( |h(\bar{X}) - \theta| \leq \varepsilon ) > 1 - \alpha. \quad (2)$$

Поскольку  $h(\mu) = h(\mu(\theta)) = \theta$  и  $\bar{X} \xrightarrow{P} \mu$ , то нам достаточно показать, что свойство (или определение) непрерывности функции:  $h(x) \rightarrow h(\mu)$  при  $x \rightarrow \mu$ , остается справедливым при замене обычной сходимости " $\rightarrow$ " на сходимость по вероятности " $\xrightarrow{P}$ ", то есть  $\bar{X} \xrightarrow{P} \mu$  влечет  $h(\bar{X}) \xrightarrow{P} h(\mu)$ . Это почти очевидно, поскольку событие, состоящее в попадании в окрестность нуля случайной величины  $|\bar{X} - \mu|$ , влечет аналогичное событие для случайной величины  $|h(\bar{X}) - h(\mu)|$ , но все же проведем строгое доказательство на языке " $\varepsilon - \delta$ ".

Так как  $\bar{X} \xrightarrow{P} \mu(\theta)$ , а  $h(\cdot)$  – непрерывная функция, то найдутся такие  $\delta = \delta(\varepsilon, \alpha)$  и  $N = N(\varepsilon, \alpha)$ , что

$$P(|\bar{X} - \mu(\theta)| < \delta) > 1 - \alpha \quad (3)$$

для всех  $n > N$  и событие  $|\bar{X} - \mu(\theta)| < \delta$  повлечет событие  $|h(\bar{X}) -$

$|h(\mu(\theta)) - h(\bar{X})| = |h(\bar{X}) - \theta| \leq \varepsilon$ . В силу этого неравенство (2) становится следствием неравенства (3). Теорема доказана.

Анализ доказательства показывает, что теорема состоятельности остается справедливой в случае векторного параметра  $\theta$ , если воспользоваться определением непрерывности векторной функции от векторного аргумента, связав его с расстояниями в евклидовых пространствах значений функции и ее аргумента.

Обратимся теперь к асимптотическому анализу распределения оценки  $\hat{\theta}_n = h(\bar{X})$  при  $n \rightarrow \infty$ . Понятно, что в данном случае употребление термина “оценка” применительно к функции  $h(\bar{X})$  ничего особенно не добавляет – речь идет просто об асимптотическом распределении статистики, имеющей вид функции от выборочного среднего.

**Теорема 3.2.** Если  $X^{(n)}$  – случайная выборка из распределения с конечными средним значением  $\mu$  и дисперсией  $\sigma^2$ , а функция  $h(x)$  обладает ограниченной второй производной  $h''(x)$  в некоторой окрестности точки  $x = \mu$ , то

$$\lim_{n \rightarrow \infty} P(\sqrt{n}(h(\bar{X}) - h(\mu)) < x) = \Phi\left(\frac{x}{\sigma|h'(\mu)|}\right). \quad (4)$$

где  $\Phi(\cdot)$  – функция распределения стандартного нормального закона.

Доказательство.. Понятно, что мы должны воспользоваться центральной предельной теоремой (§14 курса ТВ) применительно к статистике  $\bar{X} = \frac{1}{n} \sum_{k=1}^n X_k$ :

$$\lim_{n \rightarrow \infty} P(\sqrt{n}(\bar{X} - \mu) < x) = \Phi\left(\frac{x}{\sigma}\right). \quad (5)$$

Стандартный прием использования этой теоремы при асимптотическом анализе функций от сумм независимых, одинаково распределенных случайных величин состоит в “линеаризации” таких функций с помощью формулы Тейлора. В нашем случае мы разлагаем функцию  $h(\cdot)$  в окрестности точки  $\bar{X} = \mu$ , используя только два члена разложения:

$$h(\bar{X}) = h(\mu) + (\bar{X} - \mu)h'(\mu) + \frac{(\bar{X} - \mu)^2}{2!}h''(\mu + \lambda(\bar{X} - \mu)),$$

где  $0 \leq \lambda \leq 1$ .

Перепишем это разложение в виде

$$\sqrt{n}(h(\bar{X}) - h(\mu)) = \sqrt{n}(\bar{X} - \mu)h'(\mu) + \frac{\sqrt{n}(\bar{X} - \mu)^2}{2!}h''(\mu + \lambda(\bar{X} - \mu)),$$



представив тем самым случайную величину  $\sqrt{n}(h(\bar{X}) - h(\mu))$  (см. формулу (4)) в виде суммы двух случайных величин, первая из которых в силу формулы (5) имеет предельное нормальное распределение, указанное в правой части (4), а вторая сходится по вероятности к нулю. Действительно,  $h''(\mu + \lambda(\bar{X} - \mu))$  по условию теоремы ограничено с вероятностью сколь угодно близкой к единице, начиная с некоторого  $n$ . В силу неравенства Чебышева (предложение 6.2 курса ТВ) для любого  $\varepsilon > 0$  вероятность

$$P(\sqrt{n}(\bar{X} - \mu)^2 > \varepsilon) \leq \frac{\mathbf{E}\sqrt{n}(\bar{X} - \mu)^2}{\varepsilon} = \frac{\sqrt{n}\mathbf{D}\bar{X}}{\varepsilon} = \frac{\sigma^2}{\varepsilon\sqrt{n}} \rightarrow 0,$$

и поэтому стоящий перед  $h''/2!$  множитель, а с ним и все второе слагаемое, сходятся по вероятности к нулю. Доказательство теоремы завершается ссылкой на предложение 11.1 курса ТВ: если одна последовательность случайных величин имеет невырожденное предельное распределение  $F$ , а вторая сходится по вероятности к нулю, то предельное распределение суммы этих последовательностей совпадает с  $F$ .

Итак, если  $h(\bar{X})$  – оценка  $\hat{\theta}_n$  параметра  $\theta$  по методу моментов, то формула (4) теоремы 3.2 принимает вид

$$\lim_{n \rightarrow \infty} P\left(\sqrt{n}(\hat{\theta}_n - \theta) < x\right) = \Phi\left(\frac{x}{\sigma|h'(\mu)|}\right).$$

Приведем пример на использование аппроксимации распределения функции от выборочного среднего на практике.

**Пример 3.3.** *Оценка надежности изделия с показательным распределением долговечности.* При выпуске изделий обычно указывается их гарантийный срок службы  $t_0$ ; отказ изделия до истечения срока  $t_0$  чреват для поставщика расходами на ремонт или замену изделия. Чтобы планировать расходы на такого рода рекламации со стороны потребителя, поставщик должен знать надежность выпускаемых изделий:  $H(t_0) = P(X > t_0)$ . Величина  $H(t_0)$  указывает среднюю долю изделий, которые не откажут в течение гарантийного срока службы  $t_0$ . Чтобы оценить  $H(t_0)$  проводятся испытания  $n$  изделий, и пусть  $x_1, \dots, x_n$  – наработки на отказ испытуемых изделий, трактуемые как реализации случайной выборки  $X^{(n)}$  из распределения  $F(\cdot | \theta)$ , известного с точностью до значения параметра  $\theta$ . Наконец, пусть нам известно, что долговечность изделий подчиняется закону “отсутствия последствия”, в силу чего  $F(x | \theta) = 1 - \exp\{-x/\theta\}$  – показательное

распределение. В таком случае проблема состоит в оценке параметрической функции  $h(\theta) = \exp\{-t_0/\theta\}$ .

Так как  $\theta = \mathbf{E}X$  – средняя наработка на отказ, то естественно оценить  $\theta$  посредством статистики  $\bar{X}$  (работает метод моментов), а за оценку надежности взять статистику  $h(\bar{X}) = \exp\{-t_0/\bar{X}\}$ . Поскольку наибольшую, с экономической точки зрения, опасность представляет завышение надежности, то нас в первую очередь должна интересовать частота грубых превышений, например, на некоторую заданную величину  $\varepsilon$ . Если положить  $\Delta = \varepsilon\sqrt{n}$ , то вероятность, указывающая частоту грубых превышений, записывается в виде  $R(\theta; h(\bar{X})) = P(\sqrt{n}(h(\bar{X}) - h(\theta)) > \Delta)$ , что позволяет нам непосредственно использовать аппроксимацию (4) при испытаниях достаточно большого числа изделий (насколько большого, это – отдельный вопрос, решить который можно, например, моделируя выборки из показательного распределения с помощью метода Монте-Карло).

Имеем  $\sigma^2 = \theta^2$ ,  $h'(\theta) = t_0\theta^{-2} \exp\{-t_0/\theta\}$ , и формула (4) дает нам следующую аппроксимацию для риска оценки  $h(\bar{X})$ :

$$R(\theta; h(\bar{X})) \approx 1 - \Phi\left(\frac{\Delta\theta}{t_0} \exp\left\{\frac{t_0}{\theta}\right\}\right).$$

Нетрудно показать, что наибольшее значение риска достигается при  $\theta = t_0$  и равно  $\Phi(\Delta \cdot e)$ .

В том случае, когда оценка имеет вид функции от двух и более выборочных моментов, метод асимптотического анализа ее распределения тот же. Например, пусть  $\hat{\theta}_n = h(a_1, a_2)$ , и функция  $h$  удовлетворяет условиям, аналогичным требованиям к  $h$  в теореме 3.2. Используя формулу Тейлора, представим  $h$  в окрестности точки  $(\alpha_1, \alpha_2)$  в следующем виде:

$$h(a_1, a_2) = h(\alpha_1, \alpha_2) + (a_1 - \alpha_1)h'_1(\alpha_1, \alpha_2) + \\ (a_2 - \alpha_2)h'_2(\alpha_1, \alpha_2) + O_p(|a - \alpha|^2),$$

где  $a = (a_1, a_2)$ ,  $\alpha = (\alpha_1, \alpha_2)$ . Легко проверяется, что  $\sqrt{n}|a - \alpha|^2 \xrightarrow{P} 0$ , так что случайная величина  $\sqrt{n}(h(a_1, a_2) - h(\alpha_1, \alpha_2))$  асимптотически нормальна с параметрами, которые выражаются через первые четыре момента наблюдаемой случайной величины  $X$ .

## §4. Оценка параметров. Метод максимального правдоподобия

### Лекция 6

Мы приступаем к изучению более точного метода оценки неизвестного значения параметра. Он превосходит метод моментов и при наличии достаточных статистик дает оптимальные оценки с точки зрения квадратичного риска. Более того, при выполнении определенных условий регулярности этот метод приводит к асимптотически ( $n \rightarrow \infty$ ) оптимальным оценкам для широкого класса вероятностных моделей и практически при любых функциях потерь.

Идея метода состоит в математической формализации “разумного” поведения человека в условиях неопределенности. Представим себе ситуацию, что мы ожидаем появления одного из нескольких событий, вероятности которых нам неизвестны и нас интересуют не столько значения этих вероятностей, сколько то событие, которое происходит наиболее часто. Ситуация осложняется тем, что мы располагаем всего одним испытанием, в результате которого произошло некоторое событие  $A$ . Конечно, мы примем решение, что  $A$  обладает наибольшей вероятностью, и вряд ли можно предложить нечто более разумное, чем такое правило принятия решения.

В этом и состоит *принцип максимального правдоподобия*, который буквально пронизывает всю теорию оптимального статистического вывода. Применение этого принципа к проблеме оценки параметров приводит к следующему статистическому правилу: *если  $x^{(n)}$  – результат наблюдения случайной выборки  $X^{(n)}$ , то за оценку параметра следует брать то его значение, при котором результат  $x^{(n)}$  обладает наибольшим правдоподобием.*

Вы спросите, что такое “правдоподобие” результата  $x^{(n)}$ ? Давайте формализуем это понятие.

Если наблюдается дискретная случайная величина, то естественно называть правдоподобием результата  $x^{(n)}$  при фиксированном значении параметра  $\theta$  вероятность его наблюдения в статистическом эксперименте. Но в дискретном случае эта вероятность совпадает со значением функции плотности в точке  $x^{(n)}$ :  $P_\theta(X^{(n)} = x^{(n)}) = f_n(x^{(n)} | \theta)$ . Следовательно, оценка по методу максимального правдоподобия определяется точкой достижения максимума у функции плотности случайной выборки, то есть

$$\hat{\theta}_n(X^{(n)}) = \arg \max_{\theta \in \Theta} f_n(X^{(n)} | \theta). \quad (1)$$

Рассмотрим сразу же простой пример. Пусть  $X^{(n)}$  – выборка в схеме Бернулли, и мы оцениваем вероятность  $\theta$  успешного исхода. В этой модели

$$f(X^{(n)} | \theta) = \theta^{\sum_1^n X_k} (1 - \theta)^{n - \sum_1^n X_k}.$$

Дифференцируя эту функцию по  $\theta$  и приравнявая производную нулю, находим оценку максимального правдоподобия  $\theta = (1/n) \sum_1^n X_k$ . Это – давно знакомая нам оценка вероятности успеха в испытаниях Бернулли, которую мы получили с помощью моментов и постоянно использовали при иллюстрации закона больших чисел.

Теперь определим правдоподобие в случае выбора из непрерывного распределения с функцией плотности (по мере Лебега)  $f_n(x^{(n)} | \theta)$ ,  $x^{(n)} \in \mathbb{R}^n$ ,  $\theta \in \Theta$ . Пусть  $x^{(n)}$  – совокупность выборочных данных, то есть точка в  $n$ -мерном выборочном пространстве  $\mathbb{R}^n$ . Окружим эту точку прямоугольным параллелепипедом малого размера, скажем,  $V_\varepsilon = \prod_1^n [x_k - \varepsilon/2; x_k + \varepsilon/2]$ . В силу теоремы о среднем для кратного интеграла вероятность того, что выборочный вектор попадет в этот параллелепипед  $P(X^{(n)} \in V_\varepsilon) \sim f_n(x^{(n)} | \theta) \cdot \varepsilon^n$ , когда  $\varepsilon \rightarrow 0$ . Если трактовать эту вероятность, как правдоподобие результата  $x^{(n)}$ , которое, конечно, зависит от выбора малого  $\varepsilon$ , мы видим, что проблема максимизации правдоподобия сводится к проблеме отыскания точки достижения максимума по всем  $\theta \in \Theta$  у функции плотности  $f_n$ . Таким образом, и в случае непрерывного распределения разумно назвать правдоподобием результата  $x^{(n)}$  при фиксированном значении параметра  $\theta$  опять-таки величину функции плотности выборки, то есть  $f_n(x^{(n)} | \theta)$ , и определить оценку максимального правдоподобия той же формулой (1).

Рассмотрим пример на построение такой оценки в случае выбора из непрерывного распределения. Пусть наблюдается случайная величина  $X \sim \mathcal{N}(\mu, \sigma^2)$ , так что функция плотности выборки

$$f_n(x^{(n)} | \theta) = \frac{1}{(2\pi)^{n/2} \sigma^n} \exp \left\{ -\frac{1}{2\sigma^2} \sum_1^n (x_k - \mu)^2 \right\},$$

где  $\theta = (\mu, \sigma)$  – двумерный параметр, значение которого нам неизвестно. В соответствии с формулой (1) необходимо отыскать точку достижения максимума функции  $f_n(X^{(n)} | \mu, \theta)$  по переменным  $\mu \in \mathbb{R}$  и  $\sigma \in \mathbb{R}_+$ . Естественно, логарифм этой функции имеет те же точки экстремума, что и сама функция, но логарифмирование упрощает выкладки, поэтому ищем

максимум функции

$$\mathcal{L}(\theta | X^{(n)}) = \ln f_n(X^{(n)} | \theta) = -\frac{n}{2} \ln 2\pi - n \ln \sigma - \frac{1}{2\sigma^2} \sum_1^n (X_k - \mu)^2.$$

Составляем уравнения, определяющие точки экстремума:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mu} &= \frac{1}{2\sigma^2} \sum_1^n (X_k - \mu) = 0, \\ \frac{\partial \mathcal{L}}{\partial \sigma} &= -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_1^n (X_k - \mu)^2 = 0. \end{aligned}$$

Из первого уравнения сразу находим оценку параметра  $\mu$ :  $\hat{\mu}_n = \bar{X}$ . Подставляя  $\bar{X}$  вместо  $\mu$  во второе уравнение, находим оценку  $\sigma$ :  $\hat{\sigma}_n = S$  (выборочное стандартное отклонение). Очевидно,  $(\bar{X}, S)$  – точка максимума. Таким образом, метод максимального правдоподобия приводит к тем же оценкам  $\bar{X}$  и  $S^2$  параметров  $\mu$  и  $\sigma^2$ , что и метод моментов.

Теперь дадим строгое определение правдоподобия и рассмотрим еще несколько примеров, в которых метод максимального правдоподобия дает оценки, отличные от метода моментов.

**Определение 4.1.** Случайная функция

$$L(\theta | X^{(n)}) = \prod_{i=1}^n f(X_i | \theta)$$

на параметрическом пространстве  $\Theta$  называется *функцией правдоподобия*, а значение ее реализации  $L(\theta_0 | x^{(n)})$  при результате наблюдения  $X^{(n)} = x^{(n)}$  и фиксированном  $\theta = \theta_0$  – *правдоподобием значения  $\theta_0$  при результате  $x^{(n)}$* . Любая точка  $\hat{\theta}_n = \hat{\theta}_n(X^{(n)})$  (статистика) параметрического пространства  $\Theta$ , доставляющая абсолютный максимум функции правдоподобия, называется *оценкой максимального правдоподобия* параметра  $\theta$ .

Поскольку функция правдоподобия представляет собой произведение функций от  $\theta$ , то при отыскании ее максимума методами дифференциального исчисления удобнее иметь дело с логарифмом этой функции. Естественно, точки экстремума у *функции логарифмического правдоподобия*

$$\mathcal{L}(\theta | X^{(n)}) = \sum_{i=1}^n \ln f(X_i | \theta)$$

те же, что и у функции  $L$ , но если функция  $L(\cdot | x^{(n)})$  имеет непрерывные частные производные по компонентам  $\theta_1, \dots, \theta_k$  параметрического вектора  $\theta$ , то проще дифференцировать  $\mathcal{L}$  чем  $L$ . В этом случае система уравнений

$$\frac{\partial \mathcal{L}(\theta | X^{(n)})}{\partial \theta_i} = 0, \quad i = 1, \dots, k \quad (2)$$

называется *уравнениями правдоподобия*. Это еще одна разновидность так называемых *оценочных уравнений*, – в предыдущем параграфе мы имели дело с уравнениями метода моментов.

Любое решение системы уравнений (2), доставляющее максимум функции  $\mathcal{L}(\cdot | X^{(n)})$ , может рассматриваться как оценка  $\theta$  по методу максимального правдоподобия. Мы не будем изучать случаи, когда система (2) имеет несколько решений с возможно одинаковыми значениями функции правдоподобия в этих точках, так что требуются дополнительные априорные знания относительно вероятностной модели, позволяющие выбрать одно из этих решений. Во всех рассмотренных ниже примерах оценка максимального правдоподобия единственна.

**Пример 4.1.** Оценка параметра положения равномерного распределения  $U(0, \theta)$ . Равномерное на отрезке  $[0; \theta]$  распределение имеет функцию плотности  $f(x | \theta) = \theta^{-1}$ , если  $0 \leq x \leq \theta$ , и  $f(x | \theta) = 0$  вне этого отрезка. Следовательно, функция  $L(\theta | X^{(n)})$  отлична от нуля и равна  $\theta^{-n}$  только в области  $\theta \geq X_{(n)} = \max_{1 \leq k \leq n} X_k$ . Ее максимум по  $\theta$  достигается в граничной точке  $\theta = X_{(n)}$ , так что наибольшее значение  $X_{(n)}$  выборки  $X^{(n)}$  есть оценка максимального правдоподобия параметра  $\theta$ .

Легко видеть, что оценка  $\theta$  по методу моментов равна  $2\bar{X}$ . Эта оценка на порядок хуже оценки максимального правдоподобия с точки зрения квадратичного риска  $R(\theta; \hat{\theta}_n) = \mathbf{E}_\theta \left( \hat{\theta}_n(X^{(n)}) - \theta \right)^2$ . Простые вычисления соответствующих математических ожиданий показывают, что  $R(\theta; 2\bar{X}) = O(n^{-1})$ , в то время как  $R(\theta; X_{(n)}) = O(n^{-2})$ . Данный пример интересен тем, что здесь функция правдоподобия не имеет гладкого максимума, и именно это обстоятельство, как будет видно в дальнейшем, обеспечивает такое различное поведение риска рассматриваемых оценок.

**Пример 4.2.** Оценка параметров гамма-распределения  $G(a, \lambda)$ . У этого распределения функция плотности

$$f(x | \theta) = \frac{1}{a^\lambda \Gamma(\lambda)} x^{\lambda-1} \exp \left\{ -\frac{x}{a} \right\}, \quad x > 0, \quad \theta = (a, \lambda),$$

отлична от нуля только на положительной полуоси, и логарифмическое правдоподобие

$$\mathcal{L}(a, \lambda | X^{(n)}) = -n\lambda \ln a - n \ln \Gamma(\lambda) + (\lambda - 1) \sum_1^n \ln X_k - \frac{1}{a} \sum_1^n X_k.$$

Составляем уравнения правдоподобия:

$$\frac{\partial \mathcal{L}}{\partial a} = -\frac{n\lambda}{a} + \frac{1}{a^2} \sum_1^n X_k = 0,$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = -n \ln a - n\psi(\lambda) + \sum_1^n \ln X_k = 0,$$

где  $\psi(\lambda) = d \ln \Gamma(\lambda) / d\lambda$  – так называемая пси-функция Эйлера. Исключая из первого уравнения параметр  $a$  и подставляя полученный результат во второе, получаем трансцендентное уравнение

$$\ln \lambda - \psi(\lambda) = \ln \bar{X} - \frac{1}{n} \sum_1^n \ln X_k,$$

которое в силу свойства монотонности функции  $\ln \lambda - \psi(\lambda)$  имеет единственное решение. При численном решении этого уравнения может оказаться полезной асимптотическая ( $\lambda \rightarrow \infty$ ) формула

$$\ln \lambda - \psi(\lambda) = \frac{1}{2\lambda} + \frac{1}{12\lambda^2} + O\left(\frac{1}{\lambda^4}\right).$$

**Пример 4.3.** *Оценка параметров структурированного среднего при нормальном распределении отклика.* Данная задача весьма часто возникает при калибровке шкалы прибора. Две переменные  $x$  и  $y$  связаны линейным соотношением  $y = a + bx$ , и для градуировки значений  $y$  на шкале прибора необходимо знать значения параметров  $a$  и  $b$  этой зависимости. Однако, для каждого стандартного фиксированного значения  $x$  прибор измеряет значение  $y$  с ошибкой, так что замеры происходят в рамках вероятностной модели  $Y = a + bx + \xi$ , где ошибка измерения (случайная величина)  $\xi$  имеет нормальное распределение с нулевым средним и некоторой дисперсией  $\sigma^2$ , значение которой, как правило, также не известно. Случайная величина  $Y$  обычно называется *откликом* на значение *регрессора*  $x$ ; ее распределение при фиксированном  $x$  очевидно нормально  $(a + bx, \sigma^2)$ .

Для оценки  $a$  и  $b$  производится  $n$  измерений отклика  $y_1, \dots, y_n$  при некоторых фиксированных значениях  $x_1, \dots, x_n$  регрессора  $x$ , оптимальный выбор которых, обеспечивающий наибольшую точность и надежность калибровки, составляет самостоятельную задачу особой области математической статистики – *планирование регрессионных экспериментов*. Мы будем предполагать, что значения  $x_1, \dots, x_n$  априори фиксированы. В таком случае значения  $y_1, \dots, y_n$  представляют реализации  $n$  независимых случайных величин  $Y_1, \dots, Y_n$ , и  $Y_k \sim \mathcal{N}(a + bx_k, \sigma^2)$ ,  $k = 1, \dots, n$ . Совместная функция плотности  $Y_1, \dots, Y_n$  равна

$$f_n(y^{(n)} | a, b, \sigma) = \frac{1}{(2\pi)^{n/2} \sigma^n} \exp \left\{ -\frac{1}{2\sigma^2} \sum_1^n (y_k - a - bx_k)^2 \right\},$$

так что логарифмическая функция правдоподобия, необходимая для оценки параметров  $a$ ,  $b$  и  $\sigma$  методом максимального правдоподобия имеет вид

$$\mathcal{L}(a, b, \sigma | Y^{(n)}) = -\frac{n}{2} \ln 2\pi - n \ln \sigma - \frac{1}{2\sigma^2} \sum_1^n (Y_k - a - bx_k)^2.$$

Вычисляя производные этой функции по переменным  $a$ ,  $b$  и  $\sigma$ , получаем уравнения правдоподобия

$$\begin{aligned} \sum_1^n (Y_k - a - bx_k) &= 0, \\ \sum_1^n x_k (Y_k - a - bx_k) &= 0, \\ n\sigma^2 - \sum_1^n (Y_k - a - bx_k)^2 &= 0. \end{aligned}$$

Конечно, это очень простая система уравнений, решение которой не может вызывать каких-либо затруднений, и мы сразу пишем оценки максимального правдоподобия

$$\hat{a}_n = \bar{Y} - \frac{m_{xY}}{s_x^2} \bar{x}, \quad \hat{b}_n = \frac{m_{xY}}{s_x^2}, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_1^n (Y_k - \hat{a}_n - \hat{b}_n x_k)^2,$$

где

$$\bar{x} = \frac{1}{n} \sum_1^n x_k, \quad \bar{Y} = \frac{1}{n} \sum_1^n Y_k, \quad s_x = \frac{1}{n} \sum_1^n (x_k - \bar{x})^2, \quad S_Y = \frac{1}{n} \sum_1^n (Y_k - \bar{Y})^2,$$



$$m_{xY} = \frac{1}{n} \sum_1^n (x_k - \bar{x})(Y_k - \bar{Y}).$$

Легко видеть, что оценки по методу максимального правдоподобия параметров  $a$  и  $b$  совпадают с их оценками по *методу наименьших квадратов*. В этом методе “выравнивания” экспериментальных данных оценки ищутся из условия минимизации суммы квадратов *невязок*:  $\sum_1^n (Y_k - a - bx_k)^2$ , причем под невязкой понимается разность между откликом  $Y$  и его “теоретическим” средним значением  $a + bx$ .

**Пример 4.4.** *Оценка параметров двумерного нормального распределения: задачи регрессии и прогноза.* Оценка по методу максимального правдоподобия пяти параметров  $\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho$  двумерного нормального распределения с функцией плотности

$$f(x, y | \theta) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \cdot \exp \left\{ -\frac{1}{2(1-\rho^2)} \left( \frac{(x-\mu_1)^2}{\sigma_1^2} - 2\rho \frac{(x-\mu_1)(y-\mu_2)}{\sigma_1\sigma_2} + \frac{(y-\mu_2)^2}{\sigma_2^2} \right) \right\}$$

не представляет особой технической сложности. Эти оценки совпадают с оценками по методу моментов и, таким образом, равны выборочным аналогам тех характеристик двумерного нормального распределения, которые соответствуют указанным пяти параметрам:

$$\hat{\mu}_{1,n} = \bar{X}, \quad \hat{\mu}_{2,n} = \bar{Y}, \quad \hat{\sigma}_{1,n}^2 = S_X^2, \quad \hat{\sigma}_{2,n}^2 = S_Y^2, \quad \hat{\rho}_n = r.$$

Формулы для вычисления выборочных средних  $\bar{X}$  и  $\bar{Y}$ , выборочных дисперсий  $S_1^2$  и  $S_2^2$ , а также выборочного коэффициента корреляции  $r$  приведены в §2.

Полученные оценки часто используются для оценки параметров линейного прогноза  $Y = a + bX$  значений случайной величины  $Y$  по результатам наблюдений  $X$ . В случае нормального распределения линейный прогноз обладает свойством оптимальности с точки зрения малости средней квадратичной ошибки и совпадает с кривой средней квадратичной регрессии (см. предложение 10.3 курса ТВ)

$$y = \mu_2 + \rho \frac{\sigma_2}{\sigma_1} (x - \mu_1).$$

Однако формальная подгонка прогностической кривой с помощью прямой линии используется и вне рамок нормальной модели, и в этом случае оценки

$$\hat{a}_n = \bar{Y} - r \frac{S_2}{S_1} \bar{X}, \quad \hat{b}_n = r \frac{S_2}{S_1}$$

совпадают с оценками по *методу наименьших квадратов*: минимизируется, как и в примере 4.3, сумма квадратов невязок

$$\sum_1^n (Y_k - a - bX_k)^2.$$

Хотя оценки в обоих примерах имеют одинаковый вид, но решаемые в них статистические проблемы весьма различны: в примере 4.3 оценивались параметры некоторой функциональной зависимости с ошибками в наблюдениях отклика, в то время как в примере 4.4 решается задача выявления корреляционной связи и использования этой связи для прогноза.

## Лекция 7

Исследуем теперь асимптотические свойства оценок по методу максимального правдоподобия.

Начнем с выяснения достаточных условий состоятельности этих оценок. Такие ограничения на вероятностную модель обычно называются *условиями регулярности*, и в данном случае они имеют следующий вид.

- (R1) Параметрическое пространство  $\Theta$  есть открытый интервал на прямой  $\mathbb{R}$ .
- (R2) Носитель  $\mathcal{X}$  распределения  $P_\theta$  наблюдаемой случайной величины  $X$  не зависит от  $\theta$ , то есть все множества  $\mathcal{X} = \{x : f(x|\theta) > 0\}$  можно считать одинаковыми, каково бы ни было  $\theta \in \Theta$ .
- (R3) Распределения  $P_\theta$  различны при разных  $\theta$ , то есть при любых  $\theta_1 \neq \theta_2$ ,  $\theta_1, \theta_2 \in \Theta$ , имеет место тождество  $\mu\{x : x \in \mathcal{X}, f(x|\theta_1) = f(x|\theta_2)\} = 0$ , где  $\mu$  – мера, по которой вычисляется плотность  $f(x|\theta)$  распределения  $P_\theta$ .

Доказательство состоятельности оценок максимального правдоподобия, как и оценок по методу моментов, опирается на закон больших чисел, но

при этом используется следующее достаточно простое, но играющее большую роль в теории вероятностей, неравенство.

**Лемма 4.1.** (неравенство Йенсена) Пусть  $X$  – случайная величина с конечным математическим ожиданием. Если функция  $g(\cdot)$  дважды дифференцируема и выпукла ( $g'' > 0$ ) на некотором интервале, содержащем носитель распределения  $X$ , и математическое ожидание  $\mathbf{E}g(X)$  существует, то справедливо неравенство  $\mathbf{E}g(X) \geq g(\mathbf{E}X)$ , причем знак равенства достигается тогда и только тогда, когда распределение  $X$  сосредоточено в одной точке ( $X = \text{const.}$ ).

**Доказательство.** Так как функция  $g$  дважды дифференцируема, то справедливо следующее представление Тейлора в окрестности точки  $\mu = \mathbf{E}X$ :

$$g(X) = g(\mu) + (X - \mu)g'(\mu) + (X - \mu)^2 g''(\mu + \lambda(X - \mu))/2, \quad 0 < \lambda < 1.$$

Вычисляя математическое ожидание от обеих частей этого равенства, получаем

$$\mathbf{E}g(X) = g(\mathbf{E}X) + \mathbf{E}(X - \mu)^2 g''(\mu + \lambda(X - \mu))/2 \geq g(\mathbf{E}X).$$

Знак равенства возможен только в случае  $\mathbf{E}(X - \mu)^2 g''(\mu + \lambda(X - \mu)) = 0$ . Но поскольку  $g'' > 0$ , то последнее равенство с необходимостью влечет  $(X - \mu)^2 = 0$ , то есть  $X = \text{const.}$

Покажем теперь, что справедлива

**Теорема 4.1** (состоятельность). Если функция логарифмического правдоподобия

$$\mathcal{L}(\theta | X^{(n)}) = \sum_{k=1}^n \ln f(X_k | \theta) \quad (3)$$

имеет единственный максимум, то при выполнении условий регулярности (R1)–(R3) точка  $\hat{\theta}_n$  достижения максимума этой функции (оценка максимального правдоподобия) является состоятельной оценкой параметра  $\theta$ .

**Доказательство.** Покажем, что для любого фиксированного  $\theta_0 \in \Theta$  и любого  $\varepsilon > 0$  вероятность  $P_{\theta_0} \left( |\hat{\theta}_n - \theta_0| < \varepsilon \right) \rightarrow 1$ .

Если  $\theta_0$  – истинное значение параметра  $\theta$ , то в силу условия (R1)  $\theta_0$  – внутренняя точка  $\Theta$ . Тогда сформулированная выше задача состоит в доказательстве следующего утверждения: в некоторой  $\varepsilon$ -окрестности ( $\theta_0$  –

$\varepsilon; \theta_0 + \varepsilon)$  функция  $\mathcal{L}(\cdot | X^{(n)})$  обладает локальным максимумом с вероятностью, стремящейся к единице при  $n \rightarrow \infty$ .

Но если происходит событие

$$A_n = \left\{ \mathcal{L}(\theta_0 | X^{(n)}) > \mathcal{L}(\theta_0 \pm \varepsilon | X^{(n)}) \right\},$$

то внутри этой окрестности имеется точка максимума, и нам остается только показать, что  $P_{\theta_0}(A_n) \rightarrow 1$ , ибо  $P_{\theta_0}(|\hat{\theta}_n - \theta_0| < \varepsilon) \geq P_{\theta_0}(A_n)$ .

Используя условие (R2) и вид функции  $\mathcal{L}$  (см. (3)), представим неравенство, определяющее событие  $A_n$ , в виде

$$\frac{1}{n} \sum_{k=1}^n \ln \frac{f(X_k | \theta_0 \pm \varepsilon)}{f(X_k | \theta_0)} < 0.$$

В силу закона больших чисел Хинчина левая часть этого неравенства сходится по вероятности к

$$\mathbf{E}_{\theta_0} \ln \frac{f(X_k | \theta_0 \pm \varepsilon)}{f(X_k | \theta_0)}, \quad (4)$$

и для доказательства утверждения достаточно показать, что это математическое ожидание строго меньше нуля (кстати, докажите сами, что при справедливости условий теоремы математическое ожидание (4) всегда существует, в противном случае закон больших чисел Хинчина не применим).

Так как  $g(x) = -\ln x$  – выпуклая функция, то в силу неравенства Йенсена

$$\begin{aligned} \mathbf{E}_{\theta_0} \ln \frac{f(X | \theta_0 \pm \varepsilon)}{f(X | \theta_0)} &\leq \ln \mathbf{E}_{\theta_0} \frac{f(X | \theta_0 \pm \varepsilon)}{f(X | \theta_0)} = \\ &\ln \int_{\mathcal{X}} \frac{f(x | \theta_0 \pm \varepsilon)}{f(x | \theta_0)} \cdot f(x | \theta_0) d\mu(x) = \ln 1 = 0, \end{aligned}$$

причем равенство нулю первого члена в этой цепочке неравенств возможно лишь в случае

$$\frac{f(X | \theta_0 \pm \varepsilon)}{f(X | \theta_0)} = \text{const.},$$

то есть, поскольку интеграл от плотности равен 1, лишь в случае  $f(X | \theta_0 \pm \varepsilon) = f(X | \theta_0)$ , что невозможно в силу условия (R3). Таким образом, математическое ожидание (4) строго меньше нуля, и состоятельность оценки максимального правдоподобия доказана.

Изучим теперь асимптотическое распределение оценки максимального правдоподобия. Для этого нам потребуется ввести дополнительные условия регулярности.

(R4) Для каждой точки  $\theta_0$  параметрического пространства  $\Theta$  существует некоторая ее окрестность, в которой функция плотности  $f(x|\theta)$  трижды дифференцируема по параметру  $\theta$  и

$$\left| \frac{\partial f(x|\theta)}{\partial \theta} \right| \leq H_1(x), \quad (5)$$

$$\left| \frac{\partial^2 f(x|\theta)}{\partial \theta^2} \right| \leq H_2(x), \quad (6)$$

$$\left| \frac{\partial^3 \ln f(x|\theta)}{\partial \theta^3} \right| \leq H_3(x),$$

причем функции  $H_1$  и  $H_2$  интегрируемы по мере  $\mu$  на носителе  $\mathcal{X}$  распределения  $X$  и  $\mathbf{E}_\theta H_3(X) < \infty$  в некоторой окрестности каждой точки  $\theta$  параметрического пространства  $\Theta$ .

(R5) Функция

$$I(\theta) = \mathbf{E}_\theta \left( \frac{\partial \ln f(X|\theta)}{\partial \theta} \right)^2 = \int_{\mathcal{X}} \left( \frac{\partial \ln f(x|\theta)}{\partial \theta} \right)^2 f(x|\theta) d\mu(x) > 0,$$

каково бы ни было  $\theta \in \Theta$ .

Естественно, столь громоздкие и, на первый взгляд, странные условия требуют некоторого комментария.

Условие (R4) означает, что соответствующие производные функции плотности равномерно интегрируемы на  $\mathcal{X}$ , и поэтому можно выносить производную по  $\theta$  за знак интеграла.

Условие (R5) требует положительности очень важной, с точки зрения состоятельности статистического вывода, характеристики вероятностной модели:  $I(\theta)$  называется *информацией по Фишеру* в точке  $\theta$ , содержащейся в наблюдении случайной величины  $X$ . Если  $I(\theta) = 0$ , то возникают непреодолимые трудности с принятием корректного решения, соответствующего этой параметрической точке  $\theta$ . Понятно, что аналогичным образом можно

определить и информацию по Фишеру, содержащуюся в случайной выборке  $X^{(n)}$  :

$$I_n(\theta) = \mathbf{E}_\theta \left( \frac{\partial \ln f_n(X^{(n)} | \theta)}{\partial \theta} \right)^2.$$

Приведем несколько утверждений, касающихся свойств информации по Фишеру.

**Лемма 4.2.**  $1^0$ . При выполнении условия (R4) в части (6) для вычисления информации по Фишеру можно использовать формулу

$$I(\theta) = -\mathbf{E}_\theta \frac{\partial^2 \ln f(X | \theta)}{\partial \theta^2}.$$

$2^0$ . При выполнении условия (R4) в части (5) информация по Фишеру обладает свойством аддитивности  $I_n(\theta) = nI(\theta)$  – информация, содержащаяся в выборке, равна сумме информации, содержащихся в наблюдении каждой ее компоненты.

Доказательство..  $1^0$ . Условие (R4) в части (6) обеспечивает возможность смены порядка дифференцирования и интегрирования функции плотности, поэтому

$$\begin{aligned} \mathbf{E}_\theta \frac{\partial^2 \ln f(X | \theta)}{\partial \theta^2} &= \mathbf{E}_\theta \left( \frac{f''_{\theta\theta}(X | \theta)}{f(X | \theta)} - \left( \frac{f'_\theta(X | \theta)}{f(X | \theta)} \right)^2 \right) = \\ &= \int_{\mathcal{X}} \frac{f''_{\theta\theta}(x | \theta)}{f(x | \theta)} \cdot f(x | \theta) d\mu(x) - I(\theta) = \frac{d^2}{d\theta^2} \int_{\mathcal{X}} f(x | \theta) d\mu(x) - I(\theta) = -I(\theta). \end{aligned}$$

$2^0$ . Используя независимость и одинаковую распределенность компонент случайной выборки, получаем, что

$$\begin{aligned} I_n(\theta) &= \mathbf{E}_\theta \left( \frac{\partial \sum_{k=1}^n \ln f(X_k | \theta)}{\partial \theta} \right)^2 = \\ &= \mathbf{E}_\theta \left( \sum_{k=1}^n \left( \frac{\partial \ln f(X_k | \theta)}{\partial \theta} \right)^2 - \sum_{i \neq j} \frac{\partial \ln f(X_i | \theta)}{\partial \theta} \cdot \frac{\partial \ln f(X_j | \theta)}{\partial \theta} \right) = \\ &= \sum_{k=1}^n \mathbf{E}_\theta \left( \frac{\partial \ln f(X_k | \theta)}{\partial \theta} \right)^2 - \sum_{i \neq j} \mathbf{E}_\theta \frac{\partial \ln f(X_i | \theta)}{\partial \theta} \cdot \mathbf{E}_\theta \frac{\partial \ln f(X_j | \theta)}{\partial \theta} = nI(\theta), \end{aligned}$$

поскольку, в силу неравенства (5) в условии (R4), математическое ожидание

$$\mathbf{E}_\theta \frac{\partial \ln f(X | \theta)}{\partial \theta} = \int_{\mathcal{X}} \frac{f'_\theta(x | \theta)}{f(x | \theta)} \cdot f(x | \theta) d\mu(x) = \frac{d}{d\theta} \int_{\mathcal{X}} f(x | \theta) d\mu(x) = 0.$$

Теперь приступим к выводу асимптотического распределения оценки максимального правдоподобия скалярного параметра  $\theta$ .

**Теорема 4.2** (асимптотическая нормальность). *При выполнении условий (R1)–(R5) и наличии единственного локального максимума у функции правдоподобия корень  $\hat{\theta}_n = \hat{\theta}_n(X^{(n)})$  уравнения правдоподобия*

$$\partial \mathcal{L}(\theta | X^{(n)}) / \partial \theta = 0$$

асимптотически ( $n \rightarrow \infty$ ) нормален со средним  $\theta$  и дисперсией  $(nI(\theta))^{-1}$ , то есть

$$\lim_{n \rightarrow \infty} P_\theta \left( (\hat{\theta}_n - \theta) \sqrt{nI(\theta)} < x \right) = \Phi(x).$$

*Доказательство.* Если  $\hat{\theta}_n$  – оценка по методу максимального правдоподобия (корень уравнения правдоподобия), то имеет место тождество  $\partial \mathcal{L}(\hat{\theta}_n | X^{(n)}) / \partial \theta = 0$ . Используя условие (R4), разложим его левую часть по формуле Тейлора в окрестности истинного значения  $\theta_0$  параметра  $\theta$ :

$$\begin{aligned} \partial \mathcal{L}(\hat{\theta}_n | X^{(n)}) / \partial \theta &= \mathcal{L}'(\theta_0 | X^{(n)}) + \\ &(\hat{\theta}_n - \theta_0) \mathcal{L}''(\theta_0 | X^{(n)}) + (\hat{\theta}_n - \theta_0)^2 \mathcal{L}'''(\theta_1 | X^{(n)}) / 2 = 0, \end{aligned}$$

где производные от функции правдоподобия  $\mathcal{L}$  вычисляются по параметру  $\theta$ , а  $\theta_1 = \theta_0 + \lambda(\hat{\theta}_n - \theta_0)$ ,  $0 < \lambda < 1$ .

Разрешим полученное уравнение относительно величины  $\sqrt{n}(\hat{\theta}_n - \theta_0)$ , которая, согласно утверждению теоремы, должна иметь в пределе при  $n \rightarrow \infty$  нормальное распределение со средним 0 и дисперсией  $[I(\theta_0)]^{-1}$ :

$$\sqrt{n}(\hat{\theta}_n - \theta_0) = \frac{\mathcal{L}'(\theta_0 | X^{(n)}) / \sqrt{n}}{-\mathcal{L}''(\theta_0 | X^{(n)}) / n - (\hat{\theta}_n - \theta_0) \mathcal{L}'''(\theta_1 | X^{(n)}) / 2n}. \quad (7)$$

Числитель правой части этого представления

$$\frac{1}{\sqrt{n}} \mathcal{L}'(\theta_0 | X^{(n)}) = \frac{1}{\sqrt{n}} \sum_1^n \frac{\partial \ln f(X_k | \theta)}{\partial \theta}$$

есть нормированная на  $\sqrt{n}$  сумма независимых, одинаково распределенных случайных величин с нулевыми средними и дисперсиями  $I(\theta_0) > 0$  (см. доказательство пункта 2<sup>0</sup> леммы 4.2). Таким образом, в силу центральной предельной теоремы числитель правой части (7) асимптотически нормален с этими параметрами, и для завершения доказательства теоремы достаточно показать, что знаменатель (7) сходится по вероятности к постоянной  $I(\theta_0)$ , и сослаться на пункт (2) предложения 11.1 (теорема типа Слуцкого) курса ТВ.

В силу закона больших чисел и утверждения 1<sup>0</sup> леммы 4.2 первое слагаемое в знаменателе (7)

$$-\frac{1}{n}\mathcal{L}''(\theta_0 | X^{(n)}) = -\frac{1}{n} \sum_1^n \frac{\partial^2 \ln f(X_k | \theta_0)}{\partial \theta^2} \xrightarrow{P} -\mathbf{E}_{\theta_0} \frac{\partial^2 \ln f(X | \theta_0)}{\partial \theta^2} = I(\theta_0),$$

так что остается показать, что и второе слагаемое сходится по вероятности к нулю.

Так как при выполнении условий (R1)–(R3) оценка максимального правдоподобия состоятельна, то  $\hat{\theta}_n - \theta_0 \xrightarrow{P} 0$ . Множитель при этой разности

$$\frac{1}{n}\mathcal{L}'''(\theta_1 | X^{(n)}) = \frac{1}{n} \sum_1^n \frac{\partial^3 \ln f(X_k | \theta_1)}{\partial \theta^3}$$

в силу условия (R4), начиная с некоторого  $n$  по абсолютной величине не превосходит  $(1/n) \sum_1^n H_3(X_k)$  (это то  $n$ , при котором  $\theta_1$  попадает в окрестность точки  $\theta_0$  с любой наперед заданной вероятностью  $1 - \varepsilon$ ). Применяя к этой сумме закон больших чисел, получаем, что она сходится по вероятности к

$$\mathbf{E}_{\theta_0} H_3(X) < \infty,$$

и поэтому указанный выше сомножитель ограничен с вероятностью единица, а все второе слагаемое в знаменателе правой части (7) сходится по вероятности к нулю.

Доказанная теорема, как будет видно из основного результата следующего параграфа, устанавливает асимптотическую оптимальность оценок максимального правдоподобия с точки зрения квадратичного риска.



## §5. Эффективность оценок

Лекция 8

Обсуждая в начале нашего курса общую проблему статистического вывода, мы говорили о главной задаче математической статистики – построении решающих правил  $\delta_n = \delta_n(X^{(n)})$ , минимизирующих равномерно по всем  $\theta \in \Theta$  функцию риска  $R(\theta; \delta_n)$ . К сожалению, без дополнительных ограничений на класс решающих функций эта задача не разрешима. Действительно, рассмотрим проблему оценки параметра  $\theta$ , в которой пространство решений  $\mathcal{D}$  совпадает с параметрическим пространством  $\Theta$ , а решающая функция  $\delta_n = \hat{\theta}_n$  – оценка  $\theta$ . Возьмем в качестве оценки некоторую фиксированную точку  $\theta_0 \in \Theta$ , то есть при любом результате  $x^{(n)}$  статистического эксперимента будем принимать одно и то же решение  $d = \theta_0$ . Если функция потерь обладает тем естественным свойством, что  $L(\theta, \theta) = 0$ , каково бы ни было значение  $\theta \in \Theta$ , то риск такой оценки  $R(\theta; \theta_0) = L(\theta, \theta_0)$  при  $\theta = \theta_0$  равен нулю. Таким образом, если мы хотим построить оценку с равномерно минимальным риском в классе всевозможных оценок  $\theta$ , то мы должны найти оценку  $\theta_n^*$  с функцией риска  $R(\theta, \theta_n^*) \equiv 0$ , и понятно, что такой оценки не существует. Поэтому мы будем всегда при поиске оптимальных решений указывать класс оценок, в которых ищется оптимальное решение.

**Определение 5.1.** Оценка  $\theta_n^* = \theta_n^*(X^{(n)})$  называется *оптимальной* или *оценкой с равномерно минимальным риском* в классе  $\mathcal{K}$  оценок параметра  $\theta$ , если для любой оценки  $\hat{\theta}_n \in \mathcal{K}$  и каждого  $\theta \in \Theta$  имеет место неравенство  $R(\theta; \theta_n^*) \leq R(\theta; \hat{\theta}_n)$ .

Ниже предлагается метод нахождения оптимальных оценок скалярного параметра  $\theta$  при квадратичной функции потерь в классе несмещенных оценок:  $\mathbf{E}_\theta \hat{\theta}_n(X^{(n)}) = \theta$  при любом  $\theta \in \Theta \subseteq \mathbb{R}$ , но при дополнительных ограничениях на вероятностную модель и соответствующее семейство распределений оценки. Эти ограничения аналогичны тем условиям регулярности, которые мы накладывали на вероятностную модель при изучении асимптотических свойств оценок максимального правдоподобия. Мы покажем, что квадратичный риск любой несмещенной оценки, удовлетворяющей этим условиям, не может быть меньше  $[nI(\theta)]^{-1}$  – асимптотической дисперсии оценки максимального правдоподобия (см. теорема 4.2). Следовательно, метод максимального правдоподобия доставляет асимптотическое решение

проблемы оптимальной оценки. Более того, мы покажем, что при наличии достаточных статистик метод максимального правдоподобия может привести и к точному решению проблемы равномерной минимизации функции риска.

Сформулируем условия регулярности, при выполнении которых будет находиться нижняя (достижимая!) граница квадратичного риска оценки.

(B1) Носитель  $\mathcal{X}$  распределения  $P_\theta$  наблюдаемой случайной величины  $X$  не зависит от  $\theta \in \Theta$  (условие, совпадающее с (R2) в §4).

(B2) Информация по Фишеру  $I(\theta)$  строго положительна при любом  $\theta \in \Theta$  (условие, совпадающее с (R5) в §4).

(B3) Равенство

$$\int_{\mathcal{X}^n} f_n(x^{(n)} | \theta) d\mu_n(x^{(n)}) = 1$$

можно дифференцировать по  $\theta$  под знаком интеграла, то есть

$$\int_{\mathcal{X}^n} f'_n(x^{(n)} | \theta) d\mu_n(x^{(n)}) = 0.$$

По аналогии с (R4) в части (5) для этого достаточно потребовать существование такой интегрируемой по мере  $\mu$  функции  $H(x)$ , что в некоторой окрестности любой точки  $\theta \in \Theta$  выполняется неравенство  $|\partial f(x | \theta) / \partial \theta| \leq H(x)$ ,  $x \in \mathcal{X}$ .

(B4) Оценка  $\hat{\theta}_n = \hat{\theta}_n(X^{(n)})$  должна принадлежать классу оценок  $\mathcal{K}'$ , среднее значение которых

$$\mathbf{E}_\theta \hat{\theta}_n(X^{(n)}) = \int_{\mathcal{X}^n} \hat{\theta}_n(x^{(n)}) f_n(x^{(n)} | \theta) d\mu_n(x^{(n)})$$

можно дифференцировать по  $\theta \in \Theta$  под знаком интеграла.

Конечно, условие (B4) требует комментария. В “высокой” теории статистического вывода приводятся достаточные условия на семейство распределений  $\{P_\theta, \theta \in \Theta\}$  наблюдаемой случайной величины  $X$ , которые обеспечивают выполнение условия (B4), но формулировка этих условий и, в особенности, доказательство того, что они влекут (B4), настолько технически и концептуально сложны, что могут составить предмет специального курса. Однако все изучаемые нами в курсе ТВ вероятностные модели, за

исключением равномерного распределения, удовлетворяют этим условиям, и поэтому любая оценка их параметров принадлежит классу  $\mathcal{K}'$ .

Прежде, чем получить основной “технический” результат этого параграфа, вспомним одно замечательное неравенство из курса математического анализа. Это – неравенство Коши–Буняковского, которое в случае интегралов Лебега по вероятностной мере  $P$  называется неравенством Шварца. Пусть  $Y$  – случайная величина с распределением  $P$  и  $g, h$  – две интегрируемые с квадратом по мере  $P$  функции на области  $\mathcal{Y}$  значений  $Y$ . Для этих функций имеет место неравенство  $(\mathbf{E} g(Y)h(Y))^2 \leq \mathbf{E} g^2(Y) \cdot \mathbf{E} h^2(Y)$  или, что то же,

$$\left( \int_{\mathcal{Y}} g(y)h(y) dP(y) \right)^2 \leq \int_{\mathcal{Y}} g^2(y)dP(y) \cdot \int_{\mathcal{Y}} h^2(y)dP(y),$$

причем знак равенства достигается тогда и только тогда, когда функции  $g$  и  $h$  линейно зависимы: существуют такие постоянные  $a$  и  $b$ , что  $ag(y) + bh(y) = 0$  для почти всех  $y \in \mathcal{Y}$  по мере  $P$ .

**Теорема 5.1.** (неравенство Рао–Крамера) При выполнении условий (B1)–(B4) для квадратичного риска любой оценки  $\hat{\theta}_n \in \mathcal{K}'$  справедливо неравенство

$$\mathbf{E}_{\theta} \left( \hat{\theta}_n(X^{(n)}) - \theta \right)^2 \geq \mathbf{D}_{\theta} \hat{\theta}_n(X^{(n)}) \geq \frac{[d\gamma(\theta)/d\theta]^2}{nI(\theta)}, \quad (1)$$

где  $\gamma(\theta) = \mathbf{E}_{\theta} \hat{\theta}_n(X^{(n)})$ , причем знак равенства между риском и дисперсией оценки  $\hat{\theta}_n$  достигается на несмещенных оценках:  $\gamma(\theta) = \theta$ , а знак равенства во втором неравенстве (1) имеет место тогда и только тогда, когда существует такая параметрическая функция  $C(\theta)$ ,  $\theta \in \Theta$ , что

$$\hat{\theta}_n(X^{(n)}) - \gamma(\theta) = C(\theta) \frac{\partial \mathcal{L}(\theta | X^{(n)})}{\partial \theta} \quad (2)$$

почти наверное по мере  $P_{\theta}$ .

**Доказательство.** Продифференцируем обе части равенств

$$\int_{\mathcal{X}^n} f_n(x^{(n)} | \theta) d\mu_n(x^{(n)}) = 1,$$

$$\int_{\mathcal{X}^n} \hat{\theta}_n(x^{(n)}) f_n(x^{(n)} | \theta) d\mu_n(x^{(n)}) = \gamma(\theta)$$

по параметру  $\theta$ , занося производные в левых частях под знаки интегралов, что можно сделать благодаря условиям (В3) и (В4). Полученный результат, используя условие (В1), представим в виде

$$\int_{\mathcal{X}^n} \frac{\partial \mathcal{L}(\theta | x^{(n)})}{\partial \theta} f_n(x^{(n)} | \theta) d\mu_n(x^{(n)}) = 0,$$

$$\int_{\mathcal{X}^n} \hat{\theta}_n(x^{(n)}) \frac{\partial \mathcal{L}(\theta | x^{(n)})}{\partial \theta} f_n(x^{(n)} | \theta) d\mu_n(x^{(n)}) = \gamma'(\theta).$$

Вычтем из второго равенства первое, умножив его предварительно на  $\gamma(\theta)$  :

$$\int_{\mathcal{X}^n} \left( \hat{\theta}_n(x^{(n)}) - \gamma(\theta) \right) \frac{\partial \mathcal{L}(\theta | x^{(n)})}{\partial \theta} f_n(x^{(n)} | \theta) d\mu_n(x^{(n)}) = \gamma'(\theta).$$

Применим к левой части полученного равенства неравенство Шварца, полагая  $y = x^{(n)}$ ,  $\mathcal{Y} = \mathcal{X}^n$ ,  $g(x^{(n)}) = \hat{\theta}_n(x^{(n)}) - \gamma(\theta)$ ,  $h(x^{(n)}) = \partial \mathcal{L}(\theta | x^{(n)}) / \partial \theta$ ,  $dP(y) = f_n(x^{(n)} | \theta) d\mu_n(x^{(n)})$ . В результате получим неравенство

$$(\gamma'(\theta))^2 \leq \int_{\mathcal{X}^n} \left( \hat{\theta}_n(x^{(n)}) - \gamma(\theta) \right)^2 f_n(x^{(n)} | \theta) d\mu_n(x^{(n)}).$$

$$\int_{\mathcal{X}^n} \left( \frac{\partial \mathcal{L}(\theta | x^{(n)})}{\partial \theta} \right)^2 f_n(x^{(n)} | \theta) d\mu_n(x^{(n)}), \quad (3)$$

в котором знак равенства достигается тогда и только тогда, когда выполняется соотношение (2).

Мы получили неравенства (1), поскольку первое из них очевидно (на дисперсии достигается минимум всевозможных средних квадратичных отклонений случайной величины от постоянной). Второе неравенство в (1) есть следствие неравенства (3), ибо первый интеграл в правой части (3) равен  $\mathbf{D}_\theta \hat{\theta}_n$ , а второй интеграл определяет фишеровскую информацию  $I_n(\theta)$ , содержащуюся в выборке. Наконец, из пункта 2<sup>0</sup> леммы 4.2 следует, что  $I_n(\theta) = nI(\theta)$ .

**Следствие 5.1.** *Если  $\hat{\theta}_n$  принадлежит подклассу  $\mathcal{K} \subseteq \mathcal{K}'$  несмещенных оценок класса  $\mathcal{K}'$ , то ее квадратичный риск*

$$R(\theta; \hat{\theta}_n) = \mathbf{D}_\theta \hat{\theta}_n \geq [nI(\theta)]^{-1}, \quad (4)$$

причем знак равенства тогда и только тогда, когда выполняется равенство (2) с  $\gamma(\theta) = \theta$ .

Понятно, что это следствие есть частный случай доказанной теоремы. Оно указывает неконструктивный путь к построению несмещенных оценок с равномерно минимальным риском. Достаточно вычислить производную в правой части равенства (2) и затем подбирать статистику  $\hat{\theta}_n = \hat{\theta}_n(X^{(n)})$  и параметрическую функцию  $C(\theta)$ , для которых имеет место равенство

$$\hat{\theta}_n(X^{(n)}) - \theta = C(\theta) \frac{\partial \mathcal{L}(\theta | X^{(n)})}{\partial \theta}.$$

Обычно это можно сделать в случае статистических структур, обладающих достаточными статистиками, где, в силу теоремы факторизации (теорема 2.1 из §2), функция правдоподобия  $L(\theta | X^{(n)}) = g_\theta(T(X^{(n)}))h(X^{(n)})$ , и последнее равенство имеет вид

$$\hat{\theta}_n(X^{(n)}) - \theta = C(\theta) \frac{\partial \ln g_\theta(T(X^{(n)}))}{\partial \theta}. \quad (5)$$

Например, для показательного распределения с функцией плотности  $f(x | \theta) = \theta^{-1} \exp\{-x/\theta\}$ ,  $x > 0$ , функция

$$\ln g_\theta(X^{(n)}) = -n \ln \theta - \theta^{-1} \sum_1^n X,$$

ее производная

$$\partial \ln g_\theta(T(X^{(n)}))/\partial \theta = -n/\theta + \sum_1^n X/\theta^2,$$

и равенство (5) выполняется при  $C(\theta) = \theta^2/n$  и  $\hat{\theta}_n = \bar{X}$ . Таким образом, выборочное среднее  $\bar{X}$  есть несмещенная оценка с равномерно минимальным риском для параметра  $\theta$  показательного распределения. Напомним, что  $\bar{X}$  оценка  $\theta$  как по методу моментов, так и по методу максимального правдоподобия.

Легко понять, что если в (4) достигается знак равенства, то  $\hat{\theta}_n$  — оптимальная оценка в классе  $\mathcal{K}$ , но обратное, вообще говоря, может и не выполняться — мы не располагаем утверждением, что любая оптимальная оценка имеет квадратичный риск, равный  $[n I(\theta)]^{-1}$ . Чтобы подчеркнуть это различие и указать в дальнейшем более конструктивный метод построения оптимальных оценок, введем еще одно определение, рассмотрев более общую задачу несмещенной оценки некоторой параметрической функции  $\gamma(\theta)$ .

**Определение 5.2.** Несмещенная оценка  $\hat{\gamma}_n = \hat{\gamma}_n(X^{(n)})$  параметрической функции  $\gamma(\theta)$  называется *эффективной* в классе  $\mathcal{K}'$ , если ее квадратичный риск

$$R(\gamma; \hat{\gamma}_n) = \mathbf{E}_\theta \left( \hat{\gamma}_n(X^{(n)}) - \gamma(\theta) \right)^2 = \mathbf{D}_\theta \hat{\gamma}_n(X^{(n)}) = [\gamma'(\theta)]^2/n I(\theta),$$

то есть (см. теорему 5.1 с  $\hat{\theta}_n = \hat{\gamma}_n$ ) выполняется равенство

$$\hat{\gamma}_n(X^{(n)}) - \gamma(\theta) = C(\theta) \frac{\partial \mathcal{L}(\theta | X^{(n)})}{\partial \theta}. \quad (6)$$

Оценка  $\hat{\gamma}_n$  называется *асимптотически эффективной* в классе  $\mathcal{K}'$ , если  $\mathbf{E}_\theta \hat{\gamma}_n(X^{(n)}) \sim \gamma(\theta)$  и  $\mathbf{D}_\theta \hat{\gamma}_n(X^{(n)}) \sim [\gamma'(\theta)]^2/n I(\theta)$ , когда  $n \rightarrow \infty$ .

В силу теоремы 4.2 оценка по методу максимального правдоподобия скалярного параметра  $\theta$  (в данном случае  $\gamma(\theta) = \theta$ ) является асимптотически эффективной оценкой в классе  $\mathcal{K}'$ . Покажем, что она дает решение проблемы построения эффективной оценки в классе  $\mathcal{K}$ .

Пусть  $\hat{\theta}_n$  – оценка максимального правдоподобия параметра  $\theta$ . Определим оценку  $\gamma(\hat{\theta}_n)$  параметрической функции  $\gamma(\theta)$  с помощью подстановки вместо  $\theta$  ее оценки  $\hat{\theta}_n$ .

**Теорема 5.2.** Если  $\gamma(\hat{\theta}_n)$  есть несмещенная оценка параметрической функции  $\gamma(\theta)$  и эффективная в классе  $\mathcal{K}'$  оценка  $\gamma_n^*$  параметрической функции  $\gamma(\theta)$  существует, то при выполнении условий регулярности (R1)–(R5) и (B3)–(B4) почти наверное  $\gamma_n^*(X^{(n)}) = \gamma(\hat{\theta}_n(X^{(n)}))$ .

*Доказательство.* Если  $\gamma_n^*$  – эффективная оценка  $\gamma(\theta)$ , то она удовлетворяет равенству (6):

$$\gamma_n^*(X^{(n)}) - \gamma(\theta) = C(\theta) \partial \mathcal{L}(\theta | X^{(n)}) / \partial \theta, \quad (7)$$

каково бы ни было  $\theta \in \Theta$ . Но если  $\hat{\theta}_n$  – оценка по методу максимального правдоподобия, то  $\partial \mathcal{L}(\hat{\theta}_n | X^{(n)}) / \partial \theta = 0$ , так что равенство (7) при  $\theta = \hat{\theta}_n$  превращается в равенство  $\gamma_n^*(X^{(n)}) - \gamma(\hat{\theta}_n) = 0$  почти наверное по вероятности  $P_\theta^n$ .

Из доказанной теоремы немедленно вытекает, что выборочное среднее  $\bar{X}$  есть эффективная (следовательно, и оптимальная) несмещенная оценка параметра  $\theta$  таких распределений, как двухточечное, биномиальное при известном  $m$ , Пуассона, показательное;  $\bar{X}$  есть также несмещенная оценка с равномерно минимальным квадратичным риском среднего значения  $\mu$  нормального  $(\mu, \sigma^2)$  распределения.

## §6. Доверительные интервалы

Лекция 9

Мы рассмотрели несколько методов построения *точечных* оценок для параметров, значения которых определяют распределение наблюдаемой случайной величины. Был получен ряд утверждений о распределении таких оценок, что позволяет судить о надежности оценки при заданной точности, то есть вычислять вероятности событий вида  $|\hat{\theta}_n(X^{(n)}) - \theta| \leq \Delta$  при каждом фиксированном значении параметра  $\theta$ . Поскольку именно значение  $\theta$  нам неизвестно, то такого рода вычисления зачастую лишены практического смысла – слишком велик размах в надежности оценки  $\hat{\theta}_n$  при различных  $\theta$ , даже в случае, когда мы располагаем некоторой априорной информацией о возможной области значений этого параметра. Поэтому в ряде практических ситуаций пытаются решать обратную задачу: для фиксированной надежности, скажем,  $1 - \alpha$ , где  $\alpha$  мало, указать некоторую область значений  $\theta$ , зависящую, естественно, от выборки  $X^{(n)}$ , которая с вероятностью, не меньшей  $1 - \alpha$ , накрывает истинное, неизвестное нам значение  $\theta$ , причем такое надежностное утверждение должно выполняться при любых  $\theta \in \Theta$ . В таком случае по размерам области, которые определяются выборочными значениями  $x^{(n)}$ , можно судить о точности такой *интервальной* оценки.

**Определение 6.1.** Подмножество  $\Delta_n = \Delta_n(X^{(n)})$  параметрического пространства  $\Theta$  называется  $(1 - \alpha)$ -*доверительной областью*, если

$$P_\theta \left( \Delta_n(X^{(n)}) \ni \theta \right) \geq 1 - \alpha, \quad (1)$$

каково бы ни было значение  $\theta \in \Theta$ . Заданное (фиксированное) значение  $1 - \alpha$  называется *доверительным уровнем*, а наименьшее значение левой части неравенства (1) по всем  $\theta \in \Theta$  – *доверительным коэффициентом*. В случае  $\Theta \subseteq \mathbb{R}$  доверительная область вида  $\Delta_n = (\underline{\theta}_n(X^{(n)}); \bar{\theta}_n(X^{(n)}))$  называется *доверительным интервалом*, в котором различаются *нижний*  $\underline{\theta}_n$  и *верхний*  $\bar{\theta}_n$  *доверительные пределы*. Доверительные интервалы вида  $(\underline{\theta}_n; \infty)$  и  $(-\infty; \bar{\theta}_n)$  называются соответственно *нижней* и *верхней доверительными границами*.

Естественно, конфигурация доверительной области выбирается статистиком, сообразуясь с ее геометрической наглядностью и, главное, возможностью гарантировать доверительную вероятность. В случае скалярного

параметра доверительная область обычно выбирается в виде интервала, причем в ряде случаев, например, при оценке надежности или вероятности нежелательного события, в виде одностороннего интервала. В случае многомерного параметра обычно строятся доверительные эллипсоиды или параллелепипеды.

Следует обратить особое внимание на правильную формулировку доверительного утверждения, которая подчеркивается в неравенстве (1) записью  $\Delta_n(X^{(n)}) \ni \theta$  вместо обычного  $\theta \in \Delta_n(X^{(n)})$ . Говорить, что значение параметра  $\theta$  с вероятностью, не меньшей  $1 - \alpha$ , принадлежит области  $\Delta_n$ , значит сознательно вводить трудящихся на ниве прикладной статистики в заблуждение. Дело в том, что значение параметра  $\theta$  в данной вероятностной модели не является случайной величиной, это постоянная, свойственная исследуемому объекту, а постоянная принадлежит какой-либо области только с вероятностью единица или ноль. Вся случайность заключена в самой доверительной области  $\Delta_n(X^{(n)})$ , и поэтому правильное доверительное утверждение гласит: *область  $\Delta_n(X^{(n)})$  с вероятностью, не меньшей  $1 - \alpha$ , накрывает истинное (неизвестное) значение  $\theta$ .*

В самом начале нашего курса математической статистики в примере 1.1 с определением содержания общей серы в дизельном топливе мы строили доверительный интервал фиксированной ширины для среднего значения нормального распределения, когда занимались планированием объема испытаний, необходимого для достижения заданной точности и надежности оценки. Рассмотрим еще раз этот пример в свете введенных понятий интервальной оценки параметра.

1<sup>0</sup>. ДОВЕРИТЕЛЬНЫЙ ИНТЕРВАЛ ДЛЯ СРЕДНЕГО ЗНАЧЕНИЯ НОРМАЛЬНОГО РАСПРЕДЕЛЕНИЯ ПРИ ИЗВЕСТНОЙ ДИСПЕРСИИ. Итак, в примере 1.1 мы имели дело с выборкой  $X^{(n)}$  из нормального  $(\mu, \sigma^2)$  распределения, причем значение параметра  $\sigma$  нам было известно, так что в качестве неизвестного параметра  $\theta$  выступало  $\mu$ . Наша задача состоит в построении такого интервала  $(\underline{\mu}_n(X^{(n)}), \bar{\mu}_n(X^{(n)}))$ , что  $P_\mu(\underline{\mu}_n \leq \mu \leq \bar{\mu}_n) \geq 1 - \alpha$ , при любом значении  $\mu \in \mathbb{R}$ .

Напомним, что в этом примере оценкой  $\mu$  служило выборочное среднее  $\bar{X}$  – несмещенная оценка  $\mu$  с минимальным квадратичным риском. Эта линейная оценка обладает замечательным свойством инвариантности: распределение разности  $\bar{X} - \mu$  не зависит от  $\mu$ , и это обстоятельство подска-



зывает нам путь к построению доверительного интервала. Действительно,

$$P\left(\frac{|\bar{X} - \mu|}{\sigma}\sqrt{n} \leq \lambda\right) = 2\Phi(\lambda) - 1,$$

и если положить  $\lambda$  равным корню уравнения  $2\Phi(\lambda) - 1 = 1 - \alpha$ , то есть выбрать  $\lambda = \lambda_\alpha = \Phi^{-1}(1 - \alpha/2)$ , то интервал  $(\bar{X} - \lambda_\alpha\sigma/\sqrt{n}, \bar{X} + \lambda_\alpha\sigma/\sqrt{n})$  будет  $(1 - \alpha)$ -доверительным интервалом для среднего значения  $\mu$  нормального  $(\mu, \sigma^2)$  распределения при известной дисперсии  $\sigma^2$ .

В этом простейшем примере на построение доверительного интервала ключевым моментом было использование *инвариантной* случайной функции  $\hat{\theta}_n - \theta$  от оценки  $\hat{\theta}_n = \bar{X}$  и параметра  $\theta = \mu$ . В принципе, именно на подобном выборе *опорной* функции  $H(\hat{\theta}_n, \theta)$  с подходящей оценкой  $\hat{\theta}_n$  параметра  $\theta$  основаны исторически первые методы построения доверительных интервалов и множеств. Опорная функция  $H(\cdot, \cdot)$  подбирается таким образом, чтобы она была монотонно возрастающей функцией второго аргумента  $\theta$ , и при этом вероятность  $P_\theta\left(H(\hat{\theta}_n(X^{(n)}), \theta) \leq \lambda\right)$  для некоторых значений  $\lambda$  должна оставаться достаточно высокой (близкой к единице), каково бы ни было значение  $\theta \in \Theta$ . Мы проиллюстрируем этот метод построения доверительных интервалов с помощью подбора инвариантных опорных функций на примере нормального  $(\mu, \sigma^2)$  распределения, строя доверительные интервалы для каждого из параметров при известном и неизвестном значениях другого (“мешающего”) параметра.

**2<sup>0</sup>. ВЕРХНЯЯ ДОВЕРИТЕЛЬНАЯ ГРАНИЦА ДЛЯ ДИСПЕРСИИ НОРМАЛЬНОГО РАСПРЕДЕЛЕНИЯ ПРИ ИЗВЕСТНОМ СРЕДНЕМ.** При выборе из нормального  $(\mu, \sigma^2)$  распределения с известным средним значением  $\mu$  метод максимального правдоподобия приводит к несмещенной оценке  $\hat{\sigma}_n^2 = n^{-1} \sum_1^n (X_k - \mu)^2$  параметра  $\sigma^2$ . Используя результаты предыдущего параграфа, нетрудно показать, что  $\hat{\sigma}_n^2$  есть несмещенная оценка с равномерно минимальным риском.

Поскольку, в чем мы неоднократно убеждались,  $Y_k = (X_k - \mu)/\sigma \sim \mathcal{N}(0, 1)$ ,  $k = 1, \dots, n$ , то естественно рассмотреть в качестве опорной функцию  $H(\hat{\sigma}_n^2, \sigma^2) = \sigma^{-2} \sum_1^n (X_k - \mu)^2$ . Найдем ее распределение.

**Лемма 6.1.** Если  $Y_1, \dots, Y_n$  независимы и одинаково распределены по стандартному нормальному закону  $\mathcal{N}(0, 1)$ , то  $\sum_1^n Y_k^2$  имеет гамма-распределение  $G(n/2, 2)$ .

Доказательство.. Покажем, что  $Y^2$ , где  $Y \sim \mathcal{N}(0, 1)$ , имеет гамма-распределение  $G(1/2, 2)$ , после чего просто воспользуемся теоремой сложения для гамма-распределения (см. предложение 12.2, пункт 5<sup>0</sup> курса ТВ). Функция распределения  $Y^2$  вычисляется по формуле  $F(x) = P(Y^2 < x) = P(-\sqrt{x} < Y < \sqrt{x}) = \Phi(\sqrt{x}) - \Phi(-\sqrt{x}) = 2\Phi(\sqrt{x}) - 1$ , так что ее функция плотности

$$f(x) = \frac{d}{dx} \left[ \frac{2}{\sqrt{2\pi}} \int_{-\infty}^{\sqrt{x}} \exp \left\{ -\frac{t^2}{2} \right\} - 1 \right] = \frac{1}{2^{1/2}\Gamma(1/2)} x^{1/2-1} e^{-x/2},$$

поскольку  $\Gamma(1/2) = \sqrt{\pi}$ . Мы видим, что это – функция плотности гамма-распределения  $G(1/2, 2)$  с параметром формы  $\lambda = 1/2$  и параметром масштаба  $a = 2$ , откуда, как было замечено выше, немедленно следует утверждение леммы.

Гамма-распределение  $G(n/2, 2)$  очень часто используется в различных задачах математической статистики, и оно появилось раньше, чем гамма-распределение  $G(\lambda, a)$  общего вида, под названием *хи-квадрат распределение с  $n$  степенями свободы*. Функция распределения хи-квадрат обычно обозначается  $K_n(x)$ ,  $x > 0$ , а что касается термина “степени свободы”, то его смысл прояснится по мере других применений хи-квадрат распределения.

Теперь мы можем перейти к нашей основной задаче – построению доверительных границ для  $\sigma^2$ . Если обратиться к практической стороне этой проблемы (см. в связи с этим пример 1.1), то легко понять, что статистика должна интересоваться только верхней (а не двусторонняя) граница  $\sigma^2$ , на которую он будет ориентироваться, чтобы обезопасить себя от грубых ошибок при планировании статистического эксперимента. Таким образом, мы должны сформулировать доверительное утверждение в форме  $\sigma^2 \leq \bar{\sigma}_n^2$ . Понятно, что нижняя доверительная граница и двусторонние границы (доверительный интервал), коль скоро они кому-то потребуются, строятся аналогичным образом.

В рамках такой формулировки задачи, мы должны рассмотреть событие

$$A_\lambda = \left\{ H(\hat{\sigma}_n^2, \sigma^2) = \frac{\sum_1^n (X_k - \mu)^2}{\sigma^2} \geq \lambda \right\},$$

выбирая  $\lambda$  из условия  $P_{\mu,\sigma}(A_\lambda) = 1 - \alpha$ . Как мы только что выяснили, эта

вероятность не зависит от  $\mu$  и  $\sigma$ , и в силу леммы 6.1 постоянная  $\lambda$  определяется квантилью хи-квадрат распределения с  $n$  степенями свободы – корнем уравнения  $1 - K_n(\lambda) = 1 - \alpha$ . Следовательно, верхняя  $(1 - \alpha)$ -доверительная граница для  $\sigma^2$  равна  $\sum_1^n (X_k - \mu)^2 / K_n^{-1}(\alpha)$ , где, в соответствии с нашими стандартными обозначениями,  $K_n^{-1}(\alpha)$  есть  $\alpha$ -квантиль хи-квадрат распределения с  $n$  степенями свободы.

Используемые в рассмотренных примерах методы подбора опорных функций, основанные на принципе инвариантности статистик (оценок параметров  $\mu$  и  $\sigma^2$ ) относительно линейных преобразований, позволяют аналогичным образом подбирать такие функции и в случае неизвестных значений мешающего параметра. Так, если рассматривается задача построения доверительных границ для  $\sigma^2$  при неизвестном  $\mu$ , то естественно обратиться к оценке  $S^2 = n^{-1} \sum_1^n (X_k - \bar{X})^2$  параметра  $\sigma^2$ , замечая, что ее распределение не зависит от  $\mu$ , поскольку каждая из разностей  $X_k - \bar{X}$  инвариантна относительно сдвига, когда  $X_k$  заменяется на  $X_k - \mu$ ,  $k = 1, \dots, n$ . Если разделить эти разности на  $\sigma$ , то мы получим случайные величины, распределение которых не зависит как от  $\mu$ , так и от  $\sigma$ , и таким образом мы приходим к инвариантной опорной функции  $H(S^2, \sigma^2) = S^2 / \sigma^2$ . Для вывода распределения этой функции можно обратиться к нормальным  $(0, 1)$  случайным величинам  $Y_k = (X_k - \mu) / \sigma$ ,  $k = 1, \dots, n$ , поскольку, в чем легко убедиться,  $H(S^2, \sigma^2) = n^{-1} \sum_1^n (Y_k - \bar{Y})^2$ .

Если обратиться к задаче доверительной интервальной оценки  $\mu$  при неизвестном значении  $\sigma$ , то здесь инвариантную опорную функцию можно построить, комбинируя ее из опорных функций задач 1<sup>0</sup> и 2<sup>0</sup>. Как мы видели при решении этих задач, распределения случайных величин  $(\bar{X} - \mu) / \sigma$  и  $S / \sigma$  не зависят от  $\mu$  и  $\sigma$ , и поэтому в качестве опорной функции при интервальной оценке  $\mu$  можно использовать опорную функцию, определяемую отношением этих величин, то есть функцию  $|\bar{X} - \mu| / S$ .

Однако для построения доверительных интервалов на основе таких функций нам необходимо найти совместное распределение статистик  $\bar{X}$  и  $S^2$ . Мы получим это распределение в следующей лекции, сформулировав его в виде утверждения, известного в математической статистике как *лемма Фишера*.

## Лекция 10

**Теорема 6.1.** *В случае выбора из нормального  $(\mu, \sigma^2)$  распределения*

статистики  $\bar{X}$  и  $S^2$  независимы,  $\bar{X} \sim \mathcal{N}(\mu, \sigma^2/n)$ , а  $nS^2/\sigma^2 \sim \chi_{n-1}^2$  (имеет хи-квадрат распределение с  $n - 1$  степенью свободы).

Доказательство.. Пусть  $Y_1, \dots, Y_n$  – случайная выборка из стандартного нормального распределения  $\mathcal{N}(0, 1)$ . Покажем, что статистики

$$\bar{Y} = \frac{1}{n} \sum_{k=1}^n Y_k \quad \text{и} \quad S_Y = \sqrt{\sum_{k=1}^n (Y_k - \bar{Y})^2}$$

независимы,  $\bar{Y} \sim \mathcal{N}(0, 1/n)$ , а  $S_Y \sim \chi_{n-1}^2$ . Тогда утверждение теоремы будет следовать из того факта, что  $\sigma Y_k + \mu$  имеют то же распределение, что и  $X_k$ ,  $k = 1, \dots, n$ , и, следовательно, распределение  $\bar{X}$  совпадает с распределением  $\sigma \bar{Y} + \mu$ , а распределение  $S_Y$  – с распределением  $nS^2/\sigma^2$ .

Введем случайные величины

$$Z_k = \sum_{i=1}^n c_{ki} Y_i, \quad k = 1, \dots, n,$$

которые определяются заданием матрицы  $C = \|c_{ki}\|$  линейных преобразований случайных величин  $Y_1, \dots, Y_n$ . Пусть элементы первой строки этой матрицы  $c_{11} = \dots = c_{1n} = 1/\sqrt{n}$ , а остальные элементы матрицы  $C$  выберем так, чтобы произведение  $C$  на транспонированную матрицу  $C'$  было единичной матрицей:  $CC' = I$ . Как известно, такой выбор  $C$  возможен, и полученная таким образом матрица называется ортонормированной. Случайные величины  $Z_1, \dots, Z_n$  распределены в соответствии с  $n$ -мерным нормальным законом, для спецификации которого достаточно найти вектор средних значений этих величин и матрицу их ковариаций.

Средние значения

$$m_k = \mathbf{E}Z_k = \mathbf{E} \sum_{i=1}^n c_{ki} Y_i = \sum_{i=1}^n c_{ki} \mathbf{E}Y_i = 0, \quad k = 1, \dots, n.$$

Далее, поскольку средние значения равны нулю, ковариации этих случайных величин

$$\begin{aligned} \text{cov}(Z_k, Z_j) &= \mathbf{E}(Z_k - m_k)(Z_j - m_j) = \mathbf{E}Z_k Z_j = \mathbf{E} \sum_{i=1}^n c_{ki} Y_i \cdot \sum_{i=1}^n c_{ji} Y_i = \\ &= \mathbf{E} \sum_{i=1}^n c_{ki} c_{ji} Y_i^2 + \mathbf{E} \sum_{i \neq l}^n c_{ki} c_{jl} Y_i Y_l. \end{aligned}$$

Если занести математические ожидания под знаки сумм и вспомнить, что  $Y_1, \dots, Y_n$  независимы,  $\mathbf{E}Y_i = 0$ ,  $\mathbf{E}Y_i^2 = \mathbf{D}Y_i = 1$ , а при  $i \neq l$  средние значения  $\mathbf{E}Y_i Y_l = \mathbf{E}Y_i \mathbf{E}Y_l = 0$ , то получим, что

$$\text{cov}(Z_k, Z_j) = \sum_{i=1}^n c_{ki} c_{ji}, \quad k, j = 1, \dots, n.$$

Поскольку для ортонормированной матрицы последняя сумма равна нулю, если  $k \neq j$ , и равна единице, если  $k = j$ , то мы приходим к заключению, что  $Z_1, \dots, Z_n$  независимы и одинаково распределены по стандартному нормальному закону  $\mathcal{N}(0, 1)$ . Таким образом, ортонормированные преобразования случайных величин  $Y_1, \dots, Y_n$  не изменили их совместное распределение.

Теперь представим наши статистики  $\bar{Y}$  и  $S_Y$  в терминах случайных величин  $Z_1, \dots, Z_n$ . Поскольку

$$Z_1 = \sum_{i=1}^n c_{1i} Y_i = \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i,$$

то  $\bar{Y} = Z_1/\sqrt{n}$ . Далее, ортонормированное линейное преобразование сохраняет сумму квадратов компонент преобразуемого вектора, то есть  $\sum_1^n Z_k^2 = \sum_1^n Y_k^2$ . Следовательно, статистика  $S_Y$  в новых переменных приобретает вид

$$\frac{S_Y}{n} = \frac{1}{n} \sum_1^n Y_k^2 - \bar{Y}^2 = \frac{1}{n} \sum_1^n Z_k^2 - \frac{Z_1^2}{n} = \frac{1}{n} \sum_2^n Z_k^2.$$

Итак, распределение  $\bar{Y}$  совпадает с распределением  $Z_1/\sqrt{n}$ , а распределение  $S_Y$  – с распределением суммы квадратов  $n - 1$  независимых в совокупности и независящих от  $Z_1$  нормальных  $(0, 1)$  случайных величин. Следовательно,  $\bar{Y}$  и  $S_Y$  независимы,  $\bar{Y} \sim \mathcal{N}(0, 1/n)$ ,  $S_Y \sim \chi_{n-1}^2$  (см. лемму 6.1), и “лемма Фишера” доказана.

Установив совместное распределение выборочного среднего и выборочной дисперсии в случае выбора из нормального распределения, мы можем приступить к построению доверительных интервалов для каждого из параметров  $\mu$  и  $\sigma$  при неизвестном значении другого параметра.

3<sup>0</sup>. ВЕРХНЯЯ ДОВЕРИТЕЛЬНАЯ ГРАНИЦА ДЛЯ ДИСПЕРСИИ НОРМАЛЬНОГО РАСПРЕДЕЛЕНИЯ ПРИ НЕИЗВЕСТНОМ СРЕДНЕМ. Эта граница на-

ходится наиболее просто, поскольку распределение опорной функции

$$\frac{nS^2}{\sigma^2} = \frac{\sum_1^n (X_k - \bar{X})^2}{\sigma^2} = \sum_{k=1}^n \left( \frac{X_k - \mu}{\sigma} - \frac{1}{n} \sum_{i=1}^n \frac{X_i - \mu}{\sigma} \right)^2 = \sum_1^n (Y_k - \bar{Y})^2 \quad (2)$$

есть хи-квадрат распределение с  $n - 1$  степенью свободы (см. теорему 6.1). Следовательно, верхняя  $(1 - \alpha)$ -доверительная граница определяется квантилью  $\lambda_\alpha = K^{-1}(\alpha)$  хи-квадрат распределения – корнем уравнения

$$P(nS^2/\sigma^2 \geq \lambda) = 1 - K_{n-1}(\lambda) = 1 - \alpha,$$

и доверительное утверждение  $\sigma^2 \leq \bar{\sigma}_n^2 = nS^2/K_{n-1}^{-1}(\alpha)$  выполняется с заданной вероятностью  $1 - \alpha$ .

4<sup>0</sup>. **ДОВЕРИТЕЛЬНЫЙ ИНТЕРВАЛ ДЛЯ СРЕДНЕГО ЗНАЧЕНИЯ НОРМАЛЬНОГО РАСПРЕДЕЛЕНИЯ ПРИ НЕИЗВЕСТНОЙ ДИСПЕРСИИ.** В этой задаче мы имеем дело с двусторонними доверительными границами (доверительным интервалом), и в соответствии с выбором опорной функции  $|\bar{X} - \mu|/S$ , о которой мы говорили перед доказательством теоремы 6.1, нам потребуется знание вероятности события вида  $|\bar{X} - \mu|/S \leq \lambda$ .

В начале XIX века английский математик В.Госсет, писавший под псевдонимом “Стьюдент” (Student), нашел распределение случайной величины  $T_\nu = \xi\sqrt{\nu}/\sqrt{\chi_\nu^2}$ , где  $\xi \sim \mathcal{N}(0, 1)$ , а  $\chi_\nu^2$  – случайная величина, не зависящая от  $\xi$  и распределенная по закону хи-квадрат с  $\nu$  степенями свободы. Естественно, его исследования были связаны с проблемами статистического вывода о среднем значении  $\mu$  нормального распределения при неизвестной дисперсии, и Стьюдент искал распределение опорной функции (см. (2)) в связи с переходом в записи опорной функции в терминах  $X_k$  к  $Y_k$

$$H = \frac{\bar{X} - \mu}{S} \sqrt{n-1} = \frac{\frac{1}{\sqrt{n}} \sum_1^n Y_k}{\sqrt{\sum_1^n (Y_k - \bar{Y})^2}} \sqrt{n-1},$$

$$Y_k = \frac{X_k - \mu}{\sigma} \sim \mathcal{N}(0, 1), \quad k = 1, \dots, n,$$

которая отличается от выбранной нами опорной функции только множителем  $\sqrt{n-1}$ , и поэтому также может быть использована в построении доверительного интервала для  $\mu$  при неизвестном  $\sigma$ . То, что распределения  $T_{n-1}$  и  $H$  совпадают, следует из теоремы 6.1: случайная величина в знаменателе

$\xi = \sum_1^n Y_k/\sqrt{n} \sim \mathcal{N}(0, 1)$  не зависит от  $\sum_1^n (Y_k - \bar{Y})^2 \sim \chi_{n-1}^2$ , разделив которую на значение степени свободы  $n - 1$ , получаем  $\frac{1}{\nu}\chi_\nu^2$  с  $\nu = n - 1$ .

Найдем распределение случайной величины  $T_\nu$ , которое называется *распределением Стьюдента с  $\nu$  степенями свободы* или *t-распределением*. Совместная функция плотности независимых случайных величин  $\xi$  и  $\eta = \chi_\nu^2$  равна

$$f(x, y) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{x^2}{2}\right\} \frac{1}{2^{\nu/2}\Gamma(\nu/2)} y^{\nu/2-1} \exp\left\{-\frac{y}{2}\right\},$$

так что функция распределения случайной величины  $T_\nu$

$$S_\nu(t) = P(\xi\sqrt{\nu/\eta} < t) = \int_{x\sqrt{\nu} < \sqrt{y}t} \int f(x, y) dx dy = \int_0^\infty dy \int_{-\infty}^{t\sqrt{y/\nu}} f(x, y) dx.$$

Дифференцируя это выражение по  $t$ , находим функцию плотности распределения Стьюдента

$$\begin{aligned} s_\nu(t) &= \int_0^\infty \sqrt{y/\nu} f(t\sqrt{y/\nu}, y) dy = \\ &= \frac{1}{\sqrt{\pi\nu}2^{(\nu+1)/2}\Gamma(\nu/2)} \int_0^\infty y^{\frac{\nu+1}{2}-1} \exp\left\{-\frac{y}{2}\left(1 + \frac{t^2}{\nu}\right)\right\} dt = \\ &= \frac{1}{\sqrt{\pi\nu}} \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right)} \left(1 + \frac{t^2}{\nu}\right)^{-\frac{\nu+1}{2}}. \end{aligned}$$

Вид полученной функции плотности говорит о том, что распределение Стьюдента можно трактовать как обобщение стандартного ( $a = 0, b = 1$ ) распределения Коши  $S(a, b)$ , которое получается из распределения Стьюдента при числе степеней свободы  $\nu = 1$ . Это симметричное распределение, и поэтому  $S_\nu(-t) = 1 - S_\nu(t)$ , что позволяет нам довольно просто построить доверительный интервал для  $\mu$  с помощью квантили распределения  $S_{n-1}(\cdot)$ :

$$P(|T_{n-1}| \leq t) = S_{n-1}(t) - S_{n-1}(-t) = 2S_{n-1}(t) - 1 = 1 - \alpha,$$

откуда  $t_\alpha = S_{n-1}^{-1}(1 - \alpha/2)$ , и  $(1 - \alpha)$ -доверительный интервал для  $\mu$  определяется пределами  $\bar{X} \pm St_\alpha/\sqrt{n-1}$ .

Итак, мы построили доверительные пределы для параметров  $\mu$  и  $\sigma^2$  нормального распределения. Таблицы нормального, хи-квадрат и студентского распределений, а также квантилей этих распределений, необходимые для численной реализации доверительных оценок, смотрите в книге Большев Л.Н., Смирнов Н.В. Таблицы математической статистики, М.: Наука, 1983, которая в дальнейшем будет цитироваться как ТМС. Еще раз отметим, что возможность доверительной оценки этих параметров определялась, в основном, инвариантностью семейства нормальных распределений относительно линейной группы преобразований. Точно так же мы можем построить доверительные пределы для параметра  $\theta$  показательного распределения или для параметра масштаба гамма распределения при известном параметре формы; мы вернемся к этим задачам позднее при обсуждении проблемы оптимизации доверительной оценки. Что же касается других распределений, то здесь проблема усложняется отсутствием инвариантных опорных функций и невозможностью получить распределение оценок параметра, для которого строятся доверительные пределы, в явном виде. Тем не менее существует достаточно общий подход к данной проблеме, основанный на асимптотической нормальности распределения оценок по методу моментов или методу максимального правдоподобия.

Пусть  $\hat{\theta}_n = \hat{\theta}_n(X^{(n)})$  – асимптотически нормальная со средним  $\theta$  и дисперсией  $\sigma^2(\theta)/n$  оценка параметра  $\theta$  (например, при определенных условиях регулярности (см. теорему 4.2) оценка максимального правдоподобия асимптотически нормальна со средним  $\theta$  и дисперсией  $[nI(\theta)]^{-1}$ ). Тогда при  $n \rightarrow \infty$  вероятность

$$P_\theta \left( \frac{|\hat{\theta}_n - \theta|}{\sigma(\theta)} \sqrt{n} \leq \lambda_\alpha \right) \rightarrow 1 - \alpha,$$

при любом  $\theta \in \Theta$ , если  $\lambda_\alpha = \Phi^{-1}(1 - \alpha/2)$ , и мы получаем *асимптотически*  $(1 - \alpha)$ -доверительное множество

$$\Delta_n(X^{(n)}) = \left\{ \theta : \frac{|\hat{\theta}_n - \theta|}{\sigma(\theta)} \sqrt{n} \leq \lambda_\alpha \right\} \cap \Theta.$$

Если  $\Delta_n$  есть интервал на прямой  $\mathbb{R}$ , то мы решили задачу интервальной оценки параметра  $\theta$ . Если же это некоторое вычурное и непригодное



к употреблению подмножество  $\mathbb{R}$ , то можно пойти на дальнейшие упрощения асимптотического утверждения, заменив в определении  $\Delta_n$  параметрическую функцию  $\sigma(\theta)$  на ее оценку  $\sigma(\hat{\theta}_n)$ . Достаточно потребовать непрерывность функции  $\sigma(\theta)$ ,  $\theta \in \Theta$ , чтобы, ссылаясь на теорему Слуцкого (предложение 11.1 курса ТВ с  $\xi_n \propto \sigma(\hat{\theta}_n) \xrightarrow{P} \sigma(\theta)$ ), утверждать, что при  $n \rightarrow \infty$

$$P_\theta \left( \hat{\theta}_n - \frac{\lambda_\alpha \sigma(\hat{\theta}_n)}{\sqrt{n}} \leq \theta \leq \hat{\theta}_n + \frac{\lambda_\alpha \sigma(\hat{\theta}_n)}{\sqrt{n}} \right) \rightarrow 1 - \alpha.$$

Проиллюстрируем “работу” этого метода на двух полезных в практическом отношении примерах.

5<sup>0</sup>. АСИМПТОТИЧЕСКИ ДОВЕРИТЕЛЬНЫЙ ИНТЕРВАЛ ДЛЯ ВЕРОЯТНОСТИ УСПЕХА В ИСПЫТАНИЯХ БЕРНУЛЛИ. В схеме испытаний Бернулли – выборе из распределения бинарной случайной величины  $X$ , принимающей значение 1 (“успех”) с вероятностью  $p$  и значение 0 (“неудача”) с вероятностью  $1 - p$ , оптимальной несмещенной оценкой  $p$  является выборочное среднее  $\bar{X} = n^{-1} \sum_1^n X_k$  или, что то же, относительная частота успешных исходов в  $n$  испытаниях. Статистика  $n\bar{X}$  имеет биномиальное распределение  $B(n, p)$ , и это позволяет насчитать таблицы доверительных пределов для  $p$  при различных значениях доверительного уровня  $1 - \alpha$ , объема выборки  $n$  и числа успешных исходов  $n\bar{x}$  (см., например, ТМС). Что же дает асимптотический подход к построению доверительных интервалов?

Выборочное среднее асимптотически нормально со средним  $p$  и дисперсией  $p(1 - p)/n$ . Следовательно,  $(1 - \alpha)$ -доверительная область  $\Delta_n = \left\{ p : 0 \leq p \leq 1, |\bar{X} - p| \leq \lambda_\alpha \sqrt{p(1 - p)/n} \right\}$ . Разрешая неравенства в фигурных скобках относительно  $p$ , получаем доверительный интервал

$$\frac{n}{n + \lambda_\alpha^2} \left( \bar{X} + \frac{\lambda_\alpha^2}{2n} \pm \lambda_\alpha \sqrt{\frac{\bar{X}(1 - \bar{X})}{n} + \frac{\lambda_\alpha^2}{4n^2}} \right),$$

который при больших объемах испытаний  $n$  мало отличается от доверительного интервала  $\bar{X} \pm \lambda_\alpha \sqrt{\bar{X}(1 - \bar{X})/n}$ , полученного заменой  $\sigma^2(p) = p(1 - p)$  на ее оценку  $\bar{X}(1 - \bar{X})$ :

6<sup>0</sup>. АСИМПТОТИЧЕСКИ ДОВЕРИТЕЛЬНЫЙ ИНТЕРВАЛ ДЛЯ ПАРАМЕТРА ИНТЕНСИВНОСТИ РАСПРЕДЕЛЕНИЯ ПУАССОНА. Распределение Пуассона  $P(\theta)$  с функцией плотности (по считающей мере)  $f(x | \theta) = P_\theta(X = x) = \theta^x e^{-\theta} / x!$ ,  $x = 0, 1, 2, \dots$ , индексируется положительным параметром  $\theta$ , оптимальная несмещенная оценка которого по выборке  $X^{(n)}$  объема  $n$ , как и в предыдущем примере, определяется выборочным средним  $\bar{X}$ . Для распределения Пуассона также справедлива теорема сложения:  $n\bar{X} \sim P(n\theta)$ , и на основе этого можно построить точные доверительные пределы для  $\theta$ , таблица которых имеется в упомянутом сборнике ТМС. Но оценка  $\bar{X}$  асимптотически нормальна  $(\theta, \theta/n)$ , что позволяет определить асимптотически доверительную область  $\Delta_n = \{\theta : \theta > 0, |\bar{X} - \theta| \leq \lambda_\alpha \sqrt{\theta/n}\}$ . Решение неравенств в фигурных скобках относительно  $\theta$  дает асимптотически доверительный интервал

$$\bar{X} + \frac{\lambda_\alpha^2}{2n} \pm \lambda_\alpha \sqrt{\frac{\bar{X}}{n} + \frac{\lambda_\alpha^2}{4n^2}}.$$

Наконец, заменяя  $\sigma^2(\theta) = \theta$  ее оценкой  $\bar{X}$ , получаем также асимптотически доверительный, но, как показывают числовые расчеты, менее точный интервал  $\bar{X} \pm \lambda_\alpha \sqrt{\bar{X}/n}$ .

На этом мы заканчиваем изложение простейших методов построения доверительных и асимптотически доверительных интервалов на основе подбора опорных функций. О проблеме оптимального интервального оценивания мы поговорим позднее, изучив теорию оптимальной проверки гипотез – высказываний о возможных значениях параметра  $\theta$ . Оставшиеся лекции будут посвящены именно этой теории.

## §7. Статистическая проверка гипотез (критерии значимости)

### Лекция 11

В приложениях математической статистики существует обширный класс задач, в которых требуется проверить истинность некоторого высказывания относительно исследуемого объекта или выбрать одно из альтернативных решений, которое определит дальнейшее поведение статистика по отношению к этому объекту. Например, при аттестации партии дизельного топлива по общему содержанию серы мы должны не только дать точечную оценку данной характеристики топлива, но и принять решение о качестве выпускаемого продукта, которое повлечет за собой одно из следующих действий – или отослать топливо потребителю, или произвести дополнительную очистку топлива от вредных примесей. Точно так же в примере 1.2 мы строили статистическое правило, позволяющее принять одно из двух решений относительно нового лечебного препарата – или признать его эффективным и внедрить в лечебную практику, или запретить его дальнейшее использование. В исследованиях, подобных опытам Менделя, часто надо проверить гипотезу относительно предполагаемого значения вероятности наследования доминантного признака. Селекционер, работающий над получением нового вида пшеницы, должен подкрепить свое заключение о превосходстве нового вида над тем, который уже используется в сельскохозяйственной практике, с помощью сопоставления данных об урожайности этих видов. И так далее, и тому подобное, – вы сами можете привести примеры таких задач по выбору одного из ряда альтернативных решений.

В нашем курсе математической статистики мы рассмотрим задачи, связанные только с выбором одного из двух решений. Пусть мы высказываем некоторое суждение (или предпринимаем действие) об исследуемом объекте, и пусть  $d_0$  – решение об истинности этого суждения, в то время как  $d_1$  – решение о его ложности. Таким образом, пространство решений  $\mathcal{D}$  в данной статистической проблеме состоит из точек:  $\mathcal{D} = \{d_0, d_1\}$ .

Для выбора одного из решений мы наблюдаем случайную выборку  $X^{(n)}$  из некоторого распределения  $P_\theta$ , значение параметра  $\theta$  которого нам неизвестно. Пусть  $\Theta$  – область возможных значений  $\theta$ , которую мы назвали параметрическим пространством. В соответствии с принятой нами в §1 идеологией статистического вывода мы сопоставляем каждому решению  $d \in \mathcal{D}$  определенное подмножество  $\Theta_d$  пространства  $\Theta$ , то есть интерпретируем каждое решение в терминах высказываний об истинном значении

параметра  $\theta$ . В нашей статистической проблеме выбора одного из двух решений положим  $\Theta_i = \Theta_{d_i}$ ,  $i = 0, 1$ , и введем ряд понятий и определений, используемых при решении этой проблемы.

Утверждение  $H_0 : \theta \in \Theta_0$  называется *нулевой гипотезой*, а утверждение  $H_1 : \theta \in \Theta_1$  – *альтернативной гипотезой* или (коротко) *альтернативой*. Гипотеза  $H_i$  называется *простой*, если соответствующее  $\Theta_i$  состоит из одной точки параметрического пространства  $\Theta$ ; в противном случае  $H_i$  называется *сложной* гипотезой;  $i = 0, 1$ . Так, в примере 1.2 с испытанием нового лечебного препарата параметр  $\theta$  означал вероятность успешного лечения каждого пациента, и нулевая гипотеза  $H_0 : \theta = 1/2$  о “нейтральности” препарата есть простая гипотеза, в то время как альтернативная гипотеза  $H_1 : \theta > 1/2$  об его эффективности – сложная гипотеза.

Правило, по которому принимается или отвергается нулевая гипотеза  $H_0$ , называется *критерием*. Иногда добавляется – *критерий согласия* (с нулевой гипотезой), особенно, когда альтернатива  $H_1$  определена не совсем четко и под  $H_1$  подразумевается “все остальное”. В случае полного равноправия гипотез говорят о критерии *различения гипотез*. Критерий определяется заданием особого подмножества  $S$  выборочного пространства  $X^n$ , которое называется *критической областью*: если выборочные данные  $x^{(n)}$  попадают в эту область, то нулевая гипотеза  $H_0$  отклоняется и принимается альтернативное решение – справедлива  $H_1$ . Область  $A = S^c = X^n \setminus S$  называется *областью принятия* нулевой гипотезы. Нам будет удобно проводить спецификацию критической области в виде ее индикаторной функции  $\varphi = \varphi(X^{(n)})$ , которая называется *критической функцией* или, поскольку она определяет статистическое правило проверки гипотезы, просто *критерием*. Итак, функция  $\varphi(X^{(n)})$  есть бинарная случайная величина, принимающая значение 1, если произошло событие  $X^{(n)} \in S$ , и значение 0, если произошло противоположное событие  $X^{(n)} \in A$ . Понятно, что математическое ожидание  $\mathbf{E}\varphi(X^{(n)})$  означает вероятность отклонения гипотезы  $H_0$ .

В рассматриваемой статистической проблеме величина риска, связанная с отклонением верной гипотезы, обычно соотносится с функцией потерь типа 1 – 0: потери считаются равными 1, если принята гипотеза  $H_i$ , а в действительности  $\theta \in \Theta_{1-i}$ ,  $i = 0, 1$ ; если же принята  $H_i$  и  $\theta \in \Theta_i$ ,  $i = 0, 1$ , то потери полагаются равными нулю. Легко видеть, что величина риска при любом значении параметра  $\theta$  может быть определена с помощью функции  $m(\theta) = \mathbf{E}_\theta\varphi(X^{(n)}) = P_\theta(X^{(n)} \in S)$ , которая называется *функцией мощности* критерия  $\varphi$ . Эта функция указывает, как часто мы отклоняем нулевую

гипотезу, когда  $\theta$  – истинное значение параметра, и хорошим следует считать тот критерий, у которого функция  $m(\theta)$  принимает близкие к нулю значения в области  $\Theta_0$  и близкие к единице – в области  $\Theta_1$ . В связи с этим вводятся две компоненты функции риска:  $\alpha(\theta) = m(\theta)$  при  $\theta \in \Theta_0$  и  $\beta(\theta) = 1 - m(\theta)$  при  $\theta \in \Theta_1$ . Функция  $\alpha(\theta)$ ,  $\theta \in \Theta_0$  называется *вероятностью ошибки первого рода* – она указывает относительную частоту отклонения гипотезы  $H_0$ , когда она в действительности верна ( $\theta \in \Theta_0$ ). Функция  $\beta(\theta)$ ,  $\theta \in \Theta_1$  называется *вероятностью ошибки второго рода* – она указывает относительную частоту принятия гипотезы  $H_0$ , когда она ложна (верна альтернативная гипотеза  $H_1 : \theta \in \Theta_1$ ). Заметим, что функция мощности  $m(\theta)$  в области  $\Theta_1$  трактуется как вероятность отклонения гипотезы  $H_0$ , когда в действительности выбор идет из распределения с альтернативным значением  $\theta \in \Theta_1$ , и поэтому часть  $m(\theta)$  при  $\theta \in \Theta_1$  называется *мощностью* критерия  $\varphi$ .

Легко понять, что при фиксированном объеме наблюдений  $n$  невозможно одновременно минимизировать вероятности обеих ошибок, – для уменьшения вероятности ошибки первого рода  $\alpha(\theta) = P_\theta(X^{(n)} \in S)$ ,  $\theta \in \Theta_0$ , необходимо уменьшить критическую область  $S$ , что приведет к увеличению области  $A$  принятия нулевой гипотезы и, следовательно, к увеличению вероятности ошибки второго рода  $\beta(u) = P_u(X^{(n)} \in A)$ ,  $u \in \Theta_1$ . Здесь возникает такая же ситуация, что и в проблеме построения оценки параметра  $\theta$  с равномерно минимальным риском, – такие оценки существуют только в определенном классе статистических правил, например, в классе несмещенных оценок. Однако, даже и помимо задачи проверки гипотез с минимальной вероятностью ошибки, и намного раньше создания общей теории наиболее мощных критериев в статистической практике сложился следующий подход к управлению риском критерия.

Предположим, что отклонение гипотезы  $H_0$ , когда она в действительности верна, приводит к более тяжким последствиям, чем ее принятие при справедливости альтернативы. В таком случае мы заинтересованы в первую очередь контролировать вероятность ошибки первого рода. С этой целью заранее фиксируется (выбирается) некоторый уровень  $\alpha$ , выше которого вероятность ошибки первого рода не допустима, и критическая область  $S$  (критерий  $\varphi$ ) определяется таким образом, что  $\alpha(\theta) \leq \alpha$ , каково бы ни было  $\theta \in \Theta_0$ . Это ограничение  $\alpha$  на вероятность ошибки первого рода называется *уровнем значимости*, а сам критерий  $\varphi$ , для которого выполняется это ограничение, – *критерием уровня  $\alpha$* . Наибольшее значение вероятности

ошибки первого рода

$$\bar{\alpha} = \sup_{\theta \in \Theta_0} \alpha(\theta)$$

называется *размером* критерия  $\varphi$ , и если  $\bar{\alpha} = \alpha$ , то говорят о *критерии  $\varphi$  размера  $\alpha$* .

В этом выборе ограничения именно на вероятность ошибки первого, а не второго рода проявляется типичная асимметрия в практической ценности гипотезы и альтернативы. Например, если проверяется эффективность нового лекарственного препарата, то нулевой гипотезе должно соответствовать решение о его неэффективности, ибо, отклонив эту гипотезу, когда она верна, мы внедрим в лечебную практику бесполезное или вредное лекарство, что приведет к более тяжким последствиям, чем отклонение в действительности эффективного препарата. Но если мы ищем золото, анализируя состав кернов при бурении предполагаемого месторождения, то естественно принять за нулевую гипотезу утверждение о наличии золота, ибо отклонив ее, когда она верна, мы потеряем намного больше, чем стоимость нескольких дополнительных анализов, удостоверяющих, что золото в разбуренной местности отсутствует.

Следует также обратить особое внимание на общую методологию проверки гипотез, отражаемую в выборе малого значения уровня  $\alpha$ . Если наши выборочные данные попадают в область  $S$  с исключительно малой вероятностью, то естественно предположить, что то утверждение, которое привело к этому маловероятному событию, не соответствует истине и отклонить его. Поступая таким образом, мы будем терять в действительности верную гипотезу  $H_0$  крайне редко – не более, чем в  $100\alpha\%$  случаев.

Простейший метод построения критериев значимости состоит в использовании состоятельных оценок тестируемого параметра  $\theta$ . Рассмотрим простейший случай:  $\theta$  – скалярный параметр, вероятностная модель не содержит других (мешающих) параметров и проверяется простая гипотеза  $H_0 : \theta = \theta_0$  при альтернативе  $H_1 : \theta \neq \theta_0$ , где  $\theta_0$  – некоторое, априори фиксированное значение параметра  $\theta$  (например, в опытах Менделя проверяется гипотеза: вероятность  $\theta$  наследования доминантного признака равна  $\theta_0 = 3/4$ ). Если  $\hat{\theta}_n = \hat{\theta}_n(X^{(n)})$  – состоятельная оценка  $\theta$ , дисперсия которой стремится к нулю при  $n \rightarrow \infty$  как  $O(1/\sqrt{n})$ , то естественно определить критическую область посредством неравенства  $\sqrt{n}|\hat{\theta}_n(X^{(n)}) - \theta_0| > C$ . Вероятность ошибки первого рода такого критерия  $\alpha(\theta_0, C) = P_{\theta_0}(\sqrt{n}|\hat{\theta}_n(X^{(n)}) - \theta_0| > C)$ , и приравнивая эту вероятность заданному уровню значимости  $\alpha$ ,

находим *критическую константу*  $C = C(\alpha)$  как квантиль распределения случайной величины  $\sqrt{n}|\hat{\theta}_n(X^{(n)}) - \theta_0|$ ; такой выбор  $C$  приводит к критерию уровня  $\alpha$ . Если  $\theta (\neq \theta_0)$  – некоторое альтернативное значение параметра, то, в силу состоятельности оценки,  $\sqrt{n}|\hat{\theta}_n(X^{(n)}) - \theta_0| \xrightarrow{P} \infty$ , и поэтому вероятность ошибки второго рода  $\beta(\theta) = P_\theta(\sqrt{n}|\hat{\theta}_n(X^{(n)}) - \theta_0| \leq C(\alpha)) \rightarrow 0$ , когда  $n \rightarrow \infty$ . Таким образом, мы получаем критерий заданного уровня  $\alpha$ , обладающий к тому же свойством *состоятельности* – его вероятность ошибки второго рода стремится к нулю при неограниченном возрастании объема выборки  $n$ .

Сформулируем теперь основную задачу теории статистической проверки гипотез: *требуется найти такой критерий  $\varphi$  уровня  $\alpha$ , который равномерно по всем  $\theta \in \Theta_1$  максимизирует мощность  $m(\theta)$  или, что то же, равномерно по  $\theta \in \Theta_1$  минимизирует вероятность ошибки второго рода  $\beta(\theta)$* . Мы укажем метод построения таких *равномерно наиболее мощных критериев* заданного уровня  $\alpha$  в следующем параграфе, а пока обратимся к иллюстрациям введенных понятий и построению наиболее часто используемых на практике критериев, касающихся проверки гипотез о значениях параметров нормального распределения.

1<sup>0</sup>. ПРОВЕРКА ГИПОТЕЗЫ О ВЕЛИЧИНЕ СРЕДНЕГО ЗНАЧЕНИЯ НОРМАЛЬНОГО РАСПРЕДЕЛЕНИЯ ПРИ ИЗВЕСТНОЙ ДИСПЕРСИИ. Рассмотрим сначала наиболее часто встречающуюся в практических применениях математической статистики задачу проверки сложной гипотезы  $H_0 : \mu \leq \mu_0$  при сложной альтернативе  $H_1 : \mu > \mu_0$  о среднем значении  $\mu$  нормального  $(\mu, \sigma^2)$  распределения при известном значении дисперсии  $\sigma^2$ . Выборочное среднее  $\bar{X} = n^{-1} \sum_1^n X_k$  есть оптимальная оценка неизвестного значения  $\mu$ , и поэтому, в соответствии с только что предложенным методом построения состоятельных критериев, рассмотрим критерий, отвергающий нулевую гипотезу  $H_0$ , когда  $\sqrt{n}(\bar{X} - \mu_0) > C$ , или, что то же,  $\bar{X} > C$ , поскольку значения  $\mu_0$  и  $n$  фиксированы и известны. Постоянная  $C$  должна выбираться по заданному уровню значимости  $\alpha$ , ограничивающему максимальное значение вероятности ошибки первого рода.

Так как при выборе из нормального  $(\mu, \sigma^2)$  распределения статистика  $\bar{X} \sim \mathcal{N}(\mu, \sigma^2/n)$ , то функция мощности этого критерия

$$m(\mu) = P_\mu(\bar{X} > C) = 1 - \Phi\left(\frac{C - \mu}{\sigma} \sqrt{n}\right) = \Phi\left(\frac{\mu - C}{\sigma} \sqrt{n}\right).$$

Легко видеть, что  $m(\mu)$  – строго возрастает с ростом  $\mu$ , так что размер критерия

$$\bar{\alpha} = \max_{\mu \leq \mu_0} m(\mu) = m(\mu_0) = 1 - \Phi \left( \frac{C - \mu_0}{\sigma} \sqrt{n} \right).$$

Приравнивая размер критерия уровню значимости  $\alpha$ , находим критическое значение  $C(\alpha) = \mu_0 + \Phi^{-1}(1 - \alpha)\sigma/\sqrt{n}$ .

Вероятность ошибки второго рода нашего критерия размера  $\alpha$

$$\begin{aligned} \beta(\mu) &= P_{\mu}(\bar{X} \leq C(\alpha)) = \Phi \left( \frac{C(\alpha) - \mu}{\sigma} \sqrt{n} \right) = \\ &= \Phi \left( \frac{\mu_0 - \mu}{\sigma} \sqrt{n} + \Phi^{-1}(1 - \alpha) \right), \quad \mu > \mu_0, \end{aligned} \quad (1)$$

убывает с ростом  $\mu$  по мере ее отхода от граничного значения  $\mu_0$ . Наибольшее значение  $\beta(\mu)$  достигается в точке  $\mu = \mu_0$  и равно  $1 - \alpha$ . Это значение не зависит от размера выборки  $n$ , и поэтому требуются дополнительные соображения при планировании объема наблюдений. Обычно используется метод введения так называемой *зоны безразличия* – интервала  $(\mu_0, \mu_1)$ , который выбирается из тех соображений, что при истинном значении  $\mu \in (\mu_0, \mu_1)$  принятие нулевой гипотезы  $H_0$  не приводит к слишком тяжелым последствиям. Однако при истинном  $\mu \geq \mu_1$  вероятность принятия  $H_0$  должна быть под контролем и не превосходить некоторого предписанного значения  $\beta$ . Это обстоятельство позволяет спланировать объем выборки  $n$ , определив его из неравенства  $\beta(\mu_1) \leq \beta$ . Используя формулу (1) для  $\beta(\mu)$ , находим, что *объем выборки*  $n = n(\alpha, \beta, \mu_0, \mu_1)$ , *необходимый для различения гипотез*  $\mu \leq \mu_0$  и  $\mu \geq \mu_1$  *с заданными ограничениями*  $\alpha$  и  $\beta$  *на вероятности ошибок первого и второго рода*, равен наименьшему целому  $n$ , удовлетворяющему неравенству

$$n \geq \frac{[\Phi^{-1}(1 - \alpha) + \Phi^{-1}(1 - \beta)]^2}{(\mu_1 - \mu_0)^2} \sigma^2.$$

Аналогичным методом строится критерий для проверки простой гипотезы  $H_0 : \mu = \mu_0$  при сложной альтернативе  $H_1 : \mu \neq \mu_0$ . В этой задаче естественно определить критическую область посредством неравенства  $|\bar{X} - \mu_0| > C$ . Функция мощности такого критерия

$$m(\mu) = 1 - \left[ \Phi \left( \frac{C + \mu_0 - \mu}{\sigma} \sqrt{n} \right) - \Phi \left( \frac{-C + \mu_0 - \mu}{\sigma} \sqrt{n} \right) \right]$$



строго убывает при  $\mu < \mu_0$ , возрастает при  $\mu > \mu_0$  и при  $\mu = \mu_0$  равна вероятности ошибки первого рода. Таким образом, критическая константа  $C = C(\alpha)$  определяется по заданному уровню значимости  $\alpha$  из уравнения

$$m(\mu_0) = 1 - \left[ \Phi \left( \frac{C}{\sigma} \sqrt{n} \right) - \Phi \left( \frac{-C}{\sigma} \sqrt{n} \right) \right] = 2 \left[ 1 - \Phi \left( \frac{C}{\sigma} \sqrt{n} \right) \right] = 1 - \alpha,$$

откуда  $C(\alpha) = \Phi^{-1}(1 - \alpha/2)\sigma/\sqrt{n}$ .

## Лекция 12

2<sup>0</sup>. ПРОВЕРКА ГИПОТЕЗЫ О ВЕЛИЧИНЕ ДИСПЕРСИИ НОРМАЛЬНОГО РАСПРЕДЕЛЕНИЯ ПРИ НЕИЗВЕСТНОМ СРЕДНЕМ ЗНАЧЕНИИ. Это типичная задача контроля за величиной случайной ошибки в параллельных наблюдениях некоторой характеристики исследуемого объекта. Так как превышение характеристики случайной погрешности  $\sigma$  над некоторым номиналом  $\sigma_0$  в случае, когда мы утверждаем  $\sigma \leq \sigma_0$ , влечет более серьезные последствия, чем неоправданные претензии к слишком большому разбросу в данных, то следует принять за нулевую гипотезу  $\sigma > \sigma_0$ . Проверка этой гипотезы проводится при естественной альтернативе  $H_1 : \sigma \leq \sigma_0$ , причем мы не знаем значения мешающего параметра  $\mu$  – среднего значения нормального распределения, из которого производится выбор.

Как нам известно, выборочная дисперсия  $S^2 = n^{-1} \sum_1^n (X_k - \bar{X})^2$  есть состоятельная оценка  $\sigma^2$ , ее распределение не зависит от  $\mu$ , а случайная величина  $nS^2/\sigma^2$  имеет хи-квадрат распределение с  $n-1$  степенью свободы. Таким образом, разумно рассмотреть критерий с критической областью  $nS^2 < C$ . Функция мощности такого критерия

$$m(\sigma) = P_{\mu, \sigma} \left( \frac{nS^2}{\sigma^2} \leq \frac{C}{\sigma^2} \right) = K_{n-1} \left( \frac{C}{\sigma^2} \right)$$

монотонно убывает с ростом  $\sigma$ , поэтому наибольшее значение вероятности ошибки первого рода достигается при  $\sigma = \sigma_0$ , и критическое значение  $C(\alpha)$  критерия требуемого размера  $\alpha$  определяется из уравнения  $K_{n-1}(C\sigma_0^{-2}) = \alpha$ . Итак,  $C(\alpha) = \sigma_0^2 K_{n-1}^{-1}(\alpha)$ ; вероятность ошибки второго рода

$$\beta(\sigma) = P_{\mu, \sigma} (nS^2 > C(\alpha)) = 1 - K_{n-1} \left( \frac{\sigma_0^2}{\sigma^2} K_{n-1}^{-1}(\alpha) \right), \quad \sigma \leq \sigma_0,$$

монотонно убывает по мере отхода истинного значения  $\sigma$  от номинала  $\sigma_0$ .

3<sup>0</sup>. ПРОВЕРКА ГИПОТЕЗЫ О ВЕЛИЧИНЕ СРЕДНЕГО ЗНАЧЕНИЯ НОРМАЛЬНОГО РАСПРЕДЕЛЕНИЯ ПРИ НЕИЗВЕСТНОЙ ДИСПЕРСИИ (ОДНОВЫБОРОЧНЫЙ КРИТЕРИЙ СТЬЮДЕНТА). Вы, наверное, обратили внимание, что при построении критериев значимости мы по существу используем методы построения доверительных множеств? Это, действительно, так – между задачами доверительной оценки и проверки гипотез существует много общего, и, решив одну задачу, мы сразу же получаем решение другой. В конце этого параграфа мы формализуем этот параллелизм, а пока будем использовать его на интуитивном уровне: предлагается использовать для проверки гипотез о среднем значении нормального распределения статистику Стьюдента.

Рассмотрим сначала задачу проверки сложной гипотезы  $H_0 : \mu \leq \mu_0$  при сложной альтернативе  $H_1 : \mu > \mu_0$ . Так как выборочное среднее  $\bar{X}$  есть состоятельная оценка значения  $\mu$ , то статистика Стьюдента

$$T = \frac{\bar{X} - \mu_0}{S} \sqrt{n - 1}$$

опосредственно, через выборочные данные, характеризует удаленность истинного среднего значения  $\mu$  от границы  $\mu_0$ , разделяющей гипотезу и альтернативу. Поэтому предлагается отвергать нулевую гипотезу  $\mu \leq \mu_0$ , если  $T > C$ , выбирая  $C$ , как обычно, по заданному уровню значимости  $\alpha$ . Для решения последней задачи необходимо исследовать поведение функции мощности  $m(\mu; \sigma) = P_{\mu, \sigma}(T > C)$  критерия  $T > C$ . Если мы покажем, что  $m(\mu; \sigma)$  есть монотонно возрастающая функция аргумента  $\mu$  при любом фиксированном значении аргумента  $\sigma$ , то наибольшее значение вероятности ошибки первого рода  $\alpha(\mu; \sigma) = m(\mu; \sigma)$ ,  $\mu \leq \mu_0$  при каждом фиксированном  $\sigma$  будет достигаться в точке  $\mu = \mu_0$ . Следовательно, размер критерия в таком случае будет равен (см. пункт 4<sup>0</sup> предыдущего параграфа)  $\bar{\alpha} = m(\mu_0; \sigma) = P_{\mu_0, \sigma}(T > C) = 1 - S_{n-1}(C)$ , где  $S_\nu(\cdot)$  – функция распределения Стьюдента с  $\nu$  степенями свободы. Таким образом, мы получим свободный от неизвестного значения  $\sigma$  критерий  $T > C(\alpha)$  требуемого размера  $\alpha$  с критической константой  $C(\alpha) = S_{n-1}^{-1}(1 - \alpha)$ . Это и есть то статистическое правило, которое обычно называется *критерием Стьюдента* или *t-критерием*.

Покажем теперь, что вероятность (функция мощности)  $P_{\mu, \sigma}(T > C)$  монотонно возрастает с ростом  $\mu$  при любых фиксированных значениях  $\sigma$

и  $C$ . С этой целью представим статистику  $T$  в следующем виде:

$$T = \frac{\bar{X} - \mu}{S} \sqrt{n-1} + \frac{\mu - \mu_0}{\sigma} \cdot \frac{\sigma}{S} \sqrt{n-1}.$$

Если  $\mu$  – среднее значение нормального распределения, из которого происходит выбор, то первое слагаемое в этом представлении есть студентовская случайная величина с  $n - 1$  степенью свободы. Второе слагаемое есть произведение параметрической функции  $\Delta(\mu) = (\mu - \mu_0)\sqrt{n-1}/\sigma$  на положительную случайную величину  $\sigma/S$ , распределение которой не зависит от  $\mu$  и  $\sigma$ . При фиксированном  $\sigma$  функция  $\Delta(\mu)$  возрастает с ростом  $\mu$  и при этом все второе слагаемое возрастает, что влечет увеличение вероятности события перескока статистикой  $T$  порога  $C$ , то есть вероятности события  $T > C$ .

Итак, мы построили критерий проверки *односторонней* гипотезы  $\mu \leq \mu_0$  при односторонней альтернативе  $\mu > \mu_0$ . Функция мощности этого критерия зависит от  $\mu$  и  $\sigma$  только через параметрическую функцию  $\Delta = (\mu - \mu_0)\sqrt{n-1}/\sigma$ , которая называется *параметром нецентральности*. Распределение статистики  $T = (\bar{X} - \mu_0)\sqrt{n-1}/S$  при произвольных  $\mu$  и  $\sigma$ , через которое выражается функция мощности критерия Стьюдента, называется *нецентральным* распределением Стьюдента с  $n - 1$  степенью свободы; таблицы этого распределения, зависящего от параметра нецентральности  $\Delta$ , можно найти в ТМС.

Понятно, что построение критерия проверки простой гипотезы  $\mu = \mu_0$  при *двусторонней* (сложной) альтернативе  $\mu \neq \mu_0$  не вызывает принципиальных затруднений. Это критерий с критической областью  $|T| > C$ , где критическая константа  $C = C(\alpha) = S_{n-1}^{-1}(1 - \alpha/2)$ .

4<sup>0</sup>. СРАВНЕНИЕ СРЕДНИХ ЗНАЧЕНИЙ ДВУХ НОРМАЛЬНЫХ РАСПРЕДЕЛЕНИЙ С ОБЩЕЙ НЕИЗВЕСТНОЙ ДИСПЕРСИЕЙ (ДВУХВЫБОРОЧНЫЙ КРИТЕРИЙ СТЬЮДЕНТА). Пусть  $X$  и  $Y$  – независимые случайные величины, причем  $X \sim \mathcal{N}(\mu_1, \sigma^2)$ , а  $Y \sim \mathcal{N}(\mu_2, \sigma^2)$ , так что  $\mathbf{D}X = \mathbf{D}Y$ . По двум независимым выборкам  $X^{(n)} = (X_1, \dots, X_n)$  и  $Y^{(m)} = (Y_1, \dots, Y_m)$  (возможно, разного объема) требуется проверить гипотезу *однородности*  $H_0 : \mu_1 = \mu_2$  при альтернативе  $H_1 : \mu_1 > \mu_2$ . Типичный пример такой задачи – выявление эффекта нового метода лечения на группе из  $n$  пациентов посредством сравнения с контрольной группой из  $m$  пациентов, лечение которых проводится по старой методике.

Эта задача является для нас несколько новой, поскольку до сих пор мы

имели дело только с одной выборкой. Тем не менее, она сводится к той, что мы только что рассмотрели в  $\mathcal{Z}^0$ , с помощью следующих построений.

Рассмотрим сначала разность выборочных средних  $\bar{X} - \bar{Y}$ . Эта статистика имеет нормальное распределение со средним  $\mathbf{E}(\bar{X} - \bar{Y}) = \mu_1 - \mu_2$  и дисперсией  $\mathbf{D}(\bar{X} - \bar{Y}) = \mathbf{D}\bar{X} + \mathbf{D}\bar{Y} = \sigma^2(n^{-1} + m^{-1})$ . Следовательно, при справедливости нулевой гипотезы  $\mu_1 = \mu_2$  случайная величина

$$\xi = \frac{\bar{X} - \bar{Y}}{\sigma} \sqrt{\frac{nm}{n+m}}$$

имеет стандартное нормальное распределение  $\mathcal{N}(0, 1)$ . Далее, нормированные выборочные дисперсии  $nS_X^2/\sigma^2$  и  $mS_Y^2/\sigma^2$  независимы и распределены по закону хи-квадрат с  $n - 1$  и  $m - 1$  степенями свободы соответственно. Так как для хи-квадрат распределения, как частного случая гамма-распределения, имеет место теорема сложения, то случайная величина  $\eta = (nS_X^2 + mS_Y^2)/\sigma^2$  имеет хи-квадрат распределение с  $n + m - 2$  степенями свободы. Таким образом, мы приходим к двухвыборочной статистике Стьюдента

$$T_{n,m} = \frac{\xi}{\sqrt{\eta/(n+m-2)}} = \frac{\bar{X} - \bar{Y}}{\sqrt{nS_X^2 + mS_Y^2}} \sqrt{\frac{nm(n+m-2)}{n+m}},$$

распределение которой при справедливости нулевой гипотезы есть распределение Стьюдента с  $n + m - 2$  степенями свободы.

Как и в случае одновыборочного критерия Стьюдента в  $\mathcal{Z}^0$  нетрудно показать, что при любых фиксированных  $C$  и  $\sigma$  функция мощности двухвыборочного критерия Стьюдента  $T_{n,m} > C$  есть монотонно возрастающая функция параметра нецентральности  $\Delta = (\mu_1 - \mu_2)\sqrt{n+m-2}/\sigma$ , так что критическая константа  $C$  определяется по заданному уровню значимости из уравнения  $P(T_{n,m} > C) = 1 - S_{n+m-2}(C) = \alpha$  и равна квантили распределения Стьюдента:  $C(\alpha) = S_{n+m-2}^{-1}(1 - \alpha)$ . Понятно, что при альтернативе  $H_1: \mu_1 \neq \mu_2$  критическая константа  $C(\alpha) = S_{n+m-2}^{-1}(1 - \alpha/2)$ .

При использовании этого критерия следует обратить особое внимание на предположение о равенстве дисперсий наблюдаемых случайных величин:  $\sigma_X^2 = \sigma_Y^2$ . Задача сравнения средних двух нормальных распределений с неравными дисперсиями и с гарантированным ограничением  $\alpha$  на вероятность ошибки первого рода называется *проблемой Беренса–Фишера*. Известно лишь асимптотическое решение этой проблемы при больших  $n$  и  $m$ .

5<sup>0</sup>. СРАВНЕНИЕ ДИСПЕРСИЙ ДВУХ НОРМАЛЬНЫХ РАСПРЕДЕЛЕНИЙ ПРИ НЕИЗВЕСТНЫХ СРЕДНИХ (КРИТЕРИЙ ФИШЕРА). Независимые выборки  $X^{(n)}$  и  $Y^{(m)}$  берутся из соответствующих нормальных распределений  $\mathcal{N}(\mu_1, \sigma_1^2)$  и  $\mathcal{N}(\mu_2, \sigma_2^2)$ , относительно параметров которого проверяется гипотеза  $\sigma_1^2 = \sigma_2^2$  при альтернативе  $\sigma_1^2 > \sigma_2^2$  с мешающими параметрами  $\mu_1$  и  $\mu_2$ .

В этой задаче естественно рассмотреть критерий, основанный на статистике  $F = nS_X^2/mS_Y^2$ , которая распределена как

$$\frac{\chi_{n-1}^2}{\chi_{m-1}^2} \cdot \frac{\sigma_1^2}{\sigma_2^2}.$$

Функция мощности критерия  $F > C$  (который называется *критерием Фишера* или *F-критерием*)

$$m \left( \frac{\sigma_1^2}{\sigma_2^2} \right) = P_{\mu_1, \mu_2, \sigma_1, \sigma_2}(F > C) = P \left( \frac{\chi_{n-1}^2}{\chi_{m-1}^2} > C \cdot \frac{\sigma_2^2}{\sigma_1^2} \right)$$

есть монотонно возрастающая функция отношения дисперсий  $\sigma_1^2/\sigma_2^2$ . Для ее вычисления необходимо знать распределение отношения двух независимых случайных величин, распределенных по закону хи-квадрат с  $n - 1$  и  $m - 1$  степенями свободы. Это так называемое *распределение Фишера*  $F_{n-1, m-1}$ , плотность которого

$$f_{n-1, m-1}(x) = \frac{\Gamma \left( \frac{n+m-2}{2} \right)}{\Gamma \left( \frac{n-1}{2} \right) \Gamma \left( \frac{m-1}{2} \right)} \cdot \frac{x^{\frac{n-1}{2}-1}}{(x+1)^{\frac{n+m-2}{2}}}, \quad x > 0,$$

вычисляется столь же просто, как это мы делали при выводе распределения Стьюдента. Таблицы распределения Фишера можно найти в ТМС. Критическая константа  $C$  критерия Фишера заданного размера  $\alpha$  определяется как квантиль этого распределения:  $C(\alpha) = F_{n-1, m-1}^{-1}(1 - \alpha)$ .

Мы завершим иллюстрацию методов построения критериев с помощью состоятельных оценок тестируемого параметра примером, в котором не всегда размер критерия совпадает с заданным уровнем значимости.

6<sup>0</sup>. ПРОВЕРКА ГИПОТЕЗЫ О ВЕРОЯТНОСТИ УСПЕХА В ИСПЫТАНИЯХ БЕРНУЛЛИ. Рассмотрим задачу проверки гипотезы  $p = p_0$  против альтернативы  $p < p_0$  о вероятности  $p$  успешного исхода в испытаниях Бернулли. Пример такой задачи – проверка гипотезы о вероятности наследования доминантного признака в опытах Менделя, когда альтернативная модель

предписывает этой вероятности меньшее значение. Предлагаемый ниже метод решения позволяет строить критерии проверки такой гипотезы при альтернативах  $p > p_0$  или  $p \neq p_0$  посредством простой замены неравенства, определяющего критическую область, на обратное или двустороннее.

Итак, если мы располагаем выборкой  $X^{(n)}$  из двухточечного распределения  $B(1, p)$ , то относительная частота успешных испытаний (выборочное среднее)  $\bar{X}$  является несмещенной оценкой  $p$  с минимальной дисперсией. В соответствии с предложенной выше идеологией проверки гипотез с помощью оценок тестируемого параметра мы должны отвергать гипотезу  $p = p_0$  в пользу  $p < p_0$ , если  $\bar{X} - p_0 < C$ . Поскольку статистика  $T = n\bar{X} = \sum_1^n X_k$  имеет биномиальное распределение  $B(n, p)$ , а значение  $p_0$  задано, то для вычисления функции мощности удобнее записать критическую область в виде  $T < C$ . Но статистика  $T$  принимает только целочисленные значения  $0, 1, \dots, n$ , поэтому бессмысленно рассматривать дробные значения критических констант. Таким образом, мы приходим к наиболее удобной форме записи критической области в виде  $T < C$ , где  $C$  принимает значения  $1, 2, \dots, n$

Функция мощности такого критерия

$$m(p) = P_p(T < C) = \sum_{k=0}^{C-1} p^k (1-p)^{n-k},$$

и поскольку проверяется простая гипотеза, то критическая константа  $C$  должна определяться по заданному уровню значимости  $\alpha$  из неравенства

$$m(p_0) = \sum_{k=0}^{C-1} p_0^k (1-p_0)^{n-k} \leq \alpha. \quad (2)$$

Очевидно, что чем больше  $C$ , тем больше мощность критерия, и поэтому  $C(\alpha)$  следует выбирать как наибольшее целое число, удовлетворяющее неравенству (2). Размер критерия с таким  $C(\alpha)$  не обязательно равен  $\alpha$ , так что мы можем получить критерий уровня  $\alpha$ , но не размера  $\alpha$  (в предыдущих примерах с тестовыми статистиками, имеющими распределение непрерывного типа, мы имели критерии размера  $\alpha$ ). Более того, если  $p_0$  настолько мало, что  $(1-p_0)^n > \alpha$ , то не существует таких  $C$ , при которых имеет место неравенство (2). В таком случае мы должны принимать нулевую гипотезу при любом результате статистического эксперимента, обеспечивая тем самым нулевой размер такого критерия “уровня  $\alpha$ ”.

При больших объемах выборки  $n$  можно использовать нормальные аппроксимации биномиального распределения, получая таким образом критерий, размер которого асимптотически ( $n \rightarrow \infty$ ) равен  $\alpha$ . Статистика  $T$  асимптотически нормальна со средним  $np$  и дисперсией  $np(1-p)$ , поэтому неравенство (2) для определения критической константы имеет асимптотический аналог

$$\Phi \left( \frac{C - np_0}{\sqrt{np_0(1-p_0)}} \right) \leq \alpha,$$

откуда  $C(\alpha) \approx np_0 - \Phi^{-1}(1 - \alpha)/\sqrt{np_0(1-p_0)}$ . Легко понять, что такой метод построения критериев асимптотического уровня  $\alpha$  применим для любой критической области, в задании которой используется асимптотически нормальная оценка тестируемого параметра (см. пояснения в предыдущем параграфе перед пунктом 5<sup>0</sup>).

Этот пример показывает, что в случае дискретных распределений задача построения равномерно наиболее мощных критериев значительно усложняется, поскольку один из двух критериев одного и того же уровня  $\alpha$  может иметь большую мощность только потому, что он имеет больший размер. Мы столкнемся с этой проблемой в следующем параграфе, но следует заметить, что современная теория наиболее мощных критериев обходит этот неприятный момент за счет расширения понятия статистического правила, вводя так называемые *рандомизированные* критерии. К сожалению, я не располагаю временем познакомить вас с этим замечательным объектом теории статистического вывода.

Мы закончим этот параграф, как и было обещано, формулировкой *принципа двойственности* между задачами проверки гипотез и доверительного оценивания. Пусть  $A(\theta_0) \subset \mathcal{X}^n$  – область принятия некоторого критерия уровня  $\alpha$ , тестирующего гипотезу  $H_0 : \theta = \theta_0$ , и пусть  $A(\theta_0)$  определена при любом  $\theta_0 \in \Theta$ . Для каждого результата  $x^{(n)}$  наблюдения случайной выборки  $X^{(n)}$  введем подмножество  $\Delta(x^{(n)})$  параметрического пространства  $\Theta$ , положив  $\Delta(x^{(n)}) = \{\theta : x^{(n)} \in A(\theta)\}$ . Тогда  $\Delta(X^{(n)})$  есть  $(1 - \alpha)$ -доверительное множество для параметра  $\theta$ , поскольку  $P_\theta(\Delta(X^{(n)}) \ni \theta) = P_\theta(X^{(n)} \in A(\theta)) \geq 1 - \alpha$ .

Например, критерий Стьюдента проверки гипотезы  $\mu = \mu_0$  о среднем значении нормального  $(\mu, \sigma^2)$  распределения с неизвестной дисперсией  $\sigma^2$

имеет область принятия (см. п. 3<sup>0</sup> данного параграфа)

$$A(\mu_0) = \left\{ X^{(n)} : \frac{|\bar{X} - \mu_0|}{S} \sqrt{n-1} \leq S_{n-1}^{-1}(1 - \alpha/2) \right\}.$$

Подставим в это неравенство вместо фиксированного  $\mu_0$  параметр  $\mu$  и разрешим неравенство относительно  $\mu$ . В результате получим доверительное утверждение (см. п. 4<sup>0</sup> предыдущего параграфа)

$$\bar{X} - St_\alpha/\sqrt{n-1} \leq \mu \leq \bar{X} + St_\alpha/\sqrt{n-1},$$

в котором  $t_\alpha = S_{n-1}^{-1}(1 - \alpha/2)$ .

Вы сами можете сопоставить доверительные интервалы, построенные в §6, с критериями из §7. При этом сопоставлении можно вывести полезное правило, касающееся доверительной оценки скалярного параметра  $\theta$ . Если имеется состоятельный критерий проверки гипотезы  $\theta = \theta_0$  при двусторонней альтернативе  $\theta \neq \theta_0$ , то его области принятия соответствует двусторонний доверительный интервал. Если же альтернативная гипотеза носит односторонний характер, то при альтернативе  $\theta < \theta_0$  мы получаем верхнюю доверительную границу, а при  $\theta > \theta_0$  – нижнюю.

Естественно, принцип двойственности применим и к доверительным интервалам, как статистическим правилам проверки гипотез: гипотеза  $\theta \in \Theta_0$  отвергается тогда и только тогда, когда  $(1 - \alpha)$ -доверительная область принадлежит подмножеству  $\Theta_1$ , и такое статистическое правило (критерий) гарантирует заданное ограничение  $\alpha$  на вероятность ошибки первого рода.



## §8. Равномерно наиболее мощные критерии

Лекция 13

Метод построения критериев заданного уровня  $\alpha$ , который равномерно по всем альтернативным значениям параметра  $\theta$  максимизирует мощность критерия, существенно опирается на следующее, почти очевидное утверждение, которое в теории проверки гипотез обычно называется *леммой Неймана–Пирсона*.

Рассмотрим вероятностную модель, состоящую всего из двух распределений  $P_0$  и  $P_1$ , с общим носителем  $\mathcal{X}$  и функциями плотности  $f_0(x)$  и  $f_1(x)$ ,  $x \in \mathcal{X}$ . По выборке  $X^{(n)}$  проверяется простая гипотеза  $H_0$ : выборка взята из распределения  $P_0$  при простой альтернативе  $H_1$ : выборке соответствует распределение  $P_1$ . Определим критическую функцию  $\varphi^*(X^{(n)})$  как индикаторную функцию критической области

$$L(X^{(n)}) = \prod_{k=1}^n \frac{f_1(X_k)}{f_0(X_k)} > C.$$

Статистика  $L$  называется *статистикой отношения правдоподобия*, а критерий  $\varphi^*$  – *критерием отношения правдоподобия* или *критерием Неймана–Пирсона*. Критерий  $\varphi^*$  отвергает нулевую гипотезу, если правдоподобие альтернативы  $f_{1,n}(X^{(n)}) = \prod_1^n f_1(X_k)$  в  $C$  раз превосходит правдоподобие нулевой гипотезы  $f_{0,n}(X^{(n)}) = \prod_1^n f_0(X_k)$ . Этот критерий обладает следующим замечательным свойством.

**Теорема 8.1.** *Критерий отношения правдоподобия  $\varphi^*$  является наиболее мощным критерием в классе всех критериев проверки простой гипотезы при простой альтернативе, размер которых не превосходит размера критерия  $\varphi^*$ . Если критерий  $\varphi^*$  имеет размер  $\alpha$ , то он обладает наибольшей мощностью в классе всех критериев уровня  $\alpha$ .*

*Доказательство.* Пусть  $\varphi = \varphi(X^{(n)})$  – любой другой критерий, размер которого

$$\mathbf{E}_0 \varphi(X^{(n)}) \leq \mathbf{E}_0 \varphi^*(X^{(n)}). \quad (1)$$

Требуется показать, что тогда критерий  $\varphi^*$  имеет большую мощность, чем критерий  $\varphi$ , то есть  $\mathbf{E}_1 \varphi^*(X^{(n)}) \geq \mathbf{E}_1 \varphi(X^{(n)})$ .

Рассмотрим интеграл

$$\int_{\mathcal{X}^n} \left[ \varphi^*(x^{(n)}) - \varphi(x^{(n)}) \right] \left[ f_{1,n}(x^{(n)}) - C f_{0,n}(x^{(n)}) \right] d\mu_n(x^{(n)}) = \\ \mathbf{E}_1 \varphi^*(X^{(n)}) - \mathbf{E}_1 \varphi(X^{(n)}) - C \left[ \mathbf{E}_0 \varphi^*(X^{(n)}) - \mathbf{E}_0 \varphi(X^{(n)}) \right].$$

Достаточно показать, что этот интеграл неотрицателен, и тогда первое утверждение теоремы будет следовать из неравенства:

$$\mathbf{E}_1 \varphi^*(X^{(n)}) - \mathbf{E}_1 \varphi(X^{(n)}) - C \left[ \mathbf{E}_0 \varphi^*(X^{(n)}) - \mathbf{E}_0 \varphi(X^{(n)}) \right] \geq 0$$

которое влечет (см. (1))

$$\mathbf{E}_1 \varphi^*(X^{(n)}) - \mathbf{E}_1 \varphi(X^{(n)}) \geq C \left[ \mathbf{E}_0 \varphi^*(X^{(n)}) - \mathbf{E}_0 \varphi(X^{(n)}) \right] \geq 0.$$

Покажем, что функции  $\varphi^*(x^{(n)}) - \varphi(x^{(n)})$  и  $f_{1,n}(x^{(n)}) - C f_{0,n}(x^{(n)})$ , произведение которых интегрируется, одновременно положительны или отрицательны при любых  $x^{(n)} \in \mathcal{X}^n$ . Действительно, если  $\varphi^*(x^{(n)}) - \varphi(x^{(n)}) > 0$ , то это влечет  $\varphi^*(x^{(n)}) = 1$ , поскольку критическая функция равна единице, если она не равна нулю. Но, по определению критерия отношения правдоподобия, равенство  $\varphi^*(x^{(n)}) = 1$  возможно лишь в случае  $f_{1,n}(x^{(n)}) - C f_{0,n}(x^{(n)}) > 0$ . Точно также устанавливается, что неравенство  $\varphi^*(x^{(n)}) - \varphi(x^{(n)}) < 0$  влечет  $f_{1,n}(x^{(n)}) - C f_{0,n}(x^{(n)}) < 0$ .

Итак, критерий  $\varphi^*$  наиболее мощен в классе всех критериев, размер которых не превосходит размера  $\varphi^*$ . Если же  $\mathbf{E}_0 \varphi^*(X^{(n)}) = \alpha$ , то это утверждение, очевидно, влечет его наибольшую мощность в классе всех критериев уровня  $\alpha$ .

Применение этой теоремы к построению равномерно наиболее мощных критериев мы проиллюстрируем на одном частном примере, из которого будет виден общий подход к данной задаче.

**Пример 8.1. Проверка надежности при показательном распределении долговечности.** В примере 3.3 мы рассматривали проблему оценки надежности изделия с показательным распределением долговечности. Напомним, случайная величина  $X$ , реализация  $x$  которой соответствует промежутку времени от начала работы до момента отказа некоторого изделия, называется долговечностью, и по функции распределения  $F(x)$ ,  $x \geq 0$  случайной

величины  $X$  можно рассчитать надежность  $H(t)$  изделий, соответствующую гарантийному времени  $t$ :  $H(t) = P(X \geq t) = 1 - F(t)$ .

Пусть долговечность  $X$  распределена по показательному закону с функцией распределения  $F(x | \theta) = 1 - \exp\{-x/\theta\}$ , значение параметра  $\theta$  которой не известно. Мы должны удостовериться, что надежность выпускаемых изделий достаточно высока:  $H(t) \geq P_0$ , где  $P_0$  –наименьшая допустимая доля изделий, которые должны прослужить гарантийный срок  $t$ .

Это типичная задача проверки гипотез, решение которой начинается с определения нулевой гипотезы  $H_0$ . При этом следует помнить, что в статистическом критерии контролируется вероятность отклонения  $H_0$ , когда она в действительности верна. В нашей конкретной проблеме спецификация нулевой гипотезы во многом зависит от того, что повлечет за собой отказ изделия. Если мы выпускаем бытовые приборы, то отказ изделия до гарантийного срока  $t$  повлечет издержки на ремонт, которые могут быть незначительными по сравнению со стоимостью изделия. В таком случае естественно выбрать в качестве нулевой гипотезы утверждение о надежности изделий – отклонив эту гипотезу, когда она верна, мы потеряем дорогостоящую продукцию, ремонт которой нам обошелся бы значительно дешевле, чем ее уничтожение или продажа по бросовой цене. Если же отказ изделия приводит к катастрофическим последствиям, например, к гибели людей, то здесь рассуждать нечего, и за нулевую гипотезу следует брать утверждение о “ненадежности”. Отклонив такую гипотезу, когда она в действительности верна, мы столкнемся с неприемлемо большой долей отказов до истечения гарантийного срока, и поэтому риск от принятия “плохих” изделий должен быть контролируем. Остановимся на этом варианте и приступим к построению равномерно наиболее мощного критерия проверки гипотезы “ненадежности”  $H_0 : H(t) < P_0$  при альтернативе  $H_1 : H(t) \geq P_0$ , когда  $H(t) = \exp\{-t/\theta\}$ .

В терминах значений параметра  $\theta$  нулевая гипотеза принимает вид  $H_0 : \theta < \theta_0 = -t/\ln \alpha$ . Зафиксируем некоторое альтернативное значение  $\theta_1 > \theta_0$ , и рассмотрим задачу проверки простой гипотезы  $H'_0 : \theta = \theta_0$  при простой альтернативе  $H'_1 : \theta = \theta_1$ . Наиболее мощный критерий проверки простой гипотезы при простой альтернативе имеет критическую область вида (см. теорему 8.1)

$$L(X^{(n)}) = \prod_{k=1}^n \frac{f_1(X_k)}{f_0(X_k)} = \frac{\theta_0}{\theta_1} \exp \left\{ \left( \frac{1}{\theta_0} - \frac{1}{\theta_1} \right) \sum_1^n X_k \right\} > C,$$

где критическая константа  $C$  определяется по заданному уровню значимости  $\alpha$  из условия  $P_{\theta_0}(L(X^{(n)}) > C) \leq \alpha$ . Поскольку статистика  $T_n = \sum_1^n X_k$  имеет гамма-распределение  $G(n, \theta_0)$ , то для определения  $C$  в последнем неравенстве следует положить знак равенства. Кроме этого, статистика отношения правдоподобия  $L(X^{(n)})$  есть монотонная функция статистики  $T_n$ , поэтому критическую область  $L(X^{(n)}) > C$  можно записать в эквивалентной форме  $T_n > C$  и находить новое  $C$  из равенства  $P_{\theta_0}(T_n > C) = 1 - G_n(C/\theta_0) = \alpha$  (собственно говоря, нам все равно, какое  $C$  определять, но на практике, вне сомнения, удобнее иметь дело с критической областью  $T_n > C$ ).

Итак,  $C(\alpha) = \theta_0 \cdot G_n^{-1}(1 - \alpha)$ , где  $G_n^{-1}(\cdot)$  – квантиль стандартного гамма-распределения  $G(n, 1)$ , и критерий  $\varphi^*(X^{(n)}) = I_{\{T_n > C(\alpha)\}}(X^{(n)})$  заданного размера  $\alpha$  является наиболее мощным в классе всех критериев уровня  $\alpha$ , проверяющих гипотезу  $H_0'$  при альтернативе  $H_1'$ . Это означает, что для любого другого критерия  $\varphi$  с  $\mathbf{E}_{\theta_0} \varphi(X^{(n)}) \leq \alpha$  выполняется неравенство

$$\mathbf{E}_{\theta_1} \varphi(X^{(n)}) \leq \mathbf{E}_{\theta_1} \varphi^*(X^{(n)}). \quad (2)$$

Но критерий  $\varphi^*$  не зависит от выбора альтернативного значения  $\theta_1$  параметра  $\theta$  – критическая константа  $C(\alpha) = \theta_0 \cdot G_n^{-1}(1 - \alpha)$ ! Следовательно, неравенство (2) справедливо при любых  $\theta_1 > \theta_0$ , и мы приходим к заключению, что критерий  $\varphi^*$  есть равномерно наиболее мощный критерий в классе всех критериев уровня  $\alpha$ , проверяющих простую гипотезу  $H_0' : \theta = \theta_0$  при сложной альтернативе  $H_1 : \theta > \theta_0$ .

Далее, функция мощности критерия  $\varphi^*$ , как критерия различения исходных сложных гипотез  $H_0 : \theta < \theta_0$  и  $H_1 : \theta \geq \theta_0$ , равна  $m(\theta) = \mathbf{E}_{\theta} \varphi^*(X^{(n)}) = P_{\theta}(T_n > C(\alpha)) = 1 - G_n(G_n^{-1}(1 - \alpha)\theta_0/\theta)$ ,  $\theta > 0$ . Это – возрастающая функция  $\theta$ , поэтому максимум вероятности ошибки первого рода (размер критерия) равен  $m(\theta_0) = 1 - G_n(G_n^{-1}(1 - \alpha)) = \alpha$ . Таким образом, критерий  $\varphi^*$  есть критерий размера  $\alpha$  проверки гипотезы  $H_0$  при альтернативе  $H_1$ , обладающий равномерно наибольшей мощностью в классе всех критериев  $\varphi$  с ограничением  $\mathbf{E}_{\theta_0} \varphi(X^{(n)}) = \alpha$ . Но в таком случае он будет равномерно наиболее мощным и в более узком классе критериев уровня  $\alpha$ , то есть критериев  $\varphi$ , удовлетворяющих ограничению  $\mathbf{E}_{\theta} \varphi(X^{(n)}) \leq \alpha$  при любом  $\theta < \theta_0$ .

Более того, нетрудно убедиться, что критерий  $\varphi^*$  обладает минимальной вероятностью ошибки первого рода  $\alpha(\theta) = m(\theta)$ ,  $\theta \leq \theta_0$  в классе всех

критериев уровня  $\alpha$ . Для этого достаточно поменять местами нулевую гипотезу и альтернативу и выбрать уровень значимости, равный  $1 - \alpha$ .

В этом примере построение равномерно наиболее мощного критерия стало возможным благодаря особому свойству статистической структуры показательного распределения: *статистика  $L(X^{(n)})$  отношения правдоподобия есть монотонная функция статистики  $T_n = \sum_1^n X_k$* . Это – частный случай статистических структур, обладающих достаточной статистикой  $T$ , ибо в силу теоремы факторизации у таких структур  $L(X^{(n)}) = g_{\theta_1}(T)/g_{\theta_0}(T)$  зависит от  $X^{(n)}$  только через значения  $T(X^{(n)})$ . Дополнительное свойство монотонности отношения правдоподобия относительно  $T$  обеспечивает существование и возможность конструктивного построения равномерно наиболее мощного критерия, причем критическая область такого критерия обязательно имеет вид  $T > C$  или  $T < C$ . Например, критерий  $\sum_1^n X_k > C$  при соответствующем выборе  $C$  по заданному уровню значимости  $\alpha$  будет равномерно наиболее мощным критерием в классе всех критериев уровня  $\alpha$  проверки гипотезы  $\theta < \theta_0$  при альтернативе  $\theta \geq \theta_0$ , когда  $\theta$  есть среднее значение нормального распределения (дисперсия предполагается известной) или параметр масштаба гамма-распределения (параметр формы известен). Но если  $\theta$  – параметр таких распределений, как двухточечное или Пуассона, то критерий  $\varphi^*$  с критической областью  $\sum_1^n X_k > C$  обладает равномерно наибольшей мощностью только в классе тех критериев, размер которых не больше размера  $\varphi^*$ .

Другие критерии, которые мы рассматривали в предыдущем параграфе, также обладают свойством равномерной наибольшей мощности, и при доказательстве этого также используется лемма Неймана–Пирсона, но методика доказательства совершенно другая и требует разработки методов построения критериев, обладающих свойством инвариантности – независимости от мешающих параметров. Но это уже совсем другая область теории проверки гипотез, поговорить о которой у нас не хватает времени. Я лучше расскажу вам о некоторых дополнительных ухищрениях в практических применениях статистических критериев, которые позволяют с большей степенью наглядности оценить степень согласия проверяемой гипотезы с выборочными данными.

Все рассматриваемые нами критерии заданного уровня  $\alpha$  обладают тем свойством, что их критические области можно записать в виде  $T(X^{(n)}) > C(\alpha)$ , где  $T$  – некоторая статистика, характеризующая расхождение вы-

борочных данных с предполагаемыми значениями параметра. Увеличение уровня значимости  $\alpha$  приводит к уменьшению  $C(\alpha)$ , и мы получаем систему вложенных друг в друга критических областей. Это замечательное свойство наших критериев позволяет несколько изменить методологию их практического использования. До сих пор мы фиксировали уровень значимости  $\alpha$ , находили по нему критическую константу  $C(\alpha)$  и сравнивали ее с выборочным значением  $t = T(x^{(n)})$  статистики  $T = T(X^{(n)})$ . Поступим теперь следующим образом. Получив выборочные данные  $x^{(n)}$ , вычислим значение  $t = T(x^{(n)})$  и рассмотрим критерий  $T(X^{(n)}) > t$ . Размер такого критерия  $\alpha_{\text{кр.}} = P_0(T(X^{(n)}) > t)$  называется *критическим уровнем значимости*, который трактуется как вероятность получить столь же большие расхождения между выборочными данными и нулевой гипотезой, как и для выборочных данных  $x^{(n)}$ .

Естественно, мы по-прежнему можем работать с заданным уровнем значимости  $\alpha$ , отклоняя нулевую гипотезу, если  $\alpha_{\text{кр.}} < \alpha$ , и принимая ее в противном случае. Кстати, принимая гипотезу, не следует утверждать, что она верна. На этот счет существует более деликатное выражение: “выборочные данные согласуются с выдвинутой гипотезой,” ибо, как говорил один из создателей математической статистики сэр Д.Фишер, “гипотезы не проверяются, а разве лишь отвергаются”. Так вот, в свете этого высказывания более разумно просто сообщать полученный критический уровень значимости, сопровождая его следующим комментарием, который можно считать международным статистическим стандартом. Если  $\alpha_{\text{кр.}} \leq 0.01$ , то говорят, что расхождение между гипотезой и выборочными данными *высоко значимо*, если  $0.01 < \alpha_{\text{кр.}} \leq 0.05$ , то просто – *значимо*, если же  $0.05 < \alpha_{\text{кр.}} \leq 0.10$  – *почти значимо*, и в случае  $\alpha_{\text{кр.}} > 0.10$  – *не значимо*. Заметим также, что в некоторых применениях критериев значимости (особенно, в медицине)  $\alpha_{\text{кр.}}$  называют *достоверностью*. Существуют и другие, совершенно фантастические названия  $\alpha_{\text{кр.}}$ , которые я не буду здесь приводить в силу их крайне неприличного звучания.

Поговорим теперь об оптимальных свойствах доверительных границ, соответствующих равномерно наиболее мощным критериям. Рассмотрим только случай верхней доверительной границы  $\bar{\theta}_n = \bar{\theta}_n(X^{(n)})$ .

**Определение 8.1** Верхняя  $(1 - \alpha)$ -доверительная граница  $\bar{\theta}_n$  называется *равномерно наиболее точной*, если она равномерно по всем  $\theta$  и  $\theta'$ , удовлетворяющим неравенству  $\theta' > \theta$ , минимизирует вероятность  $P_\theta(\bar{\theta}_n(X^{(n)}) \geq \theta')$ .

Таким образом, в случае равномерно наиболее точной границы  $\bar{\theta}_n$  интервал  $(-\infty; \bar{\theta}_n]$  с заданной вероятностью  $1 - \alpha$  накрывает истинное значение параметра  $\theta$ , но он с минимальной вероятностью накрывает любые значения  $\theta$ , лежащие правее истинного.

Если мы проверяем гипотезу  $H : \theta = \theta_0$  при альтернативе  $K(\theta_0) : \theta < \theta_0$ , и область принятия  $A(\theta_0)$  равномерно наиболее мощного критерия размера  $\alpha$  обладает тем свойством, что подмножество  $\Delta_n(x^{(n)}) = \{\theta : x^{(n)} \in A(\theta)\}$  параметрического пространства  $\Theta \subseteq \mathbb{R}$  есть интервал  $(-\infty : \bar{\theta}_n(x^{(n)})]$ , то  $\bar{\theta}_n(X^{(n)})$  есть равномерно наиболее точная верхняя  $(1 - \alpha)$ -доверительная граница. Все объясняется довольно просто: вероятность  $P_\theta(\bar{\theta}_n(X^{(n)}) \geq \theta') = P_\theta(X^{(n)} \in A(\theta'))$  есть вероятность ошибки второго рода у критерия проверки гипотезы  $H : \theta = \theta'$  при альтернативе  $K(\theta') : \theta < \theta'$ . Равномерно наиболее мощный критерий естественно обладает равномерно минимальной вероятностью ошибки второго рода. Все построенные нами в §6 доверительные границы обладают оптимальными свойствами с точки зрения малой вероятности накрытия тех значений параметра, которые не соответствуют истине.

## §9. Проверка модельных предположений. Критерии согласия

Лекция 14

Рассмотренные нами методы построения оптимальных решающих функций в проблемах оценки параметров и проверки параметрических гипотез существенно опирались на такие особые свойства вероятностных моделей, как существование достаточных статистик, монотонность отношения правдоподобия относительно некоторой статистики, независимость выборок и прочее. Оценить же последствия от использования конкретных решающих функций (найти функцию риска статистического правила) вообще не представляется возможным без знания вероятностной модели. Отсюда возникает необходимость разработки общих методов тестирования (проверки) предлагаемой вероятностной модели  $\mathcal{P} = \{P_\theta, \theta \in \Theta\}$  по данным случайной выборки, или нескольких выборок, которые предположительно извлекаются из некоторых распределений семейства  $\mathcal{P}$ . Значимые расхождения между модельными и эмпирическими распределениями вынуждают статистика пересмотреть посылки, положенные в основу построения вероятностной модели, и тем самым избежать больших потерь от использования заведомо плохих решающих правил (точнее правил, которые оптимальны не для той модели).

Понятно, что речь идет о проверке статистических гипотез без особой спецификации альтернатив к нулевой гипотезе. Статистические правила проверки модельных предположений обычно называются *критериями согласия*, и в математической статистике сложился некоторый традиционный набор таких критериев, обладающих большой универсальностью. Это критерии, с помощью которых можно проверять не только принадлежность распределения наблюдаемой случайной величины к определенному семейству, но и тестировать некоторые более “грубые” черты модели, как то независимость компонент наблюдаемого случайного вектора (векторной случайной величины), возможность объединения нескольких выборок в одну (проверка гипотезы однородности выборок) и множество других предположений, касающихся структуры выборочных данных. Мы познакомимся в этом параграфе с набором универсальных статистических процедур, объединяемых общим названием *критерии хи-квадрат*. Об одном из них мы уже упоминали в §2 в связи с построением гистограммы выборки; это –



1<sup>0</sup>. КРИТЕРИЙ СОГЛАСИЯ ХИ-КВАДРАТ. Решается статистическая проблема проверки гипотезы о виде распределения наблюдаемой случайной величины  $X$  (возможно, векторной). Начнем с простейшего случая, когда построение вероятностной модели привело к полной спецификации распределения, то есть проблема состоит в проверке простой гипотезы  $H$ : распределение  $X$  на измеримом пространстве  $(X, \mathcal{A})$  ее значений есть  $P(A)$ ,  $A \in \mathcal{A}$ .

Построение критерия согласия выборочных данных с распределением  $P$  начинается с разбиения пространства  $X$  на  $r \geq 2$  частей  $A_1, \dots, A_r$ ;  $X = \sum_{i=1}^r A_i$ . Рекомендации по выбору числа  $r$  и способу разбиения носят довольно расплывчатый характер, и если не уточнять возможные альтернативы к  $P$ , то, как вы сами понимаете, таких рекомендаций не может быть в принципе. Главное, разбиение не должно определяться выборочными значениями, надо стремиться к областям одинаковой конфигурации и размера, не следует делать слишком подробное разбиение. Например, если  $X = \mathbb{R}$  (наблюдается действительная случайная величина), то прямая  $\mathbb{R}$  разбивается на  $r$  интервалов вида  $(-\infty, a]$ ,  $(a, a + \Delta]$ ,  $(a + \Delta, a + 2\Delta]$ ,  $\dots$ ,  $(a + (r - 3)\Delta, a + (r - 2)\Delta]$ ,  $(a + (r - 2)\Delta, +\infty)$ , так что длина внутренних интервалов постоянна и равна  $\Delta$ . Конечно выбор  $r$  зависит от объема выборки  $n$ , но даже при исключительно больших  $n$  не делается более 15-20 разбиений; этого вполне достаточно, чтобы в гистограмме отразить всю специфику формы тестируемого распределения.

После разбиения  $X$  проводится сортировка выборочных данных по областям разбиений и подсчитываются количества  $\nu_1, \dots, \nu_r$ ,  $\sum_{i=1}^r \nu_i = n$ , данных, попавших в соответствующие области  $A_1, \dots, A_r$ . Вычисляются “теоретические” вероятности  $p_i = P(A_i)$ ,  $i = 1, \dots, r$  попадания выборочных данных в эти области и вычисляется значение  $x^2$  тестовой статистики

$$X^2 = \sum_{i=1}^r \frac{(\nu_i - np_i)^2}{np_i}.$$

Гипотеза  $H$  отвергается, если  $x^2 > C$ , где критическая константа  $C$  выбирается по заданному уровню значимости  $\alpha$  как наименьшее число, удовлетворяющее неравенству  $P(X^2 > C) \leq \alpha$ . Естественно, на практике используют критический уровень значимости  $\alpha_{кр.} = P(X^2 > x^2)$ , сопровождая его комментариями типа тех, которые были приведены в предыдущем параграфе после введения понятия критического уровня значимости. Однако точное распределение статистики  $X^2$  найти в явном виде не представляется

возможным; предельное распределение  $X^2$  при  $n \rightarrow \infty$  установил К.Пирсон в самом начале XX века.

**Теорема 9.1.** *Если число разбиений  $r \geq 2$  фиксировано, а объем выборки  $n \rightarrow \infty$ , то распределение  $X^2$  сходится к распределению хи-квадрат с  $r - 1$  степенью свободы.*

Доказательство.. Очевидно, для вывода предельного распределения  $X^2$  следует в первую очередь обратиться к совместному распределению частот  $\nu_1, \dots, \nu_r$ ,  $\sum_1^n \nu_i = n$ . Это мультиномиальное распределение  $\mathcal{M}(r, n, p)$ , (см. §9 курса ТВ) с функцией плотности

$$f(x_1, \dots, x_r) = P(\nu_1 = x_1, \dots, \nu_r = x_r) = \frac{n!}{x_1! \cdots x_r!} p_1^{x_1} \cdots p_r^{x_r},$$

сосредоточенное на целочисленной решетке  $\sum_1^n x_i = n$ . Теорема 9.1 из курса ТВ утверждает, что совместное распределение первых  $r - 1$  частот  $\nu_1, \dots, \nu_{r-1}$  аппроксимируется  $r - 1$ -мерным нормальным распределением. Естественно, предельное распределение всего вектора частот  $\nu_1, \dots, \nu_r$  при соответствующей нормировке на их средние значения и стандартные отклонения будет вырожденным, ибо  $\sum_1^n \nu_i = n$ . Вырожденные распределения лучше всего исследовать с помощью характеристических функций, ибо такие распределения можно записать в явном виде, только переходя к системе координат на той гиперповерхности, где сосредоточено такое распределение, и это чрезвычайно усложняет технику асимптотического анализа распределений. Итак, найдем совместную характеристическую функцию  $\nu_1, \dots, \nu_r$ .

Вспомним схему мультиномиальных испытаний. Мы наблюдаем выборку  $Y_1, \dots, Y_n$  из распределения случайного вектора  $Y = (X_1, \dots, X_r)$ , все компоненты которого, за исключением одной (скажем,  $X_j$ ), могут принимать только нулевые значения, в то время как  $X_j = 1$ . Каждая компонента  $Y_i$  выборки  $Y^{(n)} = (Y_1, \dots, Y_n)$  есть независимая копия  $Y$ , так что  $Y_i = (X_{1i}, \dots, X_{ri})$  и  $X_{ji}$  – копия (в смысле одинаковости распределения)  $X_j$ ,  $j = 1, \dots, r$ ,  $i = 1, \dots, n$ . В таких обозначениях

$$\nu_j = \sum_{i=1}^n X_{ji}, \quad j = 1, \dots, r.$$

Если мы найдем характеристическую функцию  $\varphi_Y(\mathbf{t})$ ,  $\mathbf{t} = (t_1, \dots, t_r)$ , наблюдаемого вектора  $Y$ , то характеристическая функция  $\varphi_\nu(\mathbf{t})$  вектора

частот  $\nu = (\nu_1, \dots, \nu_r)$  будет вычисляться по формуле  $\varphi_\nu(\mathbf{t}) = \varphi_Y^n(\mathbf{t})$ , ибо характеристическая функция суммы независимых случайных величин равна произведению характеристических функций слагаемых (пункт 3<sup>0</sup> теоремы 12.1 курса ТВ). Но характеристическая функция вектора  $Y$  (напомним,  $\sum_1^r X_j = 1$ )

$$\varphi_Y(\mathbf{t}) = \mathbf{E} \exp \left\{ \mathbf{i} \sum_1^r t_j X_j \right\} = \sum_1^r p_j e^{\mathbf{i} t_j},$$

и поэтому

$$\varphi_\nu(\mathbf{t}) = \left( \sum_1^r p_j e^{\mathbf{i} t_j} \right)^n.$$

Теперь перейдем к асимптотическому анализу характеристической функции вектора  $X$  нормированных частот

$$X_j = \frac{\nu_j - np_j}{\sqrt{np_j}}, \quad j = 1, \dots, r,$$

сумма квадратов компонент которого составляет тестовую статистику  $X^2$  (извините, что использую букву  $X$  в новом смысле, но не хочется вводить для обозначения случайных величин новые символы). Характеристическая функция случайного вектора, компоненты которого подвергнуты линейному преобразованию, вычисляется по формуле, аналогичной пункту 2<sup>0</sup> теоремы 12.1:

$$\varphi_X(\mathbf{t}) = \exp \left\{ -\mathbf{i} \sum_1^r t_j \sqrt{np_j} \right\} \left( \sum_1^r p_j \exp \left\{ \frac{\mathbf{i} t_j}{\sqrt{np_j}} \right\} \right)^n.$$

Разложим логарифм этой функции в ряд Маклорена по степеням  $t_1, \dots, t_r$ , как это делалось при доказательстве центральной предельной теоремы:

$$\begin{aligned} \ln \varphi_X(\mathbf{t}) &= -\mathbf{i} \sqrt{n} \sum_1^r t_j \sqrt{p_j} + \\ & n \ln \left[ 1 + \frac{\mathbf{i}}{\sqrt{n}} \sum_1^r t_j \sqrt{p_j} - \frac{1}{2n} \sum_1^r t_j^2 + O(n^{-3/2}) \right] = \\ & = -\frac{1}{2} \sum_1^r t_j^2 + \frac{1}{2} \left( \sum_1^r t_j \sqrt{p_j} \right)^2 + O(n^{-1/2}). \end{aligned}$$

Таким образом, характеристическая функция предельного распределения вектора  $X$  нормированных частот есть

$$\lim_{n \rightarrow \infty} \varphi_X(\mathbf{t}) = \exp \left\{ -\frac{1}{2} \left[ \sum_1^r t_j^2 - \left( \sum_1^r t_j \sqrt{p_j} \right)^2 \right] \right\}.$$

Это – характеристическая функция  $r$ -мерного нормального распределения с нулевыми средними и матрицей ковариаций  $\Lambda = \mathbf{I} - \mathbf{p}\mathbf{p}'$ , где  $\mathbf{I}$  – единичная матрица, а  $\mathbf{p} = (\sqrt{p_1}, \dots, \sqrt{p_r})$  – вектор столбец.

Рассмотрим квадратичную форму

$$Q(\mathbf{t}) = \sum_1^r t_j^2 - \left( \sum_1^r t_j \sqrt{p_j} \right)^2,$$

коэффициенты которой определяют ковариации компонент вектора  $Z = (Z_1, \dots, Z_r)$ , распределенного по нормальному закону. Если произвести ортогональное преобразование  $\mathbf{A}$  вектора  $\mathbf{t}$ , полагая  $\mathbf{u} = \mathbf{A}\mathbf{t}$  и фиксируя последнюю строку матрицы  $\mathbf{A}$  таким образом, чтобы в новом векторе  $\mathbf{u} = (u_1, \dots, u_r)$  компонента  $u_r = \sum_1^r t_j \sqrt{p_j}$ , то мы получим квадратичную форму (вспомните аналогичные ортогональные преобразования нормального вектора при выводе распределения выборочной дисперсии в лемме Фишера)

$$Q(\mathbf{t}) = \sum_1^r t_j^2 - \left( \sum_1^r t_j \sqrt{p_j} \right)^2 = \sum_1^r u_j^2 - u_r^2 = \sum_1^{r-1} u_j^2.$$

Таким образом, существует ортогональное преобразование  $Y = \mathbf{B}Z$  вектора  $Z$ , после которого  $Y_1, \dots, Y_{r-1}$  независимы и одинаково нормально распределены со средними, равными нулю, и единичными дисперсиями, а  $Y_r$  имеет нулевое среднее и нулевую дисперсию, то есть  $Y_r = 0$  почти наверное. Все это, конечно, следствие вырожденности нормального распределения вектора  $Z$  – оно сосредоточено на гиперплоскости  $\sum_1^r Z_j \sqrt{p_j} = 0$ .

Изучим теперь предельное распределение статистики  $X^2 = \sum_1^r X_j^2$ . Поскольку предельное распределение вектора  $X$  совпадает с распределением вектора  $Z$ , то предельное распределение статистики  $X^2$  определяется распределением квадратичной формы  $\sum_1^r Z_j^2$ . Как известно, ортогональные преобразования не меняют суммы квадратов, поэтому  $\sum_1^r Z_j^2 = \sum_1^r Y_j^2 =$

$\sum_1^{r-1} Y_j^2$ . Следовательно, предельное распределение статистики  $X^2$  есть распределение суммы квадратов  $r - 1$  независимых случайных величин, имеющих общее стандартное нормальное распределение. По определению это – хи-квадрат распределение с  $r - 1$  степенями свободы. Теорема Пирсона доказана.

### Лекция 15

Рассмотрим теперь более сложную статистическую проблему, в которой проверяется гипотеза о принадлежности распределения  $P$  наблюдаемой случайной величины некоторому параметрическому семейству  $\mathcal{P} = \{P_\theta, \theta \in \Theta \subseteq \mathbb{R}^s\}$ , индексированному  $s$ -мерным параметром  $\theta = (\theta_1, \dots, \theta_s)$ . В таком случае

$$X^2(\theta) = \sum_{i=1}^r \frac{(\nu_i - np_i(\theta))^2}{np_i(\theta)}$$

не может называться статистикой и ее нельзя использовать для проверки сложной гипотезы  $H : P \in \mathcal{P}$ . Естественно воспользоваться какой-либо оценкой  $\hat{\theta}_n = \hat{\theta}_n(X^{(n)})$  параметра  $\theta$  и рассмотреть тестовую статистику

$$\hat{X}^2 = X^2(\hat{\theta}_n) = \sum_{i=1}^r \frac{(\nu_i - np_i(\hat{\theta}_n))^2}{np_i(\hat{\theta}_n)}.$$

Понятно, что распределение статистики  $\hat{X}^2$  может зависеть от метода оценки параметра  $\theta$ . Однако, если определить оценку  $\hat{\theta}_n$  из условия минимума случайной функции  $X^2(\theta)$ , то, как показал Фишер, при определенных условиях регулярности, которым удовлетворяют все рассмотренные нами в курсе ТВ вероятностные модели, *предельное распределение статистики  $\hat{X}^2$  есть хи-квадрат распределение с  $r - s - 1$  степенями свободы*. Если же  $\hat{\theta}_n$  – оценка  $\theta$  по методу максимального правдоподобия, то предельное распределение  $\hat{X}^2$ , также при условиях регулярности типа тех, что обеспечивали асимптотическую нормальность  $\hat{\theta}_n$ , имеет функцию распределения  $K(x)$ , для которой справедлива двусторонняя оценка

$$K_{r-1}(x) \leq K(x) \leq K_{r-s-1}(x),$$

при любом  $x > 0$ .

Доказательство этих утверждений достаточно громоздко и мы не будем им заниматься из-за недостатка времени. Идейная сторона проблемы

нам ясна, и коль скоро нам сообщили распределение тестовой статистики, то мы можем использовать его для расчета критического уровня значимости. В случае оценки максимального правдоподобия, когда мы располагаем двусторонней оценкой  $\alpha_{кр.}$ , рекомендуется при отклонении гипотезы ориентироваться на  $\alpha_{кр.} = 1 - K_{r-1}(x^2)$  ( $> 1 - K_{r-s-1}(x^2)$ ), а в случае ее принятия – на  $\alpha_{кр.} = 1 - K_{r-s-1}(x^2)$  ( $< 1 - K_{r-1}(x^2)$ ), чтобы уменьшить риск от принятия неправильного решения.

Критерий хи-квадрат является наиболее универсальным статистическим методом тестирования вероятностной модели, поскольку предельное распределение статистики не зависит от распределения наблюдаемой случайной величины даже в том случае, когда это распределение зависит от параметров, значение которых неизвестно. Критерий Колмогорова  $\sqrt{n}D_n = \sqrt{n} \sup_x |F_n(x) - F(x)| > C$ , о котором говорилось в начале §2, можно использовать только для проверки простой гипотезы  $F(\cdot) = F_0(\cdot)$  о виде функции распределения. Если  $F_0(x|\theta)$  зависит от параметра  $\theta$  и в статистику  $\sqrt{n}D_n$  вместо  $F(x)$  подставляется  $F_0(x|\hat{\theta}_n(X^{(n)}))$ , то распределение модифицированной таким образом статистики  $\sqrt{n}D_n$  зависит как от вида функции  $F_0$ , так и от параметра  $\theta$ . Существует, правда, несколько случаев особой связи между  $x$  и  $\theta$  в записи функции  $F_0$ , при наличии которой распределение тестовой статистики не зависит от  $\theta$ . Это, например, такие функции распределения с параметрами масштаба и сдвига, как нормальное и показательное. Для тестирования таких распределений составляются специальные таблицы критических констант и критических уровней значимости. Следует заметить, что прямое использование критерия Колмогорова с оценками неизвестных значений параметров является наиболее распространенной ошибкой в практических приложениях методов тестирования вероятностных моделей.

Обратимся теперь к проверке гипотез, касающихся не столько вида распределения наблюдаемых случайных величин, сколько их особых свойств, наличие которых позволяет значительно упростить вероятностную модель и добиться ее более четкой спецификации.

2<sup>0</sup>. КРИТЕРИЙ НЕЗАВИСИМОСТИ ХИ-КВАДРАТ (ТАБЛИЦЫ СОПРЯЖЕННОСТИ ПРИЗНАКОВ). Следующая задача выявления зависимости между определенными признаками наблюдаемых объектов часто возникает в практических приложениях математической статистики. Предположим, что мы случайно выбрали  $n$  особей из некоторой этнической популяции, и хо-

тим выяснить, существует ли зависимость между цветом волос и цветом глаз. Мы различаем  $s \geq 2$  уровней первого признака (например, блондин, брюнет, шатен и рыжий) и  $r \geq 2$  уровней второго (например, карие, серые, голубые и зеленые). Все  $n$  особей разбиваются на  $sr$  групп в соответствии с наличием тех или иных уровней каждого признака, и составляется следующая таблица частот особей в каждой группе.

Признаки	1	2	$\dots$	$s$	Сумма
1	$\nu_{11}$	$\nu_{12}$	$\dots$	$\nu_{1s}$	$\nu_{1\cdot}$
2	$\nu_{21}$	$\nu_{22}$	$\dots$	$\nu_{2s}$	$\nu_{2\cdot}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$r$	$\nu_{r1}$	$\nu_{r2}$	$\dots$	$\nu_{rs}$	$\nu_{r\cdot}$
Сумма	$\nu_{\cdot 1}$	$\nu_{\cdot 2}$	$\dots$	$\nu_{\cdot s}$	$n$

Такие таблицы, в которых суммы

$$\nu_{i\cdot} = \sum_{j=1}^s \nu_{ij}, \quad \nu_{\cdot j} = \sum_{i=1}^r \nu_{ij},$$

называются *таблицами сопряженности признаков*. Требуется проверить нулевую гипотезу о том, что переменные признаки, по которым построена таблица, независимы. Построим вероятностную модель, соответствующую такого рода табличным данным и составим статистику  $X^2$  для проверки гипотезы независимости.

Пусть  $p_{ij}$  – вероятность того, что случайно отобранная особь имеет  $i$ -ый уровень по первому признаку и  $j$ -ый – по второму,  $i = 1, \dots, r$ ,  $j = 1, \dots, s$ . Гипотеза независимости означает, что  $p_{ij} = p_{i\cdot} \cdot p_{\cdot j}$ , где

$$p_{i\cdot} = \sum_{j=1}^s p_{ij}, \quad p_{\cdot j} = \sum_{i=1}^r p_{ij}$$

при любых  $i = 1, \dots, r$  и  $j = 1, \dots, s$ . Для проверки гипотезы независимости предлагается использовать тестовую статистику

$$X^2 = \sum_{i,j} \frac{(\nu_{ij} - n p_{i\cdot} \cdot p_{\cdot j})^2}{n p_{i\cdot} \cdot p_{\cdot j}}, \quad (1)$$

в которой суммирование распространяется на все  $rs$  групп таблицы сопряженности признаков. Понятно, что  $X^2$  является тестовой статистикой

только в случае известных значений  $r + s - 2$  параметров  $p_{i\cdot}$  и  $p_{\cdot j}$ ,  $i = 1, \dots, r$ ,  $j = 1, \dots, s$  (напомним,  $\sum_1^r p_{i\cdot} = \sum_1^s p_{\cdot j} = 1$ , так что с помощью этих соотношений два из  $r + s$  параметров, например,  $p_{r\cdot}$  и  $p_{\cdot s}$ , можно выразить через остальные  $r + s - 2$  параметров). В этом случае  $X^2$  имеет в пределе ( $n \rightarrow \infty$ ) хи-квадрат распределение с  $rs - 1$  степенями свободы.

Конечно, вся проблема состоит в том, что эти параметры неизвестны. Оказывается, оценки максимального правдоподобия

$$\hat{p}_{i\cdot} = \frac{\nu_{i\cdot}}{n}, \quad \hat{p}_{\cdot j} = \frac{\nu_{\cdot j}}{n}, \quad i = 1, \dots, r, \quad j = 1, \dots, s,$$

этих параметров асимптотически эквивалентны оценкам по методу минимума статистики  $X^2$ , и поэтому подстановка в правую часть (1) этих оценок приводит к статистике

$$\hat{X}^2 = n \sum_{i,j} \frac{(\nu_{ij} - \nu_{i\cdot} \nu_{\cdot j} / n)^2}{\nu_{i\cdot} \nu_{\cdot j}} = n \left( \sum_{i,j} \frac{\nu_{ij}^2}{\nu_{i\cdot} \nu_{\cdot j}} - 1 \right),$$

предельное распределение которой есть хи-квадрат с  $rs - (r + s - 2) - 1 = (r - 1)(s - 1)$  степенями свободы.

Естественно, статистику  $\hat{X}^2$  можно использовать для проверки независимости компонент двумерного вектора  $(X, Y)$ , и при этом таблица сопряженности представляет частотные данные для построения гистограммы двумерной выборки  $(X_1, Y_1), \dots, (X_n, Y_n)$ . Соответствующим образом нормированная статистика  $X^2$  может служить мерой зависимости признаков (или компонент  $X$  и  $Y$  случайного вектора).

**3<sup>0</sup>. КРИТЕРИЙ ОДНОРОДНОСТИ ХИ-КВАДРАТ.** Анализируются данные  $s \geq 2$  независимых мультиномиальных схем испытаний с одинаковым числом  $r \geq 2$  возможных исходов и соответствующими объемами  $n_1, \dots, n_s$  наблюдений в каждой схеме. Проверяется гипотеза *однородности*: все схемы испытаний имеют одинаковый вектор вероятностей  $\mathbf{p} = (p_1, \dots, p_r)$ ,  $\sum_1^r p_i = 1$ , появления соответствующих исходов, причем значения компонент вектора  $\mathbf{p}$  не известны. Обозначая  $\nu_{ij}$  частоту появления  $i$ -го исхода в  $j$ -ом испытании, представим данные наблюдений в виде таблицы, аналогичной таблице сопряженности признаков



исх. \ схем.	1	2	...	s	Сумма
1	$\nu_{11}$	$\nu_{12}$	...	$\nu_{1s}$	$\nu_{1.}$
2	$\nu_{21}$	$\nu_{22}$	...	$\nu_{2s}$	$\nu_{2.}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
r	$\nu_{r1}$	$\nu_{r2}$	...	$\nu_{rs}$	$\nu_{r.}$
Сумма	$n_1$	$n_2$	...	$n_s$	$n$

Составим сначала статистику хи-квадрат для случая известного вектора вероятностей  $\mathbf{p}$ :

$$X^2 = \sum_{j=1}^s \sum_{i=1}^r \frac{(\nu_{ij} - n_j p_i)^2}{n_j p_i}.$$

Внутренняя сумма

$$X_j^2 = \sum_{i=1}^r \frac{(\nu_{ij} - n_j p_i)^2}{n_j p_i}$$

представляет статистику хи-квадрат для  $j$ -ой схемы мультиномиальных испытаний, и поэтому имеет в пределе ( $n_j \rightarrow \infty$ ) хи-квадрат распределение с  $r - 1$  степенями свободы. Статистика  $X^2$  есть сумма  $s$  независимых статистик, каждая из которых имеет предельное хи-квадрат распределение, так что, в силу теоремы сложения, предельное распределение  $X^2$  есть хи-квадрат распределение с  $(r - 1)s$  степенями свободы.

В случае неизвестных значений вероятностей исходов, которые при справедливости нулевой гипотезы одинаковы для всех схем испытаний, используем их оценки  $\hat{p}_i = \nu_{i.}/n$ ,  $i = 1, \dots, r$ , (всего оценивается  $r - 1$  параметр). Подстановка этих оценок в  $X^2$  дает статистику

$$\hat{X}^2 = n \sum_{i,j} \frac{(\nu_{ij} - n_j \nu_{i.}/n)^2}{n_j \nu_{i.}} = n \left( \sum_{i,j} \frac{\nu_{ij}^2}{n_j \nu_{i.}} - 1 \right),$$

предельное распределение которой есть хи-квадрат с  $(r - 1)s - (r - 1) = (r - 1)(s - 1)$  степенями свободы. Замечателен тот факт, что мы получили тестовую статистику такого же вида и с тем же предельным распределением, что и при проверке гипотезы независимости признаков.

Естественно, построенный критерий можно использовать для проверки гипотезы однородности распределений, из которых извлекаются  $s \geq 2$  выборки. Выборочные данные при этом подвергаются группировке в соответствии с одинаковым для всех выборок разбиением пространства  $X$  на  $r \geq 2$  областей.

## Л И Т Е Р А Т У Р А

1. Большев Л.Н., Смирнов Н.В. Таблицы математической статистики.– М.: Наука, 1983.
2. Боровков А.А. Математическая статистика.– М.: Наука, 1984.
3. Боровков А.А. Теория вероятностей.– М.: Наука, 1986.
4. Козлов М.В., Прохоров А.В. Введение в математическую статистику.– М.: Изд-во МГУ, 1987.
5. Крамер Г. Математические методы статистики.– М.: Мир, 1975.
6. Чистяков В.П. Курс теории вероятностей.– М.: Наука, 1982.
7. Ширяев А.Н. Вероятность.– М.: Наука, 1980.