

# ДИСПЕРСИОННЫЙ АНАЛИЗ

Методические разработки  
по специальному курсу

## Часть I

ОЦЕНКИ МЕТОДА НАИМЕНЬШИХ КВАДРАТОВ.  
КРИТЕРИЙ ФИШЕРА. ОДНОФАКТОРНЫЙ  
ДИСПЕРСИОННЫЙ АНАЛИЗ.

# П Е Ч А Т А Е Т С Я

ПО РЕШЕНИЮ СЕКЦИИ НАУЧНО-МЕТОДИЧЕСКОГО СОВЕТА

КАЗАНСКОГО УНИВЕРСИТЕТА

С о с т а в и т е л ь

доцент С.В.Симушкин

*Дисперсионный анализ (ДА)* представляет собой набор статистических методов, предназначенных для обработки экспериментальных данных, зависящих от большого количества одновременно действующих факторов. Целью анализа является оценка влияния факторов на результирующий признак и выявление наиболее значимых из них. В качестве примера применения методов ДА можно привести классическую задачу сравнения урожайности нескольких сортов пшеницы, высеянных на участках земли с различным химическим составом, с применением различных типов удобрений. Методы ДА позволяют ответить на вопросы о зависимости урожайности от каждого из трех рассматриваемых факторов (сорт, земля, удобрение), о возможном взаимодействии факторов и на ряд других вопросов.

Существенным моментом теории ДА является представление регрессии отклика статистического эксперимента в виде линейной комбинации функций факторов с неизвестными коэффициентами. В первых двух главах данного курса рассматривается оценка по методу наименьших квадратов неизвестных коэффициентов и выводится ее распределение в предположении нормальности распределения ошибки. В последующих главах строится критерий отношения правдоподобий для проверки различных гипотез относительно уравнения регрессии. Рассматриваются задачи одно-, двух- и многофакторного ДА. Наконец, в последних главах предлагаются методы построения планов статистического эксперимента, оптимизирующие некоторые характеристики дисперсии оценок коэффициентов регрессии.

# Оглавление

I	Оценки влияния факторов	<b>4</b>
§ 1	Классическая модель линейной регрессии . . . . .	4
§ 2	Оценки метода наименьших квадратов . . . . .	6
§ 3	Функции, допускающие оценку. Теорема Гаусса-Маркова . . . . .	9
§ 4	Каноническая форма основных предположений . . . . .	11
II	Дисперсионный анализ нормальных совокупностей	<b>14</b>
§ 1	Распределение оценок и ошибок . . . . .	14
§ 2	Доверительные эллипсоиды . . . . .	15
§ 3	Критерий отношения правдоподобий. Критерий Фишера . . . . .	17
III	Модели дисперсионного анализа	<b>22</b>
§ 1	Однофакторный дисперсионный анализ . . . . .	22
§ 2	Сравнения . . . . .	25
§ 3	Непараметрические критерии однородности . . . . .	29
§ 4	Полный двухфакторный анализ . . . . .	37
§ 5	Полный многофакторный анализ с взаимодействиями . . . . .	52
IV	Модели со случайными факторами	<b>58</b>
§ 1	Однофакторный дисперсионный анализ . . . . .	58
§ 2	Полная классификация по двум признакам . . . . .	61
V	Факторные планы	<b>65</b>
§ 1	Дробные реплики полного факторного плана . . . . .	65
§ 2	Латинские планы . . . . .	69
§ 3	Блочные схемы . . . . .	71
VI	Оптимальные планы регрессионных экспериментов	<b>73</b>
§ 1	Планы эксперимента и их информационная матрица . . . . .	74
§ 2	Критерии оптимальности . . . . .	78
§ 3	Теоремы эквивалентности . . . . .	81
§ 4	Полиномиальная и тригонометрическая регрессии . . . . .	88
§ 5	Оптимальные планы первого порядка . . . . .	89

# Г л а в а I

## ОЦЕНКИ ВЛИЯНИЯ ФАКТОРОВ

### § 1. Классическая модель линейной регрессии

Классическая модель регрессии предполагает, что учитываемые в эксперименте факторы  $x_1, \dots, x_p$  оказывают линейное воздействие на отклик – результат эксперимента  $y$  :

$$y = x_1\beta_1 + x_2\beta_2 + \dots + x_p\beta_p + \varepsilon,$$

где  $\beta_1, \dots, \beta_p$  – неизвестные параметры, в той или иной мере отражающие степень воздействия факторов,  $\varepsilon$  – случайная ошибка эксперимента,  $x_1, \dots, x_p$  – вектор известных коэффициентов, задаваемых обычно исследователем. Чаще всего эти коэффициенты равны 1 или 0 в зависимости от того, присутствует соответствующий фактор в эксперименте или нет. Иногда, однако, они могут принимать и непрерывный ряд значений, выражая степень предпочтения того или иного фактора. Целью *ДД* является получение выводов относительно значений неизвестных параметров и их комбинаций, а также относительно распределения ошибки  $\varepsilon$ .

Если было произведено  $n (> p)$  экспериментов, то результаты измерений могут быть представлены в матричном виде

$$\vec{y} = \mathbf{X}' \vec{\beta} + \vec{\varepsilon},$$

где  $\vec{y} = (y_1, \dots, y_n)'$  – вектор-столбец наблюдаемых значений,  $\vec{\beta} = (\beta_1, \dots, \beta_p)'$  – столбец неизвестных параметров,  $\vec{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)'$  – вектор-столбец ненаблюдаемых ошибок,

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \dots & \dots & \dots & \dots \\ x_{p1} & x_{p2} & \dots & x_{pn} \end{pmatrix} -$$

матрица коэффициентов, задаваемых экспериментатором.

**О п р е д е л е н и е 1.** Матрица  $\mathbf{X}$  называется *матрицей плана*.

Различают три типа моделей *Дд*.

*Модель I* – модель с постоянными факторами, в которой все параметры  $\{\beta_j\}$  постоянны. Иными словами, для этой модели факторы могут разве лишь присутствовать или отсутствовать в эксперименте, но их влияние всюду, где они есть, одинаково.

Если некоторый фактор  $\beta_j$  постоянный и входные коэффициенты  $x_{ji} = 1$  для всех  $i = 1, \dots, n$  (то есть фактор присутствует во всех экспериментах), тогда  $\beta_j$  называется *аддитивным фактором*. Легко понять, что такой фактор представляет собой средний ожидаемый отклик в предположении, что все остальные неаддитивные факторы не влияют на результат эксперимента. Поэтому аддитивный фактор может быть только один и его можно интерпретировать как среднее значение случайной составляющей модели  $\varepsilon$ . Этим, в частности, объясняется введенное ниже условие  $\mathbf{E} \varepsilon = 0$  на ошибку эксперимента.

*Модель II* – модель со случайными факторами, где все параметры  $\{\beta_j\}$  случайны, кроме, может быть, одного, который является аддитивным.

*Модель III* – смешанная модель.

Рассмотрим сначала ситуацию, описываемую моделью I *Дд*.

Статистические выводы в значительной степени зависят от предположений относительно распределения ошибок  $\vec{\varepsilon}$ . Минимальные требования, которым должно подчиняться это распределение, заключаются, во-первых, в отсутствии систематической ошибки –

$$\mathbf{E} \varepsilon_i = 0, \quad i = 1, \dots, n,$$

во-вторых, в некоррелируемости ошибок в различных экспериментах –

$$\mathbf{E} \varepsilon_j \varepsilon_k = 0 \quad \text{при } j \neq k,$$

и, в-третьих, в одинаковой распределенности ошибок. Последнее условие обычно сводится к требованию отсутствия изменчивости дисперсии ошибок от эксперимента к эксперименту:  $\mathbf{D} \varepsilon_i = \mathbf{E} \varepsilon_i^2 = \sigma^2, \quad i = 1, \dots, n$ .

Таким образом, основные предположения *Дд* можно записать в виде

$$\Omega: \quad \begin{cases} \vec{y} = \mathbf{X}'\vec{\beta} + \vec{\varepsilon}, \\ \mathbf{E} \vec{\varepsilon} = \vec{0}, \quad \text{Cov}(\vec{\varepsilon}) = \mathbf{E} \vec{\varepsilon} \vec{\varepsilon}' = \sigma^2 \mathbf{I}. \end{cases} \quad (\text{I.1})$$

**З а м е ч а н и е i.** Из основных предположений следует, что

$$\mathbf{E} \vec{y} = \mathbf{X}'\beta \quad \text{и} \quad \text{Cov}(\vec{y}) = \sigma^2 \mathbf{I}. \quad (\text{I.2})$$

## § 2. Оценки метода наименьших квадратов

Мера качества оценок  $\vec{b}$  факторных параметров  $\vec{\beta}$  должна, естественно, зависеть от расхождений  $e_i = y_i - (x_{1i}b_1 + \dots + x_{pi}b_p)$ ,  $i = 1, \dots, n$ , между наблюдаемыми значениями отклика  $\vec{y}$  и ожидаемыми значениями  $\mathbf{X}'\vec{b}$ . На практике наиболее популярна мера качества, равная сумме квадратов ошибок –

$$\mathcal{G}(\vec{b}, \vec{y}) = \sum_{i=1}^n e_i^2 = \|\vec{y} - \mathbf{X}'\vec{b}\|^2,$$

и называемая обычно средне-квадратической ошибкой (СКО).

**О п р е д е л е н и е 2.** Вектор  $\vec{\beta}^*$ , для которого СКО

$$\|\vec{y} - \mathbf{X}'\vec{\beta}^*\|^2 = \min_{\vec{b} \in \mathcal{R}^p} \|\vec{y} - \mathbf{X}'\vec{b}\|^2,$$

называется *оценкой по методу наименьших квадратов* (ОМНК). Величина СКО

$$SS_e = \|\vec{y} - \mathbf{X}'\vec{\beta}^*\|^2 \quad (\text{I.3})$$

называется *суммой квадратов ошибок*. Иногда, для того, чтобы подчеркнуть, что минимум ошибки найден при основных предположениях, мы будем обозначать  $SS_e$  через  $\mathcal{G}_\Omega$ .

Покажем, что ОМНК всегда существует.

Если ОМНК существует, то, как известно из курса математического анализа, все производные СКО по компонентам вектора  $\vec{b}$  будут равны нулю в точке  $\vec{b} = \vec{\beta}^*$ . Произведя элементарные алгебраические выкладки, легко показать, что если ОМНК  $\vec{\beta}^*$  существует, то она обязана удовлетворять системе линейных уравнений

$$\sum_{i=1}^n \sum_{j=1}^p x_{ki}x_{ji}\vec{\beta}_i^* = \sum_{i=1}^n x_{ki}y_i, \quad k = \overline{1, p}, \quad (\text{I.4})$$

или в матричной форме

$$\mathcal{S}\vec{\beta}^* = \mathbf{X}\vec{y}, \quad (\text{I.5})$$

где матрица  $\mathcal{S} = \mathbf{X}\mathbf{X}'$  – так называемая *информационная матрица* плана.

**О п р е д е л е н и е 3.** Уравнения (I.4) и (I.5) называются *нормальными уравнениями*.

Дальнейшие рассуждения будут существенно опираться на следующую геометрическую интерпретацию.

Обозначим через  $\vec{\xi}_j$   $j$ -ый столбец матрицы  $\mathbf{X}'$  (– столбец коэффициентов при  $j$ -ом факторе). Если  $\text{rang } \mathbf{X} = r (\leq p)$ , то линейное пространство, порожденное системой векторов  $\{\vec{\xi}_j\}_{j=1}^p$  (– то есть множество линейных комбинаций  $c_1 \vec{\xi}_1 + \dots + c_p \vec{\xi}_p$ ), будет иметь размерность  $r$ . Обозначим это пространство  $V_r$ . Тогда, очевидно,  $\mathbf{X}' \vec{b} = b_1 \vec{\xi}_1 + \dots + b_p \vec{\xi}_p \in V_r$  и  $\|\vec{y} - \mathbf{X}' \vec{b}\|$  равно расстоянию от вектора  $\vec{y}$  до элемента  $\mathbf{X}' \vec{b}$  пространства  $V_r$ .

Мы часто будем пользоваться критерием перпендикулярности векторов к пространству  $V_r$ .

**Л е м м а I.1.** Вектор  $\vec{a} \perp V_r$  тогда и только тогда, когда

$$\vec{a}' \mathbf{X}' = \vec{0}' \quad \text{или} \quad \mathbf{X} \vec{a} = \vec{0}.$$

*Доказательство.* Вектор  $\vec{a} \perp V_r$  только в том случае, если он перпендикулярен всем векторам  $\vec{\xi}_j$ , порождающим  $V_r$ :  $\vec{a}' \vec{\xi}_j = 0$ ,  $j = \overline{1, p}$ . Поскольку  $\vec{\xi}_j$  образуют столбцы матрицы  $\mathbf{X}'$ , то отсюда получаем, что перпендикулярность  $\vec{a} \perp V_r$  эквивалентна соотношению

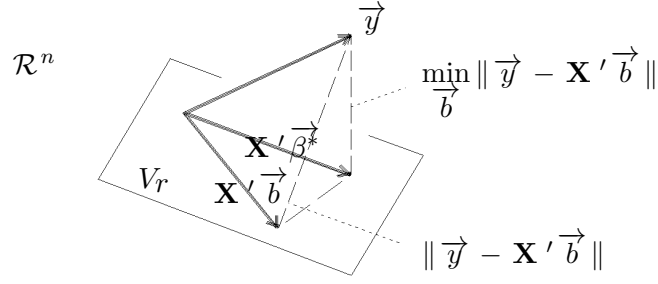
$$\vec{a}' \mathbf{X}' = \vec{a}' (\vec{\xi}_1, \dots, \vec{\xi}_p) = (\vec{a}' \vec{\xi}_1, \dots, \vec{a}' \vec{\xi}_p) = \vec{0}'.$$

Второе равенство получается из первого простым транспонированием.

⊗

**Т е о р е м а I.1.** ОМНК всегда существует и удовлетворяет нормальным уравнениям. Любое решение нормальных уравнений будет ОМНК.

*Доказательство.* Вектор  $\vec{y}$  принадлежит пространству  $\mathcal{R}^n \supset V_r$ . Как известно, проекция  $\text{Proj}_{V_r} \vec{y}$  вектора  $\vec{y}$  на пространство  $V_r$  является единственным элементом из  $V_r$ , минимизирующим расстояние до вектора  $\vec{y}$ . Поскольку  $\text{Proj}_{V_r} \vec{y} \in V_r$ , то существует такой вектор  $\vec{\beta}^*$ , что  $\text{Proj}_{V_r} \vec{y} = \mathbf{X}' \vec{\beta}^*$  и, следовательно,  $\vec{\beta}^*$  есть ОМНК (см.рис.).



Пусть теперь  $\vec{\beta}^*$  – некоторая ОМНК. Тогда  $\mathbf{X}'\vec{\beta}^* = \text{Proj}_{V_r} \vec{y}$ , т.к. проекция единственна. Следовательно,

$$\mathbf{X}'\vec{\beta}^* - \vec{y} \perp V_r.$$

В силу леммы I.1 это возможно только тогда, когда

$$\mathbf{X}(\mathbf{X}'\vec{\beta}^* - \vec{y}) = \vec{0},$$

что эквивалентно системе нормальных уравнений (I.5).  $\otimes$

**Замечание ii.** Несмотря на то, что проекция  $\text{Proj}_{V_r} \vec{y}$  определяется единственным образом, ОМНК может быть не единственной. Следующее утверждение содержит условие единственности ОМНК и, кроме того, здесь вычисляется ее ковариационная матрица.

**Следствие 1.** Если  $\text{rang } \mathbf{X} = p$ , то

- а) существует единственная ОМНК  $\vec{\beta}^* = \mathcal{S}^{-1} \mathbf{X}' \vec{y}$ ;
- б)  $\vec{\beta}^*$  есть линейная несмещенная оценка вектора  $\vec{\beta}$ ;
- в) ковариационная матрица  $\text{Cov}(\vec{\beta}^*) = \sigma^2 \mathcal{S}^{-1}$ .

*Доказательство.* а) Из теории матриц известно, что  $\text{rang } \mathcal{S} = \text{rang } \mathbf{X} \mathbf{X}' = \text{rang } \mathbf{X} = p$ . Следовательно, матрица  $\mathcal{S}$  невырождена и имеет обратную матрицу. Первый пункт следствия вытекает теперь из вида нормальных уравнений (I.5).

б) Так как среднее значение  $\mathbf{E} \vec{y} = \mathbf{X}' \vec{\beta}$ , то математическое ожидание вектора  $\vec{\beta}^*$  равно

$$\mathbf{E} \vec{\beta}^* = \mathbf{E} \mathcal{S}^{-1} \mathbf{X}' \vec{y} = \mathcal{S}^{-1} \mathbf{X}' \mathbf{E} \vec{y} = \mathcal{S}^{-1} \mathbf{X}' \mathbf{X}' \vec{\beta} = \mathcal{S}^{-1} \mathcal{S} \vec{\beta} = \vec{\beta}.$$

в) ОМНК  $\vec{\beta}^* = \mathcal{S}^{-1} \mathbf{X}' \vec{y}$  есть линейная комбинация вектора  $\vec{y}$  с ковариационной матрицей  $\text{Cov}(\vec{y}) = \sigma^2 \mathbf{I}$ , поэтому

$$\text{Cov}(\vec{\beta}^*) = \mathcal{S}^{-1} \mathbf{X}' \text{Cov}(\vec{y}) (\mathcal{S}^{-1} \mathbf{X}')' = \sigma^2 \mathcal{S}^{-1} \mathbf{X}' \mathbf{X}' (\mathcal{S}^{-1})'.$$



Очевидно матрица  $\mathcal{S}$  симметрична и  $\mathbf{X} \mathbf{X}' = \mathcal{S}$ , следовательно,

$$\text{cov}(\vec{\beta}^*) = \sigma^2 \mathcal{S}^{-1} \mathcal{S} \mathcal{S}^{-1} = \sigma^2 \mathcal{S}^{-1}.$$

что и требовалось доказать.  $\otimes$

### § 3. Функции, допускающие оценку.

#### Теорема Гаусса-Маркова

В реальных статистических задачах зачастую требуется определить не сами значения параметров  $\vec{\beta}$ , входящих в уравнение регрессии, а либо оценить некоторую функцию от этих параметров, либо проверить ряд гипотез относительно  $\vec{\beta}$ . Как та, так и другая проблема формализуется обычно с помощью линейных функций вектора  $\vec{\beta}$ . Здесь будет описан класс линейных функций, для которых (в любом случае) существует единственная ОМНК и которые применяются для целей ДД.

**Определение 4.** 1) Линейная функция от  $\vec{\beta}$  с известными коэффициентами

$$\psi(\vec{\beta}) = \sum_{j=1}^p c_j \beta_j = \vec{c}' \vec{\beta}$$

называется *параметрической функцией*.

2) Говорят, что параметрическая функция  $\psi(\vec{\beta})$  *допускает оценку*, если существует такой вектор  $\vec{a} = (a_1, \dots, a_p)'$ , что

$$\mathbf{E} \vec{a}' \vec{y} = \psi(\vec{\beta}). \quad (\text{I.6})$$

**Замечание iii.** Условие I.6 означает, по-существу, что для параметрической функции  $\psi$  существует линейная несмещенная оценка  $\hat{\psi} = \vec{a}' \vec{y}$ .

**Теорема I.2.** Функция  $\psi = \vec{c}' \vec{\beta}$  допускает оценку тогда и только тогда, когда вектор-строка  $\vec{c}'$  является линейной комбинацией строк матрицы  $\mathbf{X}'$ , т.е. когда существует вектор  $\vec{a}$  такой, что  $\vec{c}' = \vec{a}' \mathbf{X}'$ .

*Доказательство.* Функция  $\psi$  допускает оценку только тогда, когда существует такой вектор  $\vec{a}$ , для которого при любом  $\vec{\beta}$

$$\vec{c}' \vec{\beta} = \mathbf{E} \vec{a}' \vec{y} = \vec{a}' \mathbf{X}' \vec{\beta},$$

что равносильно тождеству  $\vec{c}' = \vec{a}'\mathbf{X}'$ .  $\otimes$

Найдем теперь наилучшую из возможных оценок параметрической функции. Сравнение оценок будем производить по их дисперсии.

**Определение 5.** Оценку  $\hat{\psi} = \vec{a}'\vec{y}$  назовем МНК-оценкой для параметрической функции  $\psi$ , если она имеет наименьшую дисперсию среди всех линейных несмещенных оценок  $\psi$ .

**Теорема I.3. (Гаусса-Маркова)** Для любой параметрической функции  $\psi = \vec{c}'\vec{\beta}$ , допускающей оценку, существует единственная МНК-оценка. Значение этой оценки может быть вычислено по формуле  $\hat{\psi} = \vec{c}'\vec{\beta}^*$ , где  $\vec{\beta}^*$  – любая ОМНК  $\vec{\beta}$ .

*Доказательство.* Заметим сначала, что дисперсия любой линейной оценки  $\vec{a}'\vec{y}$  пропорциональна квадрату длины вектора  $\vec{a}$ :

$$\begin{aligned} \mathbf{D} \vec{a}'\vec{y} &= \text{Cov}(\vec{a}'\vec{y}) = \vec{a}' \text{Cov}(\vec{y}) \vec{a} = \\ &= \sigma^2 \vec{a}'\vec{a} = \sigma^2 \|\vec{a}\|^2. \end{aligned}$$

Пусть теперь  $\vec{a}'\vec{y}$  – несмещенная оценка  $\psi$  и  $\vec{a}^* = \text{Proj}_{V_r} \vec{a}$  – проекция вектора  $\vec{a}$  на подпространство  $V_r$ . Тогда  $\vec{a} = \vec{a}^* + \vec{b}$ , где  $\vec{b} = \vec{a} - \vec{a}^*$ , причем  $\vec{b} \perp V_r$ . Следовательно, (см. лемму I.1)  $\vec{b}'\mathbf{X}' = \vec{0}'$ . Это равенство позволяет доказать, что статистика  $\vec{a}^*\vec{y}$  также образует несмещенную оценку параметрической функции  $\psi$ :

$$\begin{aligned} \mathbf{E} \vec{a}^*\vec{y} &= \mathbf{E}(\vec{a} - \vec{b})'\vec{y} = \mathbf{E} \vec{a}'\vec{y} - \mathbf{E} \vec{b}'\vec{y} = \\ &= \psi - \vec{b}'\mathbf{E} \vec{y} = \psi - \vec{b}'\mathbf{X}'\vec{\beta} = \psi. \end{aligned}$$

Легко видеть, однако, что дисперсия  $\vec{a}^*\vec{y}$  меньше дисперсии оценки  $\vec{a}'\vec{y}$ . Действительно, в силу теоремы «Пифагора»

$$\mathbf{D} \vec{a}'\vec{y} = \sigma^2 \|\vec{a}\|^2 = \sigma^2 (\|\vec{a}^*\|^2 + \|\vec{b}\|^2) \leq \sigma^2 \|\vec{a}^*\|^2 = \mathbf{D} \vec{a}^*\vec{y},$$

причем знак равенства достигается лишь в том случае, если  $\vec{b} = \vec{0}$ , т.е., когда вектор  $\vec{a} \in V_r$ .

Покажем теперь, что для любой несмещенной оценки  $\vec{q}'\vec{y}$  функции  $\psi$  проекции  $\text{Proj}_{V_r} \vec{q}$  будут одинаковыми. В самом деле, так

как рассматриваемые оценки несмещены ( $\mathbf{E} \vec{a}^* \vec{y} = \psi = \mathbf{E} \vec{q}^* \vec{y}$ ), то  $\forall \vec{\beta}$

$$0 = \mathbf{E} \vec{a}^* \vec{y} - \mathbf{E} \vec{q}^* \vec{y} = (\vec{a}^* - \vec{q}^*)' \mathbf{X}' \vec{\beta}.$$

Отсюда следует, что  $(\vec{a}^* - \vec{q}^*)' \mathbf{X}' = \vec{0}$ , то есть  $(\vec{a}^* - \vec{q}^*) \perp V_r$ . С другой стороны, очевидно,  $(\vec{a}^* - \vec{q}^*) \in V_r$ . Это возможно только в том случае, если  $\vec{a}^* = \vec{q}^*$ , что завершает доказательство первой части теоремы.

Осталось доказать, что значение МНК-оценки  $\hat{\psi}$  может быть получено подстановкой ОМНК  $\vec{\beta}^*$  в уравнение для  $\psi$ :  $\vec{a}^* \vec{y} = \vec{c}' \vec{\beta}^*$ .

При доказательстве теоремы I.1 мы показали, что для ОМНК  $\vec{\beta}^*$  вектор  $\mathbf{X}' \vec{\beta}^* = \text{Proj}_{V_r} \vec{y}$ . Так как  $\vec{a}^* \in V_r$ , то  $(\vec{y} - \mathbf{X}' \vec{\beta}^*) \perp \vec{a}^*$  и, следовательно,

$$\vec{a}^* \vec{y} = \vec{a}^* \mathbf{X}' \vec{\beta}^*.$$

С другой стороны, так как  $\vec{a}^* \vec{y}$  есть несмещенная оценка  $\psi$ , то

$$\psi = \vec{c}' \vec{\beta} = \mathbf{E} \vec{a}^* \vec{y} = \vec{a}^* \mathbf{X}' \vec{\beta}, \quad \forall \vec{\beta}.$$

Подставляя в это соотношение оценку  $\vec{\beta}^*$  и сравнивая с предыдущим соотношением, получаем доказательство теоремы.  $\otimes$

## § 4. Каноническая форма основных предположений

Вывод распределений оценок параметров и параметрических функций в значительной степени упрощается, если от переменных  $\vec{y}$  перейти к так называемым каноническим переменным. Для этого определим ортонормированный базис  $(\vec{\alpha}_1, \dots, \vec{\alpha}_r, \vec{\alpha}_{r+1}, \dots, \vec{\alpha}_n)$  во всём пространстве  $\mathcal{R}^n$  таким образом, чтобы его первые  $r$  векторов образовывали базис в пространстве  $V_r$ . Теперь любой вектор  $\vec{y}$  из  $\mathcal{R}^n$  может быть записан в виде

$$\vec{y} = \sum_{i=1}^n z_i \vec{\alpha}_i,$$

где  $z_i = \vec{\alpha}_i' \vec{y}$ ,  $i = \overline{1, n}$ , — координаты  $\vec{y}$  при разложении по «осям»  $(\vec{\alpha}_1, \dots, \vec{\alpha}_n)$ .

**Определение 6.** Координаты  $\vec{z} = (z_1, \dots, z_n)$  разложения наблюдаемого вектора  $\vec{y}$  по базису  $(\vec{\alpha}_1, \dots, \vec{\alpha}_n)$ , где первые  $r$

базисных векторов совпадают с базисом пространства  $V_r$ , называются *каноническими переменными*.

Образует ортогональную матрицу  $P' = (\vec{\alpha}_1, \dots, \vec{\alpha}_n)$ , столбцы которой совпадают с векторами построенного базиса. Тогда  $PP' = P'P = \mathbf{I}$ , и связь между  $\vec{y}$  и  $\vec{z}$  может быть записана в матричном виде

$$\vec{z} = P \vec{y}, \quad \vec{y} = P' \vec{z}.$$

Найдем среднее значение и матрицу ковариаций канонических переменных.

Обозначим  $\zeta_i = \mathbf{E} z_i$ ,  $i = \overline{1, n}$ , – средние значения  $\vec{z}$ . Так как при  $i > r$  базисные векторы  $\vec{\alpha}_i \perp V_r$ , то для этих  $i > r$  произведение  $\vec{\alpha}_i' \mathbf{X}' = \vec{0}'$  и средние значения

$$\zeta_i = \mathbf{E} z_i = \vec{\alpha}_i' \mathbf{E} \vec{y} = \vec{\alpha}_i' \mathbf{X}' \vec{\beta} = 0.$$

Ковариационная матрица

$$\text{Cov}(\vec{z}) = \text{Cov}(P \vec{y}) = P \text{Cov}(\vec{y}) P' = \sigma^2 P P' = \sigma^2 \mathbf{I}$$

в силу ортогональности  $P$ .

Таким образом, основные предположения ДД могут быть переписаны в *канонической форме*:

$$\Omega : \begin{cases} \vec{z} = \vec{\zeta} + \vec{\varepsilon} \\ \mathbf{E} \vec{\varepsilon} = \vec{0}, \quad \zeta_i = 0, \quad i > r, \quad \text{Cov}(\vec{z}) = \sigma^2 \mathbf{I}. \end{cases} \quad (\text{I.7})$$

### Пространство ошибок

Оказывается, что в канонических переменных сумма квадратов ошибок  $SS_e$  (I.3) записывается через последние  $n - r$  координат  $\vec{z}$ . Действительно, так как

$$\vec{y} = \sum_{i=1}^n z_i \vec{\alpha}_i, \quad \text{то} \quad \text{Proj}_{V_r} \vec{y} = \sum_{i=1}^r z_i \vec{\alpha}_i.$$

Следовательно,

$$SS_e = \|\vec{y} - \text{Proj}_{V_r} \vec{y}\|^2 = \left\| \sum_{i=r+1}^n z_i \vec{\alpha}_i \right\|^2 = \sum_{i=r+1}^n z_i^2. \quad (\text{I.8})$$

Это обстоятельство служит одним из оснований для введения следующего определения.

**О п р е д е л е н и е 7.** Пространством ошибок называется линейное пространство, порожденное  $\{z_{r+1}, \dots, z_n\}$ .

### Пространство оценок

**О п р е д е л е н и е 8.** Пространство, порожденное  $\{z_1, \dots, z_r\}$ , называется пространством оценок.

Введение этого пространства связано с тем, что для параметрической функции  $\psi = \vec{c}'\vec{\beta}$ , допускающей оценку, по теореме Гаусса-Маркова МНК-оценка имеет вид  $\hat{\psi} = \vec{a}'\vec{y}$ , где  $\vec{a} \in V_r$ , т.е.  $\vec{a}'\vec{\alpha}_j = 0$  для всех  $j > r$ . Следовательно,

$$\begin{aligned}\hat{\psi} &= \vec{a}'\vec{y} = \vec{a}'P'\vec{z} = (\vec{a}'\vec{\alpha}_1, \dots, \vec{a}'\vec{\alpha}_n)\vec{z} = \\ &= (\vec{a}'\vec{\alpha}_1, \dots, \vec{a}'\vec{\alpha}_r, 0, \dots, 0)\vec{z} = \sum_{i=1}^r c_i z_i,\end{aligned}$$

где  $c_i = \vec{a}'\vec{\alpha}_i$ ,  $i = 1, \dots, r$ .

Таким образом, МНК-оценка любой параметрической функции, допускающей оценку, есть линейная комбинация только первых  $r$  канонических переменных  $(z_1, \dots, z_r)$ .

### Несмещенная оценка $\sigma^2$

Хорошей иллюстрацией использования канонических переменных является вычисление математического ожидания  $SS_e$  и построение несмещенной оценки  $\sigma^2$ .

#### Т е о р е м а I.4. Статистика

$$\mathfrak{S}^2 = \frac{SS_e}{n-r}$$

есть несмещенная оценка  $\sigma^2$ .

*Доказательство.* В силу (I.8)

$$\mathbf{E} SS_e = \sum_{i=r+1}^n \mathbf{E} z_i^2.$$

Кроме того, по основным предположениям (I.7),  $\mathbf{E} z_i^2 = \mathbf{D} z_i = \sigma^2$  при  $i > r$ , что, очевидно, доказывает теорему.  $\otimes$

## Г л а в а П

### ДИСПЕРСИОННЫЙ АНАЛИЗ НОРМАЛЬНЫХ СОВОКУПНОСТЕЙ

Чтобы построить статистические выводы относительно неизвестных параметров  $\vec{\beta}$ , необходимо конкретизировать вид распределения ошибок в основных предположениях  $\mathcal{D}_A$ . Наиболее продвинутые результаты в  $\mathcal{D}_A$  получены в предположении, что ошибки имеют нормальное распределение. Таким образом, основные предположения  $\mathcal{D}_A$  переписываются в следующем виде:

$$\Omega : \vec{y} \sim \mathcal{N}_n(X' \vec{\beta}, \sigma^2 \mathbf{I}). \quad (\text{II.1})$$

Канонические переменные, как линейные комбинации  $\vec{y}$ , также будут распределены нормально с той же ковариационной матрицей  $\sigma^2 \mathbf{I}$  и нулевыми средними  $\zeta_{r+1} = 0, \dots, \zeta_p = 0$ :

$$\Omega : \vec{z} \sim \mathcal{N}_n(\vec{\zeta}, \sigma^2 \mathbf{I}), \quad \zeta_i = 0, i > r. \quad (\text{II.2})$$

### § 1. Распределение оценок и ошибок

Найдем сначала распределение МНК-оценок параметрических функций, допускающих оценку, и распределение суммы квадратов ошибок  $SS_e$ .

Рассмотрим  $q$  параметрических функций, допускающих оценку,

$$\psi_i = \vec{c}_i' \vec{\beta} = \sum_{j=1}^p c_{ij} \beta_j, \quad i = 1, \dots, q.$$

Пусть  $\psi_i^* = \vec{a}_i' \vec{y}$  – МНК-оценка функции  $\psi_i$ , где вектор  $\vec{a}_i \in V_r$ . Обозначим через  $C$  и  $A$  матрицы, строки которых совпадают с вектор-строками  $\vec{c}_i'$  и  $\vec{a}_i'$ ,  $i = 1, \dots, q$ , соответственно. Тогда можно записать

$$\vec{\psi} = (\psi_1, \dots, \psi_q)' = C \vec{\beta} \quad \text{и} \quad \vec{\psi}^* = (\psi_1^*, \dots, \psi_q^*)' = A \vec{y}.$$

**Теорема II.1.** Если  $q$  линейно независимых параметрических функций  $\vec{\psi}$  допускают оценку, тогда

- 1) МНК-оценки  $\vec{\psi}^*$  не зависят от суммы квадратов ошибок  $SS_e$ ;
- 2)  $\vec{\psi}^* \sim \mathcal{N}_q(\vec{\psi}, \sigma^2 B)$  с некоторой известной матрицей  $B > 0$ ;
- 3)  $\frac{SS_e}{\sigma^2} \sim \chi_{n-r}^2$ .

*Доказательство.* 1) Оценки  $\vec{\psi}^*$  принадлежат пространству оценок, а  $SS_e$  есть функция переменных пространства ошибок, другими словами, они зависят от  $(z_1, \dots, z_r)$  и  $(z_{r+1}, \dots, z_n)$ , соответственно. В силу основных предположений эти совокупности случайных величин независимы, что доказывает первую часть теоремы.

2) Оценки  $\vec{\psi}^*$  несмещены и являются линейными комбинациями  $\vec{y}$ , следовательно, вектор  $\vec{\psi}^*$  также распределен нормально с вектором средних  $\mathbf{E} \vec{\psi}^* = \vec{\psi}$  (в силу несмещенности). Вычислим его ковариационную матрицу

$$\text{Cov}(\vec{\psi}^*) = \text{Cov}(A \vec{y}) = A \text{Cov}(\vec{y}) A' = \sigma^2 A A'.$$

Для завершения доказательства второго пункта теоремы осталось показать, что матрица  $B = A A'$  имеет ранг  $q$ .

Из следующей цепочки равенств

$$C \vec{\beta} = \vec{\psi} = \mathbf{E} \vec{\psi}^* = \mathbf{E} A \vec{y} = A X' \vec{\beta},$$

верной для всех  $\vec{\beta}$ , получаем, что  $C = A X'$ . Следовательно, ранг  $\text{rang}(C) \leq \text{rang}(A) \leq q$ . Однако, так как все  $q$  параметрических функций  $\vec{\psi}$  линейно независимы, то  $\text{rang}(C) = q$ . Поэтому  $\text{rang}(A) = q$  и  $\text{rang}(A A') = q$ .

3) Как уже отмечалось,  $SS_e$  принадлежит пространству ошибок и

$$SS_e = \sum_{i=r+1}^n z_i^2,$$

причем все  $z_i \sim \mathcal{N}_1(0, \sigma^2)$  и независимы. Таким образом,  $SS_e / \sigma^2$  есть сумма квадратов  $n - r$  независимых  $\mathcal{N}_1(0, 1)$  случайных величин. По определению эта сумма имеет хи-квадрат распределение с  $n - r$  степенями свободы.  $\otimes$

## § 2. Доверительные множества для функций, допускающих оценку. Критерий проверки гипотез

Утверждения теоремы II.1 позволяют строить доверительные множества для параметрических функций, допускающих оценку.

Пусть  $\psi_1, \dots, \psi_q$  —  $q$  линейно независимых параметрических функций, допускающих оценку и  $\vec{\psi}^*$  — вектор их МНК-оценок. По теореме II.1  $\vec{\psi}^* \sim \mathcal{N}_q(\vec{\psi}, \sigma^2 B)$  с некоторой матрицей  $B$  полного ранга  $q$ . Тогда, как известно, квадратичная форма

$$\frac{(\vec{\psi}^* - \vec{\psi})' B^{-1} (\vec{\psi}^* - \vec{\psi})}{\sigma^2} \sim \chi_q^2, \quad (\text{II.3})$$

то есть распределена как хи-квадрат случайная величина с  $q$  степенями свободы. Кроме того, в силу утверждений 1) и 3) теоремы II.1,  $SS_e / \sigma^2 \sim \chi_{n-r}^2$  и не зависит от квадратической формы (II.3). Следовательно, отношение

$$\frac{(\vec{\psi}^* - \vec{\psi})' B^{-1} (\vec{\psi}^* - \vec{\psi}) / q}{SS_e / (n - r)} \sim F_{q, n-r}, \quad (\text{II.4})$$

то есть имеет распределение Фишера с  $(q, n - r)$  степенями свободы. Поэтому, если  $F_{q, n-r}^\alpha$  — верхняя  $\alpha$ -квантиль соответствующего распределения Фишера, тогда

$$\mathbf{P} \left\{ \frac{(\vec{\psi}^* - \vec{\psi})' B^{-1} (\vec{\psi}^* - \vec{\psi}) / q}{\mathfrak{S}^2} \leq F_{q, n-r}^\alpha \right\} = 1 - \alpha \quad (\text{II.5})$$

с  $\mathfrak{S}^2 = SS_e / (n - r)$ . Таким образом, имеет место

**Теорема II.2.** В основных предположениях (II.1) множество тех  $\vec{\psi}$ , для которых выполняется неравенство

$$(\vec{\psi}^* - \vec{\psi})' B^{-1} (\vec{\psi}^* - \vec{\psi}) \leq q \mathfrak{S}^2 F_{q, n-r}^\alpha, \quad (\text{II.6})$$

образует  $(1 - \alpha)$ -доверительное множество (доверительный эллипсоид) для вектора линейно независимых параметрических функций  $\vec{\psi} = (\psi_1, \dots, \psi_q)$ .

Случай зависимых параметрических функций сводится к только что рассмотренному путем сокращения количества оцениваемых



функций до максимального линейно независимого числа. Остальные функции будут линейно выражаться через выбранные.

В частном случае  $q = 1$  распределение Фишера  $\mathcal{F}_{1,n-r}$  совпадает с распределением квадрата «студентовской» случайной величины с  $n - r$  степенями свободы и доверительный эллипсоид представляет собой стандартный интервал числовой оси

$$\psi^* - t_{n-r}^{\alpha/2} \hat{\sigma} \leq \psi \leq \psi^* + t_{n-r}^{\alpha/2} \hat{\sigma},$$

где  $\hat{\sigma} = \mathfrak{S}|a|$  – оценка стандартного отклонения  $\psi^*$ ,  $t_{n-r}^{\alpha/2}$  – верхняя  $\alpha/2$ -квантиль распределения Стьюдента с  $n - r$  степенями свободы.

### Критерий проверки гипотезы, построенный по доверительному эллипсоиду

Доверительный эллипсоид, предложенный в теореме II.2, можно использовать для построения критерия проверки любых простых гипотез относительно набора параметрических функций. В частности, имеет место

**Теорема II.3.** *При проверке гипотезы*

$$\mathbf{H} : \psi_1 = \psi_2 = \dots = \psi_q = 0$$

*относительно  $q$  линейно независимых параметрических функций критерий, отвергающий гипотезу  $\mathbf{H}$ , если  $\vec{\psi}^*{}' B^{-1} \vec{\psi}^* > q \mathfrak{S}^2 F_{q,n-r}^\alpha$ , имеет уровень  $\alpha$ .*

*Доказательство.* Указанным критерием гипотеза  $\mathbf{H}$  принимается тогда и только тогда, когда вектор  $\vec{\psi} = 0$  принадлежит доверительному множеству (II.6). Следовательно, вероятность принятия гипотезы, если она верна, равна вероятности накрытия доверительным эллипсоидом (II.6) истинной точки  $\vec{\psi} = 0$ . По построению, эта вероятность равна  $1 - \alpha$ .  $\otimes$

## § 3. Критерий отношения правдоподобий.

### Критерий Фишера

Как известно, наиболее мощный критерий проверки гипотез основан обычно на статистике отношения правдоподобий (ОП)

$$\lambda = \frac{\max_{\omega} f(\vec{y})}{\max_{\Omega} f(\vec{y})},$$

где  $f$  – плотность распределения наблюдений  $\vec{y}$  и максимумы берутся по параметрам, определяющим плотность и принадлежащим либо основным предположениям  $\Omega$ , либо множеству  $\omega$ , характеризующему предположения гипотезы.

В рамках рассмотренной нами геометрической интерпретации основные предположения  $\Omega$  задают  $r$ -мерное подпространство  $V_r$  (– пространство всех линейных комбинаций  $\mathbf{X}'\vec{\beta}$ ) пространства  $\mathcal{R}^n$ . Предположения  $\omega$  определяют  $(r - q)$ -мерное подпространство  $V_{r-q}$  пространства  $V_r$ . (Для примера, условия  $y = 0, z = 0$  задают одномерное подпространство – ось ОХ трехмерного пространства  $\mathcal{R}^3$ ).

Для вычисления статистики ОП выпишем плотность совместного распределения наблюдений  $f(\vec{y}) = f_1(y_1) \cdots f_n(y_n)$ . В наших предположениях каждое наблюдение получено из нормального распределения, то есть

$$f_i(y_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} \left( y_i - \sum_{j=1}^p x_{ji}\beta_j \right)^2 \right\}$$

и

$$f(\vec{y}) = \frac{1}{(2\pi\sigma^2)^{-n/2}} \exp \left\{ -\frac{\mathcal{G}(\vec{y}, \vec{\beta})}{2\sigma^2} \right\}, \quad (\text{II.7})$$

где  $\mathcal{G}(\vec{y}, \vec{\beta}) = \|\vec{y} - \mathbf{X}'\vec{\beta}\|^2$  равно как раз той сумме квадратов, которая минимизировалась при построении ОМНК.

При фиксированном  $\sigma^2$  максимум плотности (II.7) достигается в точках достижения минимума  $\mathcal{G}$ . В случае  $\Omega$  это есть точки  $\vec{\beta}^*$  (обозначим их здесь  $\vec{\beta}_\Omega$ ), для которых  $\mathbf{X}'\vec{\beta}_\Omega = \text{Proj}_{V_r} \vec{y}$ . Аналогично, точки достижения максимума плотности по  $\omega$  суть множество точек  $\vec{\beta}_\omega$ , для которых  $\mathbf{X}'\vec{\beta}_\omega = \text{Proj}_{V_{r-q}} \vec{y}$ . Соответствующие минимумы функции  $\mathcal{G}$  обозначим  $\mathcal{G}_\Omega$  и  $\mathcal{G}_\omega$ . Таким образом, для нахождения максимума плотности (по  $\Omega$  и  $\omega$ ) необходимо найти теперь максимум по  $\sigma^2$  выражений вида

$$\frac{1}{(2\pi\sigma^2)^{-n/2}} \exp \left\{ -\frac{\mathcal{G}^*}{2\sigma^2} \right\}, \quad (\text{II.8})$$

где  $\mathcal{G}^*$  равно  $\mathcal{G}_\Omega$  или  $\mathcal{G}_\omega$ . Стандартным способом легко устанавливается, что максимум (II.8) достигается при

$$\sigma^2 = \frac{\mathcal{G}^*}{n}.$$

Подставляя полученные значения в (II.7), получаем статистику ОП

$$\lambda = \left( \frac{\mathcal{G}_\omega}{\mathcal{G}_\Omega} \right)^{-n/2}. \quad (\text{II.9})$$

Для построения критерия нам необходимо теперь найти распределение этой статистики. Оказывается, что следующее взаимно однозначное преобразование статистики  $\lambda$  будет иметь хорошо известное распределение Фишера и, следовательно, более удобно с практической точки зрения.

**Теорема II.4.** *Если верна гипотеза  $\mathbf{H}$ , то статистика*

$$\mathcal{F} = \frac{n-r}{q} \cdot \frac{\mathcal{G}_\omega - \mathcal{G}_\Omega}{\mathcal{G}_\Omega} \sim F_{q, n-r}.$$

*Доказательство.* Для доказательства теоремы несколько переопределим базис, построенный при определении канонических переменных. Рассмотрим три пространства  $V_{r-q} \subset V_r \subset \mathcal{R}^n$  и выберем ортонормированный базис в пространстве  $\mathcal{R}^n$  по следующей схеме

$$\underbrace{\underbrace{\underbrace{\vec{\alpha}_1, \dots, \vec{\alpha}_q, \vec{\alpha}_{q+1}, \dots, \vec{\alpha}_r, \vec{\alpha}_{r+1}, \dots, \vec{\alpha}_n}_{\text{базис } V_{r-q}}}_{\text{базис } V_r}}_{\text{базис } \mathcal{R}^n}.$$

Как было показано ранее, средние значения канонических переменных  $\mathbf{E} z_i = 0$  для  $i = r+1, \dots, n$ . Аналогично показывается, что, если верна гипотеза  $\mathbf{H}$  (то есть регрессия  $\mathbf{X}' \vec{\beta} \in V_{r-q}$ ), то средние значения  $\mathbf{E} z_i = 0$  и для  $i = 1, \dots, r$ .

Далее, так как соответствующие ОМНК найдены через проекции на пространства  $V_r$  и  $V_{r-q}$ , то

$$\mathcal{G}_\Omega = \|\vec{y} - \text{Proj}_{V_r} \vec{y}\|^2 = \left\| \sum_{i=1}^n z_i \vec{\alpha}_i - \sum_{i=1}^r z_i \vec{\alpha}_i \right\|^2 = \sum_{i=r+1}^n z_i^2.$$

и, аналогично,

$$\mathcal{G}_\omega = \sum_{i=1}^q z_i^2 + \sum_{i=r+1}^n z_i^2.$$

Таким образом, статистика отношения правдоподобий записывается через канонические переменные в виде

$$\mathcal{F} = \frac{n-r}{q} \frac{\sum_{i=1}^q z_i^2}{\sum_{i=r+1}^n z_i^2} = \frac{n-r}{q} \frac{\sum_{i=1}^q \left(\frac{z_i}{\sigma}\right)^2}{\sum_{i=r+1}^n \left(\frac{z_i}{\sigma}\right)^2}.$$

Если верна гипотеза  $\mathbf{H}$ , то, как отмечалось выше, переменные  $z_i$  в числителе имеют среднее 0. Следовательно, числитель статистики  $\mathcal{F}$  есть сумма квадратов независимых  $\mathcal{N}_1(0, 1)$  случайных величин и поэтому имеет хи-квадрат распределение с  $q$  степенями свободы. Знаменатель  $\mathcal{F}$  (в любом случае) имеет хи-квадрат распределение с  $n - r$  степенями свободы, что в силу условия независимости канонических переменных доказывает утверждение теоремы.  $\otimes$

**Замечание i.** В знаменателе статистики  $\mathcal{F}$  стоит сумма квадратов ошибок  $SS_e$ . Среднее значение этой суммы  $SS_e / (n - r)$  (равное  $\mathfrak{S}^2$  – несмещенной оценке  $\sigma^2$ ) принято обозначать  $\overline{SS_e}$ . Сумму квадратов в числителе  $\mathcal{F}$  обозначают  $SS_H$ , а её среднее  $\overline{SS_H} = SS_H / q$ . В этих обозначениях статистика

$$\mathcal{F} = \frac{\overline{SS_H}}{\overline{SS_e}}.$$

### **F-критерий проверки линейных гипотез**

С учетом полученного распределения статистики отношения правдоподобий можно предложить следующий критерий проверки гипотез.

**Определение 1.** Критерий проверки гипотезы  $\mathbf{H} : \psi_1 = \psi_2 = \dots = \psi_q = 0$ , относительно  $q$  линейно независимых параметрических функций, отвергающий  $\mathbf{H}$  при значениях статистики

$$\mathcal{F} > F_{q, n-r}^\alpha,$$

где  $F_{q, n-r}^\alpha$  – верхняя  $\alpha$ -квантиль распределения Фишера с  $(q, n - r)$  степенями свободы, называется *F-критерием*.

**Замечание ii.** Алгебраическими методами легко показывается, что F-критерий эквивалентен построенному выше критерию, основанному на доверительном эллипсоиде для вектора параметрических функций. Не вдаваясь в подробности, заметим только, что гипотезу  $\mathbf{H}$  можно записать в эквивалентной форме  $\mathbf{H} : \zeta_1 = \dots = \zeta_q = 0$  относительно средних значений канонических переменных. Несмещенными оценками с минимальной дисперсией этих неизвестных параметров являются, очевидно, сами канонические переменные –  $\zeta_1^* = z_1, \dots, \zeta_q^* = z_q$ . Ковариационная матрица вектора  $(z_1, \dots, z_q)'$  пропорциональна единичной матрице, поэтому числитель статистики П.6 полностью совпадает с числителем статистики F-критерия. Знаменатели этих статистик равны по определению.

**Замечание iii.** Мощность F-критерия равна

$$\mathbf{P}\{\mathcal{F} > F_{q, n-r}^\alpha\},$$

где вероятность вычисляется при общих предположениях  $\Omega$ . В этих предположениях средние значения  $\zeta_i = \mathbf{E} z_i$  ( $i = 1, \dots, q$ ) уже не равны 0 и числитель статистики отношения правдоподобий будет иметь нецентральное хи-квадрат распределение, а сама статистика нецентральное распределение Фишера. Параметр нецентральности равен

$$\delta^2 = \frac{\sum_{i=1}^q \zeta_i^2}{\sigma^2}.$$

Таким образом, если неизвестные значения  $\zeta_1, \dots, \zeta_q$  и  $\sigma^2$ , входящие в  $\delta^2$ , заменить на их соответствующие оценки, то мощность F-критерия может быть приближенно вычислена с помощью таблиц нецентрального распределения Фишера.

Последнее, что необходимо теперь сделать, – сравнить мощность построенного F-критерия с мощностью других возможных критериев уровня  $\alpha$ . Утверждать, что F-критерий имеет наибольшую мощность среди всех критериев уровня  $\alpha$  нельзя, однако ему присущ ряд оптимальных свойств, которые приведем здесь без доказательства.

Во-первых, для произвольного критерия  $W$  проверки гипотезы  $\mathbf{H}$  определим мощность  $\mathfrak{b}(\theta, W)$  как вероятность отвергнуть гипотезу  $\mathbf{H}$ , когда истинное значение вектора неизвестных параметров есть  $\theta$ . Критерий уровня  $\alpha$  имеет мощность  $\mathfrak{b}(\theta, W) \leq \alpha$  для всех  $\theta$ , удовлетворяющих условиям гипотезы  $\mathbf{H}$  (коротко,  $\theta \in \mathbf{H}$ ).

**Теорема II.5.** *Среди всех критериев уровня  $\alpha$ , мощность которых зависит от неизвестных параметров только через значения параметра нецентральности  $\delta^2$ , F-критерий имеет равномерно наибольшую мощность при  $\theta \notin \mathbf{H}$ .*

Далее, пусть  $\mathcal{W}_\alpha$  – класс критериев уровня  $\alpha$ . Рассмотрим огибающую функцию мощности

$$\tilde{\mathfrak{b}}(\theta) = \sup_{W \in \mathcal{W}_\alpha} \mathfrak{b}(\theta, W)$$

как (недостижимый) эталон мощности критерия. Назовем силой критерия  $W$  максимум расхождения между огибающей функцией мощности и мощностью  $W$ :

$$\sup_{\theta \notin \mathbf{H}} (\tilde{\mathfrak{b}}(\theta) - \mathfrak{b}(\theta, W)). \quad (\text{II.10})$$

Критерий, минимизирующий (II.10), называется *наиболее сильным*.

**Теорема II.6.** *F-критерий – наиболее сильный критерий уровня  $\alpha$ .*

## Г л а в а III

### МОДЕЛИ ДИСПЕРСИОННОГО АНАЛИЗА

#### § 1. Однофакторный дисперсионный анализ

Простейшим случаем ДА является *однофакторный анализ*. Этот термин относится к сравнению средних нескольких одномерных популяций. Обозначим эти средние через  $\beta_1, \dots, \beta_p$  и предположим, что наблюдения в каждой популяции имеют нормальное распределение с одинаковой дисперсией  $\sigma^2$ . Пусть в каждой из популяций было произведено  $n_1, \dots, n_p$  наблюдений, соответственно. Тогда  $i$ -ое наблюдение в  $j$ -ой популяции может быть записано в виде

$$y_{ji} = \beta_j + \varepsilon_{ji}, \quad (\text{III.1})$$

где  $\varepsilon_{ji}$  – ошибка наблюдения. К данным подобного рода может быть применена общая теория, разработанная в предыдущих главах.

Представим соотношение III.1 в матричной форме. Для этого рассмотрим  $p$ -мерный вектор  $\vec{u}_j = (0, \dots, 0, 1, 0, \dots, 0)'$  с единственным ненулевым элементом, стоящим на  $j$ -ом месте, и матрицу

$$\mathbf{X} = (\vec{u}_1, \dots, \vec{u}_1, \dots, \vec{u}_p, \dots, \vec{u}_p),$$

у которой первые  $n_1$  столбцов совпадают с вектором  $\vec{u}_1$ , последующие  $n_2$  столбцов совпадают с  $\vec{u}_2$  и т.д., последние  $n_p$  столбцов совпадают с вектором  $\vec{u}_p$ . Тогда вектор наблюдений

$$\vec{y} = \begin{pmatrix} y_{11} \\ \vdots \\ y_{1n_1} \\ \vdots \\ y_{pn_p} \end{pmatrix} = \mathbf{X}' \vec{\beta} + \vec{\varepsilon} = \begin{pmatrix} \vec{u}_1' \\ \vdots \\ \vec{u}_1' \\ \vdots \\ \vec{u}_p' \end{pmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} + \vec{\varepsilon}$$

с соответствующим вектором ошибок  $\vec{\varepsilon}$ .

Ранее было установлено, что ОМНК  $\vec{\beta}^*$  параметров  $\vec{\beta}$  удовлетворяют нормальным уравнениям  $\mathbf{X} \mathbf{X}' \vec{\beta}^* = \mathbf{X} \vec{y}$ , причем, если  $r = \text{rang } \mathbf{X} = p$ , то  $\vec{\beta}^* = \mathcal{S}^{-1} \mathbf{X} \vec{y}$  с информационной матрицей  $\mathcal{S} = \mathbf{X} \mathbf{X}'$ . В рассматриваемом случае информационная матрица

$$\mathcal{S} = \mathbf{X} \mathbf{X}' = (\vec{u}_1, \dots, \vec{u}_1, \dots, \vec{u}_p) \begin{pmatrix} \vec{u}_1' \\ \vdots \\ \vec{u}_1' \\ \vdots \\ \vec{u}_p' \end{pmatrix} = \sum_{i=1}^{n_1} \vec{u}_1 \vec{u}_1' + \dots + \sum_{i=1}^{n_p} \vec{u}_p \vec{u}_p'.$$

Произведение  $\vec{u}_j \vec{u}_j'$  образует матрицу  $p \cdot p$ , у которой все элементы равны 0, кроме  $j$ -ого диагонального элемента, который равен 1. Таким образом, информационная матрица равна диагональной матрице

$$\mathcal{S} = \begin{pmatrix} n_1 & 0 & \dots & 0 \\ 0 & n_2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & n_p \end{pmatrix} \quad \text{и} \quad \mathcal{S}^{-1} = \begin{pmatrix} \frac{1}{n_1} & 0 & \dots & 0 \\ 0 & \frac{1}{n_2} & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & \frac{1}{n_p} \end{pmatrix}.$$

Легко видеть, что

$$\mathbf{X} \vec{y} = \begin{pmatrix} \sum_{i=1}^{n_1} y_{1i}, \dots, \sum_{i=1}^{n_p} y_{pi} \end{pmatrix}.$$

Следовательно,

$$\vec{\beta}^* = (y_{1.}, \dots, y_{p.})', \quad (\text{III.2})$$

где с целью сокращения записи введено обозначение  $y_{j.} = \frac{1}{n_j} \sum_{i=1}^{n_j} y_{ji}$ .

Сумма квадратов ошибок равна

$$\mathcal{G}_\Omega = S S_e = \sum_{j=1}^p \sum_{i=1}^{n_j} (y_{ji} - y_{j.})^2 = \sum_{ji} y_{ji}^2 - \sum_{j=1}^p n_j y_{j.}^2.$$

Последнее равенство является следствием известного утверждения (сравните с двумя способами вычисления дисперсии), которое мы будем часто использовать в дальнейшем изложении.

**Лемма III.1.** Для любого набора чисел  $z_1, \dots, z_m$  сумма

$$\sum_{i=1}^m (z_i - z_{.})^2 = \sum_{i=1}^m z_i^2 - m z_{.}^2.$$

Следуя общей теории, можно найти распределение вектора ОМНК. Заметим сначала, что так как существует обратная матрица  $\mathcal{S}^{-1}$ , то ранг нашей задачи  $r = p$ . Следовательно, в силу следствия 1 ковариационная матрица  $\text{Cov} \vec{\beta}^* = \sigma^2 \mathcal{S}^{-1} = \sigma^2 \mathbf{I}$ . Далее, параметрические функции  $\psi_j(\vec{\beta}) = \beta_j$  ( $j = 1, \dots, p$ ) допускают оценку ( $= y_{j.}$ ), отсюда, как следствие теоремы (II.1), получаем следующее утверждение, которое здесь легче доказать непосредственно.

**Теорема III.1.** *В основных предположениях с нормальным распределением ошибок ОМНК  $\beta_j^* = y_{j.} \sim \mathcal{N}_1(\beta_j, \sigma^2/n_j)$  ( $j = 1, \dots, p$ ) и независимы.*

Перейдем теперь к вопросу проверки гипотез. Обычно в данной ситуации проверяется стандартная гипотеза о равенстве средних значений всех популяций

$$\mathbf{H} : \beta_1 = \beta_2 = \dots = \beta_p.$$

В терминах равенства нулю параметрических функций эта гипотеза записывается в виде

$$\mathbf{H} : \psi_1(\vec{\beta}) = \psi_2(\vec{\beta}) = \dots = \psi_{p-1}(\vec{\beta}) = 0,$$

где функции  $\psi_j(\vec{\beta}) = \beta_j - \beta_p$ .

Для построения F-критерия необходимо найти теперь ОМНК в предположениях  $\mathbf{H}$ . Вычислим их непосредственно.

Если справедлива гипотеза  $\mathbf{H}$ , тогда  $\mathbf{X}' \vec{\beta} = (\beta, \dots, \beta)'$  и, следовательно, для нахождения ОМНК необходимо найти минимум функции

$$\| \vec{y} - \mathbf{X}' \vec{b} \|^2 = \sum_{j=1}^p \sum_{i=1}^{n_j} (y_{ji} - b)^2$$

по параметру  $b$ . Другими словами, необходимо найти константу  $b$ , дающую наилучший среднеквадратический прогноз всего вектора данных  $\vec{y}$ . Хорошо известно, что оптимальная константа равна среднему значению всех данных. Таким образом, ОМНК неизвестного параметра  $\beta$  в предположениях гипотезы равна

$$\beta_\omega = y_{..} \stackrel{\text{def}}{=} \frac{1}{n} \sum_{j=1}^p \sum_{i=1}^{n_j} y_{ji} = \frac{1}{n} \sum_{j=1}^p n_j y_{j.} .$$

Сумма квадратов ошибок (см. Лемму III.1)

$$\mathcal{G}_\omega = \sum_{j=1}^p \sum_{i=1}^{n_j} (y_{ji} - y_{..})^2 = \sum_{ji} y_{ji}^2 - n y_{..}^2$$



и

$$SS_H = \mathcal{G}_\omega - \mathcal{G}_\Omega = \sum_{j=1}^p n_j y_j^2 - n y_{..}^2 = \sum_{j=1}^p n_j (y_{j.} - y_{..})^2 .$$

В рассматриваемом случае гипотеза задается  $q = p - 1$  линейно независимыми функциями. Следовательно, статистика F-критерия

$$\mathcal{F} = \frac{n - p}{p - 1} \cdot \frac{\sum_{j=1}^p n_j (y_{j.} - y_{..})^2}{\sum_{ji} y_{ji}^2 - \sum_{j=1}^p n_j y_j^2} \sim F_{p-1, n-p} .$$

Сведем полученные результаты в единую таблицу.

Таблица однофакторного ДА

Источник разброса	$SS$	Ст.свободы	$\overline{SS}$	$\mathcal{F}$	$\alpha$
Между группами	$SS_H$	$\nu_H = p - 1$	$\overline{SS}_H$	$\frac{\overline{SS}_H}{\overline{SS}_e}$	$1 - F_{\nu_H, \nu_e}(\mathcal{F})$
Внутри групп	$SS_e$	$\nu_e = n - p$	$\overline{SS}_e$	–	–
Полная сумма квадратов	$SS_\Pi$	$n - 1$	–	–	–

## Разбиение полной суммы квадратов

Легко видеть, что справедливо соотношение

$$\sum_{ji} y_{ji}^2 - n y_{..}^2 = SS_e + SS_H ,$$

которое называется разбиением полной суммы квадратов. Левая часть этого равенства обозначается обычно  $SS_\Pi$  . С помощью этого равенства легче вычислить  $SS_e = SS_\Pi - SS_H$  .

## § 2. Сравнения

Если гипотеза о равенстве всех средних отвергается, то возникает задача получения дополнительной информации о различающихся популяциях. Для этой цели рекомендуется использовать метод сравнений. Введем

**Определение 1.** Сравнением параметров  $\beta_1, \dots, \beta_p$  называется параметрическая функция  $\psi = \sum c_j \beta_j$  с коэффициентами, удовлетворяющими условию  $\sum c_j = 0$ .

**Пример.** Если требуется сравнить две популяции (например, первую и вторую), естественно рассмотреть параметрическую функ-

цию  $\psi = \beta_1 - \beta_2$ , которая будет, очевидно, сравнением. С другой стороны, если первые  $k$  и, соответственно, последние  $p - k$  популяций можно объединить по какому-либо признаку в две группы, то для оценки различий в этих группах можно использовать сравнение  $\psi = \sum_{j=1}^k \beta_j/k - \sum_{j=k+1}^p \beta_j/(p - k)$ .

В однофакторном ДА любое сравнение допускает оценку. А именно, для сравнения  $\psi = \sum c_j \beta_j$  МНК-оценкой (по теореме Гаусса-Маркова) будет статистика  $\psi^* = \sum c_j y_j$ . Дисперсия этой оценки

$$\mathbf{D}(\psi^*) = \sum_{j=1}^p c_j^2 \mathbf{D}(y_j) = \sigma^2 \sum_{j=1}^p \frac{c_j^2}{n_j},$$

а её оценка

$$\mathfrak{D}(\psi^*) = \mathfrak{S}^2 \sum_{j=1}^p \frac{c_j^2}{n_j}. \quad (\text{III.3})$$

В обычной практике применения статистических методов сразу после отвержения основной гипотезы рассматривается некоторая вспомогательная гипотеза о равенстве нулю какого-либо сравнения, к проверке которой применяются аналогичные методы (F-критерий). Однако такой составной вывод будет уже иметь уровень значимости выше ожидаемого  $\alpha$ . Г.Шеффе предложил метод (называемый в честь автора S-методом), позволяющий оценить сразу все сравнения, связанные с основной гипотезой.

Рассмотрим набор из  $q$  линейно независимых сравнений  $\psi_1, \dots, \psi_q$ , (описывающих гипотезу H) и пространство  $\mathfrak{L}$ , порожденное этими функциями. Справедлива

**Теорема III.2.** Положим  $\Delta = qF_{q,n-r}^\alpha$ . Вероятность того, что

$$\psi^* - \sqrt{\Delta \mathfrak{D}(\psi^*)} \leq \psi \leq \psi^* + \sqrt{\Delta \mathfrak{D}(\psi^*)} \quad (\text{III.4})$$

одновременно для всех сравнений  $\psi \in \mathfrak{L}$ , равна  $1 - \alpha$ .

Связь этой теоремы с F-критерием объясняет следующая

**Теорема III.3.** F-критерий отвергает гипотезу H тогда и только тогда, когда для какого-либо сравнения  $\psi$  интервал III.4 не покрывает точку  $\psi = 0$ .

Таким образом, если гипотеза H отвергается F-критерием, то, применяя доверительный интервал III.4, можно попытаться отыскать сравнение, способствовавшее этому выводу.

Рассмотрим следующий

**Пример.** Сравниваются веса новорожденных поросят в восьми опоросах. Были получены следующие данные.

1	2	3	4	5	6	7	8
2.0	3.5	3.3	3.2	2.6	3.1	2.6	2.5
2.8	2.8	3.6	3.3	2.6	2.9	2.2	2.4
3.3	3.2	2.6	3.2	2.9	3.1	2.2	3.0
3.2	3.5	3.1	2.9	2.0	2.5	2.5	1.5
4.4	2.3	3.2	3.3	2.0		1.2	
3.6	2.4	3.3	2.5	2.1		1.2	
2.9	2.0	2.9	2.6				
2.5	1.6	3.4	2.8				
2.8		3.2					
2.1		3.2					

Требуется с 10%-ым уровнем значимости ( $\alpha = 0.1$ )

а) проверить гипотезу отсутствия различий между средними весами в восьми приплодах;

б) в предположении, что опоросы 1,3 и 4 получены от одной свиноматки, а остальные пять от другой, установить значимость различия между средними весами в этих двух группах;

в) проверить значимость различия между средними весами в больших опоросах (1,2,3,4) и меньших (5,6,7,8).

С целью упрощения вычислений преобразуем данные, умножая их на 10 (убираем десятичную точку) и вычитая 25 (сдвигаем к приближительному общему центру). В следующей таблице наряду с преобразованными данными указаны и их квадраты:

1	2	3	4	5	6	7	8
Y: Y <sup>2</sup>	Y: Y <sup>2</sup>	Y: Y <sup>2</sup>	Y: Y <sup>2</sup>	Y: Y <sup>2</sup>	Y: Y <sup>2</sup>	Y: Y <sup>2</sup>	Y: Y <sup>2</sup>
-5: 25	10: 100	8: 64	7: 49	1: 1	6: 36	1: 1	0: 0
3: 9	3: 9	11: 121	8: 64	1: 1	4: 16	-3: 9	-1: 1
8: 64	7: 49	1: 1	7: 49	4: 16	6: 36	-3: 9	5: 25
7: 49	10: 100	6: 36	4: 16	-5: 25	0: 0	0: 0	-10: 100
19: 361	-2: 4	7: 49	8: 64	-5: 25	:	-13: 169	:
11: 121	-1: 1	8: 64	0: 0	-4: 16	:	-13: 169	:
4: 16	-5: 25	4: 16	1: 1	:	:	:	:
0: 0	-9: 81	9: 81	3: 9	:	:	:	:
3: 9	:	7: 49	:	:	:	:	:
-4: 16	:	7: 49	:	:	:	:	:
10	8	10	8	6	4	6	4
46:	13:	68:	38:	-8:	16:	-31:	-6:
: 670	: 369	: 530	: 252	: 84	: 88	: 357	: 126
4.6:	1.63:	6.8:	4.75:	-1.3:	4:	-5.2:	-1.5:
: 21.16	: 2.64	: 46.24	: 22.56	: 1.77	: 16	: 26.73	: 2.25
: 211.6	: 21.13	: 462.4	: 180.5	: 10.7	: 64	: 160.2	: 9

Кроме того, в этой таблице после выборочных данных приведены:

в 1-ой строке – количества данных в популяциях –  $n_j$ ;

во 2-ой строке – суммы первых степеней –  $\sum_{i=1}^{n_j} y_{ji}$ ;

в 3-ей строке – суммы квадратов в популяциях –  $\sum_{i=1}^{n_j} y_{ji}^2$ ;

в 4-ой строке – выборочные средние в популяциях –  $y_{j.}$ ;

в 5-ой строке – их квадраты –  $y_{j.}^2$ ;

в 6-ой строке – произведения  $n_j y_{j.}^2$ .

Суммируя значения в последних шести строках, получаем

$$\begin{aligned}
 n &= 56 && \text{– сумма в 1-ой строке;} \\
 y_{..} &= \frac{136}{56} = 2.43 && \text{– во 2-ой строке;} \\
 y_{..}^2 &= 5.9 \quad ny_{..}^2 = 330.3 && \\
 SS_{\Pi} &= 2476 - 330.3 = 2145.7 && \text{– в 3-ей строке;} \\
 SS_e &= 2476 - 1119.5 = 1356.5 && \text{– в 6-ой строке;} \\
 \sigma^2 &= \overline{SS_e} = 1356.3 / (56 - 8) = 28.26 && \text{– оценка дисперсии;} \\
 SS_H &= 2145.7 - 1356.5 = 789.2 . && 
 \end{aligned}$$

Таким образом, найдя по таблицам распределения Фишера критический уровень значимости  $\alpha_{кр} = 1 - F_{7,48}(3.99) = 0.0016$ , получаем следующую таблицу:

Таблица однофакторного ДА

Источник разброса	$SS$	Ст. свободы	$\overline{SS}$	$\mathcal{F}$	$\alpha$
Между группами	789.2	7	112.7	3.99	0.0016
Внутри групп	1356.5	48	28.26		
Полная сумма квадратов	2145.7	55			

Гипотеза о равенстве средних значений во всех опросах отвергается с уровнем значимости  $0.001 < 0.1$ . Перейдем теперь к решению следующих двух поставленных задач.

Для решения второй задачи (сравнение среднего веса в 1,3 и 4 опросах с остальными) рассмотрим сравнение

$$\psi = \frac{1}{3}(\beta_1 + \beta_3 + \beta_4) - \frac{1}{5}(\beta_2 + \beta_5 + \beta_6 + \beta_7 + \beta_8).$$

ОМНК этой параметрической функции равна  $\psi^* = 5.85$

Оценка дисперсии (по формуле III.3)  $\mathfrak{D}(\psi^*) = 28.01 \cdot 0.074 = 2.07$ .

По таблице распределения Фишера со степенями свободы 7 и 48 находим верхнюю 10%-ую точку  $F_{7,48}^{0.1} = 1.85$ . Поэтому  $\Delta = 7 \cdot 1.85 =$

12.95 ,  $\sqrt{\Delta\mathfrak{D}(\psi^*)} = \sqrt{26.6} = 5.16$  и доверительный интервал для сравнения равен (0.69; 11.01). Так как этот интервал не содержит 0, то можно утверждать, что выбор свиноматки влияет на вес приплода.

Решая последнюю поставленную задачу, рассмотрим сравнение

$$\psi = \frac{1}{4}(\beta_1 + \beta_2 + \beta_3 + \beta_4) - \frac{1}{4}(\beta_5 + \beta_6 + \beta_7 + \beta_8).$$

Вычисления, подобные вышеприведенным, дают следующие результаты:

$$\psi^* = 5.45, \quad \mathfrak{D}(\psi^*) = 2.24, \quad \sqrt{\Delta\mathfrak{D}(\psi^*)} = 5.39.$$

Доверительный интервал для  $\psi - (0.06, 9.94)$  также не содержит 0 и, следовательно, можно утверждать, что в более крупных опоросах вес приплода значимо выше.

### § 3. Непараметрические критерии однородности

В том случае, когда нет оснований предполагать нормальность распределения ошибок, применяют так называемые непараметрические (или свободные от распределения) методы анализа однородности различных групп. Эти методы чаще всего основываются на ранговых статистиках, то есть не на абсолютных значениях наблюдаемого отклика, а на их относительных местах при расположении в единый ряд.

Пусть  $r_{ji}$  — ранг (место) наблюдения  $y_{ji}$  в упорядоченной по возрастанию совокупности всех наблюдаемых значений отклика  $\{y_{ji}\}$ . Очевидно, сумма всех рангов от 1 до  $n$  равна  $n(n+1)/2$ , а их общее среднее значение равно  $(n+1)/2$ . Поэтому в качестве меры отклонения  $j$ -ой группы от среднего уровня может быть взята величина  $(r_{j\cdot} - (n+1)/2)^2$ , где, как обычно,  $r_{j\cdot}$  — средний ранг наблюдений, попавших в  $j$ -ую группу. Если гипотеза однородности групп верна, то следует ожидать, что значения этой меры во всех группах будут близки к нулю.

О п р е д е л е н и е 2. Статистика

$$KW = \frac{12}{n(n+1)} \sum_{j=1}^p n_j \left( r_{j\cdot} - \frac{n+1}{2} \right)^2 = \frac{12}{n(n+1)} \sum_{j=1}^p n_j r_{j\cdot}^2 - 3(n+1)$$

называется *статистикой Краскела-Уоллиса*.

**Теорема III.4.** 1) Если наблюдения во всех  $p$  группах имеют одинаковое распределение, то есть группы однородны, тогда

$$KW \rightsquigarrow \chi_{p-1}^2, \quad n \rightarrow \infty.$$

2) Пусть  $K_{p-1}$  – функция распределения хи-квадрат с  $(p-1)$ -ой степенью свободы, тогда критерий, отвергающий гипотезу однородности при критическом уровне значимости

$$\alpha_{кр} = 1 - K_{p-1}(KW) < \alpha,$$

имеет асимптотический уровень  $\alpha$ .

**З а м е ч а н и е i.** При малых объемах выборок распределение статистики  $KW$  может быть найдено точно. Идея его построения основывается на том факте, что в предположениях гипотезы все возможные расположения рангов в группах равновероятны.

**З а м е ч а н и е ii.** На практике чаще всего встречается ситуация, когда наблюдаемые значения отклика содержат совпадающие элементы. В этом случае при малом числе совпадений рекомендуется использовать средние ранги. Так, например, если четыре равных значения отклика занимают места с 7-ого по 10-е, то ранги всех этих наблюдений выбираются равными 8.5. Если совпадений слишком много, то статистику  $KW$  необходимо модернизировать следующим образом

$$KW' = \frac{KW}{1 - \sum_{k=1}^g T_k / (n^3 - n)},$$

где  $g$  – число блоков совпадающих наблюдений,  $T_k = (t_k^3 - t_k)$ ,  $t_k$  – число совпадающих наблюдений в  $k$ -ом блоке.

**З а м е ч а н и е iii.** При  $p = 2$  статистика Краскела-Уоллиса по своему действию эквивалентна статистике Уилкоксона и статистике Манна-Уитни.

Нередко на практике известно, что группы упорядочены по возрастанию влияния фактора. В этом случае можно предложить критерий с более высокой мощностью.

Рассмотрим две произвольные группы  $j$  и  $k$ . Вычислим статистику Манна-Уитни для этих групп

$$U_{jk} = \sum_{il} \phi(y_{ji}, y_{kl}),$$

где

$$\phi(u, v) = \begin{cases} 1, & \text{если } u < v; \\ 1/2, & \text{если } u = v; \\ 0, & \text{если } u > v. \end{cases}$$

**О п р е д е л е н и е 3.** Статистика

$$I = \sum_{1 \leq j < k \leq p} U_{jk}$$

называется *статистикой Джонкхиера*.

**Теорема III.5.** 1) Если наблюдения во всех группах однородны, тогда статистика Джонкхиера асимптотически нормальна со средним

$$m = \frac{1}{4}(n^2 - \sum_{j=1}^k n_j^2)$$

и дисперсией

$$d^2 = \frac{1}{72} \left[ n^2(2n + 3) - \sum_{j=1}^p n_j^2(2n_j + 3) \right].$$

2) Критерий, отвергающий гипотезу однородности в пользу альтернативы возрастающего влияния эффектов при критическом уровне значимости

$$\alpha_{кр} = 1 - \Phi \left( \frac{I - m}{d} \right) < \alpha,$$

имеет асимптотический уровень  $\alpha$ .

## Оценка сравнений

**Определение 4.** Медиана всевозможных разностей пар выборочных данных, попавших в  $j$ -ую и  $k$ -ую группы,

$$z_{jk} = \text{med}(y_{ji} - y_{kl}, \quad i = \overline{1, n_j}, \quad l = \overline{1, n_k}),$$

называется статистикой Ходжеса-Лемана.

**Замечание iv.** Статистика Ходжеса-Лемана представляет собой непараметрическую оценку сдвига между эффектами  $j$ -ой и  $k$ -ой групп  $\tau_{jk} = \beta_j - \beta_k$ . Эта оценка предпочтительна в ситуациях, когда ничего неизвестно о распределении ошибки наблюдений. Кроме того, оценка Ходжеса-Лемана в сравнении со стандартной оценкой, равной разности выборочных средних в группах, менее чувствительна к наличию резко выделяющихся наблюдений.

К сожалению оценка Ходжеса-Лемана не удовлетворяет естественному требованию аддитивности:  $z_{jk} + z_{km} \neq z_{jm}$ , в то время как сдвиги  $\tau_{jk} + \tau_{km} = \tau_{jm}$ . Поэтому рассматривают скорректированную оценку сдвига

$$\Delta_j = \frac{1}{n} \sum_{k=1}^p n_k z_{jk}, \quad j = \overline{1, p},$$

где  $z_{kk} = 0$ . Величина  $\Delta_j$  отражает сдвиг  $j$ -ой группы относительно всех остальных групп.

С помощью этой оценки легко построить оценку любого сравнения. Для этого представим сравнение в эквивалентной форме:

$$\psi = \sum_{j=1}^p c_j \beta_j = \sum_{j=1}^p \sum_{k=1}^p \frac{c_j}{p} (\beta_j - \beta_k),$$

которое легко проверить непосредственно. Тогда оценка сравнения

$$\psi^* = \sum_{j=1}^p \sum_{k=1}^p \frac{c_j}{p} (\Delta_j - \Delta_k).$$

**Пример.** Рассмотрим задачу сравнения весов новорожденных поросят, используя критерий Краскела-Уоллиса. Расположив данные на странице 27 в порядке возрастания, поставим каждому наблюдению в соответствие его ранг с учетом повторов:

	6.5	52.5	48	42	23	37	23	18
	27.5	27.5	54.5	48	23	32	11.5	14.5
	48	42	23	42	32	37	11.5	35
	42	52.5	37	32	6.5	18	18	3
	56	13	42	48	6.5		1.5	
	54.5	14.5	48	18	9.5		1.5	
	32	6.5	32	23				
	18	4	51	27.5				
	27.5		42					
	9.5		42					
$n_j$	10	8	10	8	6	4	6	4
$r_j$	321.5	212.5	419.5	280.5	100.5	124	67	70.5
	32.15	26.56	41.95	35.06	16.75	31	11.17	17.63

Статистика Краскела-Уоллиса равна  $KW = 20.473$ . По таблице распределения хи-квадрат с 7 степенями свободы находим  $\alpha_{кр} = 0.0046$ , то есть гипотезу однородности групп следует отвергнуть.

Модернизированная статистика Краскела-Уоллиса  $KW' = 20.59$  практически не отличается от  $KW$ .

Критерий Джонкхиера к данной задаче, конечно, можно применить, однако в этом нет большого смысла, поскольку нет никаких оснований считать, что группы упорядочены по возрастанию среднего веса.



## § 4. Полный двухфакторный анализ

В предыдущем разделе были представлены методы обработки экспериментов, в которых отклик (результат эксперимента) зависел только от одной управляющей переменной. Предположим теперь, что уже два фактора А и В изменяются в эксперименте, и наша задача состоит в оценке влияния этих факторов на отклик.

Пусть фактор А разбит на  $p_1$  уровней, а фактор В на  $p_2$  уровней. Данные статистического эксперимента, таким образом, распадаются в матрицу с  $p_1$  строками и  $p_2$  столбцами. Наблюдение, попавшее в  $(i, j)$ -ую ячейку этой матрицы, обозначим через  $y_{ij}$ . Разброс данных обусловлен как чисто случайными (неучтенными) факторами, так и влиянием на результат эксперимента рассматриваемых нами факторов. Обозначим через  $\eta_{ij}$  истинное среднее значение отклика в  $(i, j)$ -ой ячейке. Степень различия между этими средними указывает на степень влияния факторов на результат. Основным моментом излагаемой далее теории 2-факторного ДД является представление среднего  $\eta_{ij}$  в виде четырех слагаемых:

$$\left. \begin{aligned} \eta_{ij} &= \mu + \alpha_i + \beta_j + \gamma_{ij}, \\ \sum_{i=1}^{p_1} \alpha_i &= 0, \quad \sum_{j=1}^{p_2} \beta_j = 0, \quad \sum_{i=1}^{p_1} \gamma_{ij} = 0, \quad \sum_{j=1}^{p_2} \gamma_{ij} = 0, \end{aligned} \right\} \quad (\text{III.5})$$

Покажем справедливость этого представления и, заодно, дадим интерпретацию входящих в него слагаемых.

Обозначим через  $A_i = \eta_{i.}$  – истинное среднее  $i$ -ого уровня фактора А (– среднее в  $i$ -ой строке). Аналогично, через  $B_j = \eta_{.j}$  обозначим истинное среднее  $j$ -ого уровня фактора В (– среднее по  $j$ -ому столбцу).

Общее среднее значение ожидаемых результатов экспериментов по всем ячейкам

$$\mu = \eta_{..} = \frac{1}{p_1 p_2} \sum_{ij} \eta_{ij} = \frac{1}{p_1} \sum_{i=1}^{p_1} A_i = \frac{1}{p_2} \sum_{j=1}^{p_2} B_j$$

называется *генеральным средним*.

Отличие среднего  $i$ -ого уровня фактора А от генерального среднего естественно взять в качестве меры влияния фактора А на отклик. Соответствующая разность  $\alpha_i = A_i - \mu$  называется *главным эффектом  $i$ -ого уровня фактора А*. Аналогично, *главным эффектом  $j$ -ого уровня фактора В* называется отклонение  $\beta_j = B_j - \mu$ .

Очевидно, оба эти параметра удовлетворяют соотношениям (III.5).

Заметим далее, что главный эффект  $\alpha_i$  есть среднее (по столбцам) величин  $\eta_{ij} - B_j$ :  $\alpha_i = \sum_j (\eta_{ij} - B_j)/p_2$ . Указанные величины можно трактовать как эффекты уровней  $A$  по отношению к  $j$ -ому уровню  $B$ . Поэтому превышение этой величины над своим средним естественно назвать *взаимодействием  $i$ -ого уровня  $A$  с  $j$ -ым уровнем  $B$* :

$$\gamma_{ij} = \eta_{ij} - B_j - \alpha_i = \eta_{ij} - A_i - B_j + \mu.$$

Легко проверить, что взаимодействие удовлетворяет всем соотношениям из представления (III.5). В частности,

$$\eta_{ij} = A_i + B_j + \gamma_{ij} - \mu = \mu + (A_i - \mu) + (B_j - \mu) + \gamma_{ij} = \mu + \alpha_i + \beta_j + \gamma_{ij}.$$

**Замечание v.** Следует сказать, что пять уравнений (III.5) единственным образом определяют компоненты  $(\mu, \alpha_i, \beta_j, \gamma_{ij})$  по значениям истинных средних  $\eta_{ij}$ . В частности, количество линейно независимых параметров уравнения (III.5) равно  $p = p_1 p_2$  — числу неизвестных средних  $\eta_{ij}$ .

**Замечание vi.** Случай  $\gamma_{ij} = 0$  при всех  $i, j$  называется случаем аддитивности эффектов или отсутствия взаимодействия, так как, по определению, в этом случае истинное среднее

$$\eta_{ij} = \mu + \alpha_i + \beta_j \quad \forall i, j.$$

Количество линейно независимых параметров при отсутствии взаимодействия сокращается до  $p_1 + p_2 - 1$ .

**Замечание vii.** Легко понять, что если общее число наблюдений  $n$  не больше общего числа параметров  $p$ , то статистическими методами невозможно выявить влияние факторов на отклик. Такая ситуация возникает, например, при проведении экспериментов только с одним наблюдением в каждой ячейке. В этом случае общее число наблюдений  $n = p$  и для проведения статистического анализа необходимы дополнительные предположения относительно вида представления (III.5). Обычно делается допущение, что отсутствует взаимодействие и рассматривается дисперсионный анализ с аддитивными эффектами.

### Двухфакторный анализ с одним наблюдением в ячейке.

Итак, рассмотрим 2-факторный ДА с одним наблюдением в каждой ячейке и предположим, что взаимодействия отсутствуют. Тогда основные предположения могут быть записаны в виде

$$\Omega : \begin{cases} y_{ij} = \mu + \alpha_i + \beta_j + e_{ij}, \\ \sum_i \alpha_i = \sum_j \beta_j = 0, \\ \{e_{ij}\} \text{ независимы и распределены } \mathcal{N}(0, \sigma^2). \end{cases}$$

Данные подобного рода можно представить в виде  $\vec{y} = \mathbf{X}'\beta + \vec{\varepsilon}$ , если в качестве матрицы  $\mathbf{X}$  взять следующую матрицу:

	$y_{11}$	$y_{12}$	$\dots$	$y_{1p_2}$	$y_{21}$	$\dots$	$y_{2p_2}$	$\dots$	$y_{p_11}$	$y_{p_12}$	$\dots$	$y_{p_1p_2}$
$\mu$	1	1	$\dots$	1	1	$\dots$	1	$\dots$	1	1	$\dots$	1
$\alpha_1$	1	1	$\dots$	1	0	$\dots$	0	$\dots$	0	0	$\dots$	0
$\vdots$			$\vdots$			$\vdots$		$\vdots$			$\vdots$	
$\alpha_{p_1}$	0	0	$\dots$	0	0	$\dots$	0	$\dots$	1	1	$\dots$	1
$\beta_1$	1	0	$\dots$	0	1	$\dots$	0	$\dots$	1	0	$\dots$	0
$\vdots$			$\vdots$			$\vdots$		$\vdots$			$\vdots$	
$\beta_{p_2}$	0	0	$\dots$	1	0	$\dots$	1	$\dots$	0	0	$\dots$	1

Первая строка и первый столбец здесь приписаны для объяснения способа построения этой матрицы: чтобы получить  $y_{ij}$ , нужно умножить соответствующий столбец на первый столбец.

Таким образом, имеем:

- 1) общее число наблюдений  $n = p_1 p_2$ ;
- 2) общее число параметров  $p = p_1 + p_2 + 1$ ;
- 3) ранг задачи  $r = \text{rang } \mathbf{X} = p_1 + p_2 - 1$  (– первая строка матрицы  $\mathbf{X}$  равна сумме  $p_2$  строк  $\beta$ , а строка  $\alpha_1$  равна сумме всех строк  $\beta$  минус сумма оставшихся строк  $\alpha$ ; остальные строки линейно независимы).

Т.к.  $r < p$ , то для нахождения ОМНК не применим способ, связанный с обращением информационной матрицы. Рассмотрим функцию

$$\mathcal{G} = \sum_{ij} (y_{ij} - \mu - \alpha_i - \beta_j)^2$$

и найдем ее производные по параметрам с учетом связей (III.5):

$$\mathcal{G}'_{\mu} = -2 \sum_{ij} (y_{ij} - \mu - \alpha_i - \beta_j) = -2 \left( \sum_{ij} y_{ij} - p_1 p_2 \mu \right).$$

Следовательно, нормальное уравнение  $\mathcal{G}'_{\mu} = 0$  дает ОМНК

$$\mu^* = \frac{1}{n} \sum_{ij} y_{ij} = y_{..}.$$

Далее,

$$\mathcal{G}'_{\alpha_i} = -2 \sum_{j=1}^{p_2} (y_{ij} - \mu - \alpha_i - \beta_j) = -2 \left( \sum_{j=1}^{p_2} y_{ij} - p_2 \mu - p_2 \alpha_i \right).$$

Отсюда, учитывая, что  $\mu^* = y_{..}$ , получаем ОМНК

$$\alpha_i^* = \frac{1}{p_2} \sum_{j=1}^{p_2} y_{ij} - \mu^* = y_{i.} - y_{..}.$$

Аналогично, ОМНК

$$\beta_j^* = \frac{1}{p_1} \sum_{i=1}^{p_1} y_{ij} - \mu^* = y_{.j} - y_{..} .$$

**З а м е ч а н и е** viii. Несмотря на то, что для данной задачи ранг  $r < p$ , мы получили единственные решения нормальных уравнений. Это обусловлено тем, что здесь на параметры наложены два дополнительных условия и число линейно независимых параметров  $p - 2$  совпадает с рангом задачи.

Таким образом, сумма квадратов ошибок, которая в данном случае, в силу того, что мы не учитываем взаимодействия, иногда называется суммой квадратов взаимодействий, равна

$$SS_e = \sum_{ij} (y_{ij} - y_{i.} - y_{.j} + y_{..})^2 .$$

Её число степеней свободы равно

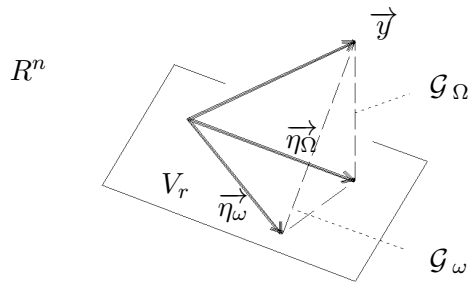
$$\nu_e = n - r = p_1 p_2 - (p_1 + p_2 - 1) = (p_1 - 1)(p_2 - 1) .$$

### Проверка стандартных гипотез

В данной ситуации чаще всего рассматриваются гипотезы об отсутствии воздействия факторов на отклик

$$\mathbf{H}_A : \alpha_i = 0, \forall i, \quad \mathbf{H}_B : \beta_j = 0, \forall j .$$

Для нахождения ОМНК при гипотезе  $\mathbf{H}_A$  необходимо минимизировать функцию  $\mathcal{G} = \sum (y_{ij} - \mu - \beta_j)^2$ . Непосредственными вычислениями легко показывается, что ОМНК  $\mu_\omega^* = \mu^* = y_{..}$  и  $\beta_{j\omega}^* = \beta_j^* = y_{.j}$ . Следовательно, сумма квадратов расхождений  $\mathcal{G}_\omega = \sum_{ij} (y_{ij} - y_{.j} - y_{..})^2$ . Для вычисления суммы квадратов  $SS_H$  воспользуемся следующей геометрической иллюстрацией:



где, как обычно,  $\vec{\eta}_\Omega$  и  $\vec{\eta}_\omega$  – оценки регрессии при основных предположениях и предположениях гипотезы, представляющие собой векторы с компонентами  $\mu^* + \alpha_i^* + \beta_j^*$  и  $\mu^* + \beta_j^*$ , соответственно. Отсюда легко видеть (в силу теоремы Пифагора), что  $\mathcal{G}_\omega - \mathcal{G}_\Omega = \|\vec{\eta}_\Omega - \vec{\eta}_\omega\|^2$  и

$$SS_H = \mathcal{G}_\omega - \mathcal{G}_\Omega = \sum_{ij} (\mu^* + \alpha_i^* + \beta_j^* - \mu^* - \beta_j^*)^2 = \sum_{ij} (\alpha_i^*)^2 = p_2 \sum_{i=1}^{p_1} (\alpha_i^*)^2 .$$

Поскольку указанная сумма квадратов связана с гипотезой относительно фактора А, её обычно обозначают  $SS_A$  и называют суммой квадратов главных эффектов фактора А. Подставляя в это соотношение выражение для  $\alpha_i^*$  и воспользовавшись утверждением леммы III.1, получаем

$$SS_A = p_2 \sum_i (y_{i.} - y_{..})^2 = p_2 \sum_i y_{i.}^2 - p_1 p_2 y_{..}^2 .$$

Хотя гипотеза  $\mathbf{H}_A$  и задается  $p_1$  параметрическими функциями, среди них всего  $q = p_1 - 1$  линейно независимых, так как на них накладывается одна линейная связь. Поэтому критическая область F-критерия при проверке гипотезы  $\mathbf{H}_A$  определяется неравенством

$$\mathcal{F}_A = \overline{SS_A} / \overline{SS_e} > F_{p_1-1, \nu_e}^\alpha,$$

где  $\overline{SS}_A = SS_A / (p_1 - 1)$ .

Аналогично, F-критерий будет отвергать гипотезу  $\mathbf{H}_B$ , если

$$\mathcal{F}_B = \overline{SS_B} / \overline{SS_e} > F_{p_2-1, \nu_e}^\alpha,$$

где  $\overline{SS}_B = SS_B / (p_2 - 1)$  и  $SS_B = p_1 \sum_j (y_{.j} - y_{..})^2 = p_1 \sum_j y_{.j}^2 - p_1 p_2 y_{..}^2 .$

### Разбиение полной суммы квадратов

Как и в однофакторном  $\mathcal{D}_A$  здесь также полная сумма квадратов, символизирующая разброс всех данных, может быть разложена в сумму, компоненты которой обусловлены соответствующими факторами.

#### Теорема III.6. Полная сумма квадратов

$$SS_{\Pi} = \sum_{ij} (y_{ij} - y_{..})^2 = SS_A + SS_B + SS_e$$

*Доказательство.* Заметим, что  $SS_e = \|\vec{y} - \vec{\eta}_{\Omega}\|^2 = \|\vec{y}\|^2 - \|\vec{\eta}_{\Omega}\|^2$ . Следовательно,

$$\begin{aligned} SS_e &= \sum_{ij} y_{ij}^2 - \sum_{ij} (\mu^* + \alpha_i^* + \beta_j^*)^2 = \sum_{ij} y_{ij}^2 - \sum_{ij} [(\mu^*)^2 + (\alpha_i^*)^2 + (\beta_j^*)^2] = \\ &= \sum_{ij} y_{ij}^2 - p_1 p_2 y_{..}^2 - SS_A - SS_B, \end{aligned}$$

где удвоенные произведения в предпоследней сумме пропали в силу дополнительных условий  $\sum \alpha_i^* = \sum \beta_j^* = 0$ .  $\otimes$

С помощью разложения полной суммы квадратов упрощается вычисление остаточной суммы квадратов  $SS_e = SS_{\Pi} - SS_A - SS_B$ .

Соберем результаты в единую таблицу.

Таблица двухфакторного ДД

Источник разброса	$SS$	Ст.свободы	$\overline{SS}$	$\mathcal{F}$	$\alpha$
Фактор А (строки)	$SS_A$	$\nu_A = p_1 - 1$	$\overline{SS}_A$	$\frac{\overline{SS}_A}{\overline{SS}_e}$	$1 - F_{\nu_A, \nu_e}(\mathcal{F})$
Фактор В (столбцы)	$SS_B$	$\nu_B = p_2 - 1$	$\overline{SS}_B$	$\frac{\overline{SS}_B}{\overline{SS}_e}$	$1 - F_{\nu_B, \nu_e}(\mathcal{F})$
Ошибки	$SS_e$	$\nu_e = (p_1 - 1)(p_2 - 1)$	$\overline{SS}_e$	–	–
Полная сумма квадратов	$SS_{\Pi}$	$p_1 p_2 - 1$	–	–	–

**Замечание ix.** Предположение о равенстве нулю взаимодействий можно интерпретировать как добавление этих взаимодействий к ошибке. Поэтому в данном случае сумму квадратов ошибок иногда называют суммой квадратов взаимодействий  $SS_{AB}$ . Забегая вперед, отметим, что при полном уравнении регрессии формула для вычисления  $SS_{AB}$  и её число степеней свободы совпадают с формулами для  $SS_e$  при отсутствии взаимодействий.

### Сравнения главных эффектов

Как и при 1-факторном ДД здесь можно уточнить вывод в случае, если гипотеза об отсутствии эффекта воздействия фактора на отклик отвергается. Равенство нулю всех главных эффектов фактора А эквивалентно условию постоянства средних значений всех уровней фактора:  $A_i = \mu$ . Поэтому для уточнения причин отвержения гипотезы  $\mathbf{H}_A$  необходимо рассмотреть всевозможные сравнения средних уровней:  $\psi = \sum c_i A_i$ . МНК-оценка  $\psi$  равна  $\psi^* = \sum c_i y_{i.}$ . Так как случайные величины  $y_{i.}$  ( $i = 1, \dots, p_1$ ) независимы, то дисперсия

$$\mathbf{D}(\psi^*) = \sum_{i=1}^{p_1} c_i^2 \mathbf{D}(y_{i.}) = \frac{\sigma^2}{p_2} \sum_{i=1}^{p_1} c_i^2. \quad (\text{III.6})$$

Для того, чтобы применить теперь теорему III.2 и формулу (III.4), достаточно заменить в формуле (III.6)  $\sigma^2$  ее оценкой  $\mathfrak{S}^2 = \overline{SS}_e$  и выбрать  $q = p_1 - 1, n - r = \nu_e$ .

**Пример.** Исследовалась поглощаемость 5-ти сортов масла при изготовлении пончиков. С этой целью в течение 6 рабочих дней на каждом из сортов масла приготавливалось в день по одной партии (24

штуки) пончиков. Были получены следующие результаты:

Сорт	1	2	3	4	5
1	164	172	163	150	164
2	177	197	177	172	169
3	168	167	144	146	145
4	146	161	165	141	149
5	172	180	166	169	170
6	196	190	178	183	167

Требуется с 10% уровнем значимости

1) проверить гипотезу об отсутствии различия между истинными средними количества поглощенного масла при использовании 5 сортов масла;

2) проверить гипотезу о неизменности среднего количества поглощенного масла по дням недели;

Вычисления иллюстрирует следующая таблица.

-6:   2:   -7:   -20:   -6:	↔ -37: : 1369
: 36  : 4  : 49  : 400  : 36	↔ : 525:
7:   27:   7:   2:   -1:	↔ 42: : 1764
: 49  : 729  : 49  : 4  : 1	↔ : 832:
-2:   -3:   -26:   -24:   -25:	↔ -80: : 6400
: 4  : 9  :676  : 576  : 625	↔ :1890:
-24:   -9:   -5:   -29:   -21:	↔ -88: : 7744
: 576  : 81  : 25  : 841  : 441	↔ :1964:
2:   10:   -4:   -1:   0:	↔ 7: : 49
: 4  : 100  : 16  : 1  : 0	↔ : 121:
26:   20:   8:   13:   -3:	↔ 64: : 4096
: 676  : 400  : 64  : 169  : 9	↔ :1318:
↓ : ↓ ↓ : ↓ ↓ : ↓ ↓ : ↓ ↓ : ↓	↓ : ↓ : ↓
3:   47:   -27:   -59:   -56:	↔ -92: :21422
:1345  :1323  :879  :1991  :1112	↔ :6650:
9:   2209:   729:   3481:   3136:	↔ 9564: :

В эту таблицу сначала заносятся преобразованные данные (сдвинутые на 170) и их квадраты (по диагонали). После этого суммируются все 10 столбцов (по строкам) и 12 строк (по столбцам) и результаты записываются, соответственно, в правые 2 колонки и нижние 2 строки, которые также суммируются. Последний столбец и последняя строка содержат квадраты соответствующих сумм первых степеней по строкам и столбцам. Процесс суммирования для наглядности показан стрелками. Таким образом, получаем следующие результаты:

1. Число неизвестных параметров  $p_1 = 6$ ,  $p_2 = 5$ .
2. Общий объем выборки  $n = 30$ .

3. Сумма всех значений  $\sum_{ij} y_{ij} = -92$  и оценка генерального среднего

$$\mu^* = y_{..} = -3.07 .$$

4.  $ny_{..}^2 = 282.13$  .

5. Сумма всех квадратов  $\sum_{ij} y_{ij}^2 = 6650$  и полная сумма квадратов

$$SS_{\Pi} = 6650 - 282.13 = 6367.87 .$$

6. Сумма квадратов главных эффектов фактора А

$$SS_A = \sum_{i=1}^{p_1} \left( \sum_{j=1}^{p_2} y_{ij} \right)^2 - ny_{..}^2 = 21422/5 - 282.13 = 4002.27 .$$

7. Сумма квадратов главных эффектов фактора В

$$SS_B = \sum_{j=1}^{p_2} \left( \sum_{i=1}^{p_1} y_{ij} \right)^2 - ny_{..}^2 = 9564/6 - 282.13 = 1311.87 .$$

8. Сумма квадратов ошибок

$$SS_e = SS_{\Pi} - SS_A - SS_B = 1053.73$$

Запишем результаты в таблицу.

Таблица двухфакторного ДА

Источник разброса	$SS$	Ст. свободы	$\overline{SS}$	$\mathcal{F}$	$\alpha$	%
По дням	4002.27	5	800.454	15.2	< 0.0001	63%
По сортам	1311.87	4	327.97	6.23	0.002	20%
Ошибки	1053.73	20	52.67	-	-	17%
Полная сумма квадратов	6367.87	29	-	-	-	100%

Критический уровень значимости  $\alpha$  найден по таблице распределения Фишера со степенями свободы (5,20) для первого фактора и (4,20) для второго. Таким образом, налицо значимое изменение поглощаемости масла как по сортам, так и по дням недели.

Последний столбец таблицы содержит долю соответствующей суммы квадратов в процентах к полной сумме квадратов. В соответствии с устоявшейся на практике терминологией можно утверждать, что на 60% изменчивость поглощаемости масла обусловлена разбросом по дням недели и на 20% влиянием сорта масла.



# ДИСПЕРСИОННЫЙ АНАЛИЗ

Методические разработки  
по специальному курсу

## Часть 2

МНОГОФАКТОРНЫЙ ДИСПЕРСИОННЫЙ АНАЛИЗ.  
ОПТИМАЛЬНОЕ ПЛАНИРОВАНИЕ ФАКТОРНЫХ  
ЭКСПЕРИМЕНТОВ

Казань – 1998

## Двухфакторный дисперсионный анализ с одним наблюдением в ячейке с учетом взаимодействий

Обсудим теперь вопрос учета взаимодействий при проведении ДА с одним наблюдением в каждой ячейке. Как уже отмечалось, это возможно только при определенных предположениях относительно неизвестных параметров. Выше рассматривался случай полного отсутствия взаимодействий. Если все же мы хотим каким-либо образом учесть возможные взаимодействия при проведении ДА, нам необходимо построить новую модель регрессии с взаимодействиями, содержащими меньшее число параметров. Одна из таких моделей предполагает, что взаимодействия есть произведения главных эффектов

$$\gamma_{ij} = G\alpha_i\beta_j,$$

где  $G$  неизвестная константа. Таким образом, регрессия принимает вид  $\eta_{ij} = \mu + \alpha_i + \beta_j + G\alpha_i\beta_j$ , то есть является нелинейной по параметрам. Для этой регрессии задача нахождения ОМНК становится нетривиальной. Можно поступить следующим образом. Если предположить, что главные эффекты известны, тогда МНК-оценка для  $G$  легко находится и равна

$$G^* = \frac{\sum_{ij} \alpha_i \beta_j y_{ij}}{\sum_i \alpha_i^2 \sum_j \beta_j^2}.$$

Теперь, положив формально  $\alpha_i = \alpha_i^*$ ,  $\beta_j = \beta_j^*$  и  $\gamma_{ij}^* = G^* \alpha_i^* \beta_j^*$ , можно рассмотреть формальную сумму квадратов взаимодействий

$$SS_{\text{взаим}} = \sum_{ij} (\gamma_{ij}^*)^2 = \frac{\left[ \sum_{ij} \alpha_i^* \beta_j^* y_{ij} \right]^2}{\sum_i (\alpha_i^*)^2 \sum_j (\beta_j^*)^2} = \frac{p_1 p_2 K^2 \left[ \sum_{ij} \alpha_i^* \beta_j^* y_{ij} \right]^2}{SS_A SS_B}.$$

Основываясь на значениях  $SS_{\text{взаим}}$ , мы можем проверить гипотезу об отсутствии взаимодействий, отвергая гипотезу при больших  $SS_{\text{взаим}}$ .

**Теорема III.7.** Если уравнение регрессии не содержит взаимодействий –  $\eta_{ij} = \mu + \alpha_i + \beta_j$ , тогда

- 1)  $SS_{\text{взаим}} / \sigma^2 \sim \chi_1^2$ ;
- 2) Остаточная сумма квадратов

$$SS_{\text{ост}} = \sum_{ij} (y_{ij} - y_{i.} - y_{.j} + y_{..})^2 - SS_{\text{взаим}} = SS_e - SS_{\text{взаим}}$$

имеет хи-квадрат распределение с  $p_1 p_2 - p_1 - p_2$  степенями свободы.

- 3)  $SS_{\text{взаим}}$  и  $SS_{\text{ост}}$  независимы.

Эту теорему можно стандартным образом использовать для проверки гипотезы об отсутствии взаимодействий. Рассмотрим статистику

$$\mathcal{F} = \frac{SS_{\text{взаим}}}{SS_{\text{ост}}},$$

которая, в силу предыдущей теоремы, имеет распределение Фишера с одной и  $p_1 p_2 - p_1 - p_2$  степенями свободы. Таким образом, гипотезу об отсутствии взаимодействий следует отвергать, если критический уровень значимости, вычисленный по распределению Фишера ( $\alpha_{\text{кр.}} = 1 - F_{1, p_1 p_2 - p_1 - p_2}(\mathcal{F})$ ) будет меньше заданного уровня  $\alpha$ .

При подтверждении гипотезы дальнейший дисперсионный анализ проводится для модели регрессии без взаимодействий. Значительные сложности начинаются, если гипотеза об отсутствии взаимодействий отвергается. Мы не будем вдаваться в подробности статистического вывода в этом случае. Заметим только, что здесь приходится ограничиваться приближенным критерием.

### Двухфакторный анализ с одинаковым числом наблюдений в ячейке

Пусть теперь в каждой  $ij$ -ой ячейке производилось по  $K_{ij}$  наблюдений (например, каждый день на каждом сорте масла изготавливалось несколько партий пончиков). Обозначим  $y_{ijk}$  —  $k$ -ое наблюдение в  $ij$ -ой ячейке. Тогда основные предположения записываются в виде

$$\Omega : \begin{cases} y_{ijk} = \eta_{ij} + e_{ijk}, \\ \{e_{ijk}\} \text{ независимы и распределены } \mathcal{N}(0, \sigma^2), \\ k = 1, \dots, K_{ij}; (ij) \in D, \end{cases}$$

где через  $D$  обозначено множество всех пар (– ячеек)  $(ij)$ , для которых было произведено хотя бы одно наблюдение ( $K_{ij} > 0$ ). Итак, мы имеем  $n = \sum_{ij} K_{ij}$  наблюдений, а общее число неизвестных параметров ( $\eta_{ij}$ ) равно количеству элементов множества  $D$ . Для нахождения ОМНК при  $\Omega$  необходимо минимизировать по  $\eta_{ij}$  функцию

$$\mathcal{G} = \sum_{ij \in D} \sum_{k=1}^{K_{ij}} (y_{ijk} - \eta_{ij})^2.$$

Очевидно, здесь существует единственный минимум, и ОМНК  $\eta_{ij}^* = y_{ij.}$  для всех непустых ячеек (легко проверить дифференцированием по  $\eta_{ij}$ ). Отсюда получаем, что, во-первых, в силу единственности ОМНК ранг задачи  $r = p$  (см. следствие 1 из теоремы о существовании ОМНК) и, во-вторых, сумма квадратов ошибок

$$SS_e = \mathcal{G}_\Omega = \sum_{ij \in D} \sum_{k=1}^{K_{ij}} (y_{ijk} - y_{ij.})^2,$$

а её число степеней свободы  $\nu_e = n - p$ .

Как отмечалось, все главные эффекты однозначно определяются значениями  $\eta_{ij}$ . Поэтому их МНК-оценки легко получить простой подстановкой ОМНК  $\eta_{ij}^*$  в уравнения для главных эффектов. Иными словами, МНК-оценки

$$\begin{aligned}\mu^* &= \eta_{..}^* = y_{...}, \\ \alpha_i^* &= \eta_{i.}^* - \eta_{..}^* = y_{i..} - y_{...}, \quad \beta_j^* = \eta_{.j}^* - \eta_{..}^* = y_{.j.} - y_{...}, \\ \gamma_{ij}^* &= \eta_{ij}^* - \eta_{i.}^* - \eta_{.j}^* + \eta_{..}^* = y_{ij.} - y_{i..} - y_{.j.} + y_{...}.\end{aligned}$$

## Проверка гипотез

Рассмотрим теперь задачу проверки стандартных гипотез  $D_A$ :

$H_A$ : все  $\alpha_i = 0$  – фактор А не влияет на отклик;

$H_B$ : все  $\beta_j = 0$  – фактор В не влияет на отклик;

$H_{AB}$ : все  $\gamma_{ij} = 0$  – между факторами отсутствует взаимодействие.

Дальнейшие выкладки резко усложнятся, если не все  $K_{ij}$  одинаковы. Поэтому, с целью упрощения, предположим, что все  $K_{ij} = K > 1$ . В этом случае общее число наблюдений  $n = Kp_1p_2$  при  $p = p_1p_2$  параметрах. Число степеней свободы  $SS_e$  равно  $\nu_e = n - p = p_1p_2(K - 1)$ .

Для каждой из гипотез необходимо минимизировать функцию

$$\mathcal{G} = \sum_{ijk} (y_{ijk} - \mu - \alpha_i - \beta_j - \gamma_{ij})^2.$$

Сначала представим эту функцию в удобной для наших целей форме.

**Лемма III.2.** *Сумма квадратов расхождений*

$$\begin{aligned}\mathcal{G} &= SS_e + p_1p_2K(\mu^* - \mu)^2 + \\ &+ p_2K \sum_i (\alpha_i^* - \alpha_i)^2 + p_1K \sum_j (\beta_j^* - \beta_j)^2 + K \sum_{ij} (\gamma_{ij}^* - \gamma_{ij})^2.\end{aligned}$$

*Доказательство.* Указанное разложение легко получается возведением в квадрат и суммированием следующего очевидного представления

$$\begin{aligned}y_{ij} - \mu - \alpha_i - \beta_j - \gamma_{ij} &= (y_{ij} - \mu^* - \alpha_i^* - \beta_j^* - \gamma_{ij}^*) + \\ &+ (\mu^* - \mu) + (\alpha_i^* - \alpha_i) + (\beta_j^* - \beta_j) + (\gamma_{ij}^* - \gamma_{ij}).\end{aligned}$$

Как обычно, здесь удвоенные суммы квадратов пропадают в силу дополнительных соотношений для главных эффектов.  $\otimes$

Чтобы найти теперь минимум  $\mathcal{G}$  при гипотезе  $\mathbf{H}_A$ , необходимо положить все  $\alpha_i = 0$  и минимизировать по оставшимся переменным с учетом дополнительных условий. Однако ясно, что этот минимум будет достигаться, если выбрать параметры так, чтобы соответствующие слагаемые в разложении  $\mathcal{G}$  обратились в нуль, то есть положить  $\mu = \mu^*$ ,  $\beta_j = \beta_j^*$ ,  $\gamma_{ij} = \gamma_{ij}^*$ . При этом, сумма квадратов

$$\mathcal{G}_{\omega_A} = SS_e + p_2 K \sum_{i=1}^{p_1} (\alpha_i^*)^2 \quad SS_A = p_2 K \sum_{i=1}^{p_1} (\alpha_i^*)^2 .$$

Гипотеза  $H_A$  описывается  $p_1$  линейными функциями от неизвестных параметров, между которыми существует одна линейная связь. Поэтому число степеней свободы  $SS_A$  равно  $p_1 - 1$ . Аналогично, для проверки гипотезы  $\mathbf{H}_B$  сумма квадратов ошибок

$$SS_B = p_1 K \sum_{j=1}^{p_2} (\beta_j^*)^2$$

с  $p_2 - 1$  степенью свободы.

Для гипотезы  $\mathbf{H}_{AB}$  сумма квадратов  $SS_{AB} = K \sum_{i=1}^{p_1} \sum_{j=1}^{p_2} (\gamma_{ij}^*)^2$ .

Число степеней свободы  $SS_{AB}$  легко найти, если расположить все функции, описывающие эту гипотезу, в матрицу  $p_1 \cdot p_2$  и заметить, что в силу линейных связей на  $\gamma_{ij}$  первая строка и первый столбец этой матрицы линейно выражаются через элементы оставшейся  $(p_1 - 1) \cdot (p_2 - 1)$  матрицы, которые между собой уже линейно независимы. Таким образом, число степеней свободы  $SS_{AB}$  равно  $(p_1 - 1)(p_2 - 1)$ . Приведем упрощенный способ проведения вычислений, необходимых для  $\mathcal{D}_A$ .

## Вычисления

Во-первых, наряду с таблицей данных полезно составить новую таблицу из средних значений каждой ячейки  $y_{ij}$  и их квадратов. Эта новая таблица похожа на таблицу данных 2-факторного  $\mathcal{D}_A$  с одним наблюдением в ячейке. Сумма квадратов ошибок  $SS_e$  для таблицы средних, вычисленная по формуле для  $\mathcal{D}_A$  с одним наблюдением в ячейке, совпадает с суммой квадратов взаимодействий, вычисленной по формуле для  $\mathcal{D}_A$  с  $K$  наблюдениями в ячейке с точностью до множителя  $K$  (вспомним, что сумма квадратов ошибок при  $\mathcal{D}_A$  с одним наблюдением в ячейке называлась суммой квадратов взаимодействий). Так же, с точностью до множителя  $K$ , совпадают формулы для  $SS_A$  и  $SS_B$  при построении  $\mathcal{D}_A$  с одним наблюдением и с  $K$  наблюдениями в ячейке. Если суммы квадратов, относящиеся к таблице средних,

снабдить штрихом, то будет иметь место разбиение

$$SS_{\Pi}' = \sum_{ij} y_{ij}^2 - p_1 p_2 y_{...}^2 = SS_A' + SS_B' + SS_e'.$$

Положим  $SS_{\text{ячеек}} = K SS_{\Pi}'$ . Тогда, как отмечалось,

$$SS_A = K SS_A', \quad SS_B = K SS_B', \quad SS_{AB} = K SS_e'$$

и

$$SS_{AB} = SS_{\text{ячеек}} - SS_A - SS_B. \quad (\text{III.7})$$

Для вычисления  $SS_e$  заметим, что полная сумма квадратов исходной таблицы  $SS_{\Pi} = SS_e + SS_{\text{ячеек}}$ .

Теперь мы можем представить результаты в виде единой таблицы.

Таблица двухфакторного  $\mathcal{D}_A$  с  $K$  наблюдениями в ячейке

Источник разброса	$SS$	Ст.свободы	$\overline{SS}$	$\mathcal{F}$	$\alpha$
Фактор А (строки)	$SS_A$	$\nu_A = p_1 - 1$	$\overline{SS}_A$	$\frac{\overline{SS}_A}{\overline{SS}_e}$	$1 - F_{\nu_A, \nu_e}(\mathcal{F})$
Фактор В (столбцы)	$SS_B$	$\nu_B = p_2 - 1$	$\overline{SS}_B$	$\frac{\overline{SS}_B}{\overline{SS}_e}$	$1 - F_{\nu_B, \nu_e}(\mathcal{F})$
Взаимодействия	$SS_{AB}$	$\nu_{AB} = (p_1 - 1)(p_2 - 1)$	$\overline{SS}_{AB}$	$\frac{\overline{SS}_{AB}}{\overline{SS}_e}$	$1 - F_{\nu_{AB}, \nu_e}(\mathcal{F})$
Ошибки	$SS_e$	$\nu_e = p_1 p_2 (K - 1)$	$\overline{SS}_e$	—	—
Полная сумма квадратов	$SS_{\Pi}$	$K p_1 p_2 - 1$	—	—	—

Рекомендуемый процесс вычислений будет состоять из следующих шагов.

1) Вычисляются средние значения по каждой ячейке  $y_{ij}$ . и заносятся в новую таблицу вместе со своими квадратами  $y_{ij}^2$ .

2) По новой таблице вычисляются  $y_{...}$ ,  $SS_{\text{ячеек}}$ ,  $SS_A$ ,  $SS_B$ .

3) Вычисляется  $SS_{AB}$  по формуле (III.7).

4) По исходной таблице находится сумма квадратов всех наблюдений и вычисляется  $SS_{\Pi}$ .

5) Вычисляется  $SS_e = SS_{\Pi} - SS_{\text{ячеек}}$ .

Попутно нами доказано стандартное разложение для полной суммы квадратов

$$SS_{\Pi} = SS_A + SS_B + SS_e.$$

**Пример** 2-факторного дисперсионного анализа с 3 наблюдениями в ячейке

На фабрике по изготовлению мясных консервов наполнение банок происходит на 6 ротационных машинах. Партии мяса поступают от 5 поставщиков. Для проверки наполняемости банок от каждой партии и каждой машины были случайно отобраны 3 наполненных банки и взвешены. Результаты представлены в следующей таблице.

Поставщик Цилиндр	1	2	3	4	5
1	501, 501, 502	504, 503, 505	506, 503, 507	503, 501, 503	501, 503, 503
2	499, 503, 499	498, 501, 500	501, 503, 505	502, 501, 500	501, 500, 501
3	501, 501, 501	502, 500, 501	502, 503, 504	501, 503, 503	503, 503, 503
4	498, 503, 500	498, 500, 501	503, 503, 504	500, 500, 502	500, 501, 501
5	501, 501, 499	502, 501, 505	500, 501, 502	501, 500, 499	498, 503, 501
6	501, 501, 500	500, 503, 500	503, 503, 504	503, 500, 502	501, 503, 502

Проверить с уровнем значимости 10% гипотезы о том, что

- а) средний вес наполненной банки не зависит от поставщика;
- б) средний вес наполненной банки не изменяется от машины к машине;
- в) отсутствует взаимодействие между типом поставщика и ротационной машиной.

## § 5. Полный многофакторный анализ с взаимодействиями

В этом разделе на примере 3-факторного ДА мы распространим понятия, введенные для 2-факторного анализа, на случай многих факторов, влияющих на отклик.

Пусть на результат эксперимента действуют три фактора А, В и С, причем, фактор А разбит на  $p_A$  уровней, фактор В на  $p_B$  уровней, фактор С на  $p_C$  уровней. Первым шагом построения статистического вывода представим регрессию в виде суммы главных эффектов и взаимодействий. С целью унификации будем обозначать все компоненты регрессии буквой  $\alpha$  с верхним индексом, указывающим на соответствующий фактор. Регрессия равна

$$\eta_{ijk} = \mu + \alpha_i^A + \alpha_j^B + \alpha_k^C + \alpha_{ij}^{AB} + \alpha_{ik}^{AC} + \alpha_{jk}^{BC} + \alpha_{ijk}^{ABC}.$$

Здесь

$$\mu = \eta_{...} \text{ — генеральное среднее;} \quad (\text{III.8})$$

$$\left. \begin{aligned} \alpha_i^A &= \eta_{i..} - \eta_{...} \\ \alpha_j^B &= \eta_{.j.} - \eta_{...} \\ \alpha_k^C &= \eta_{..k} - \eta_{...} \end{aligned} \right\} \begin{array}{l} \text{главные эффекты} \\ \text{факторов А, В, С;} \end{array} \quad (\text{III.9})$$

$$\left. \begin{aligned} \alpha_{ij}^{AB} &= \eta_{ij.} - \eta_{i..} - \eta_{.j.} + \eta_{...} \\ \alpha_{ik}^{AC} &= \eta_{i.k} - \eta_{i..} - \eta_{..k} + \eta_{...} \\ \alpha_{jk}^{BC} &= \eta_{.jk} - \eta_{.j.} - \eta_{..k} + \eta_{...} \end{aligned} \right\} \begin{array}{l} \text{двухфакторные} \\ \text{взаимодействия;} \end{array} \quad (\text{III.10})$$

$$\alpha_{ijk}^{ABC} = \eta_{ijk} - \eta_{ij.} - \eta_{i.k} - \eta_{.jk} + \eta_{i..} + \eta_{.j.} + \eta_{..k} - \eta_{...} \quad (\text{III.11})$$

— трехфакторное взаимодействие .

Из этих формул легко вывести схему построения всех взаимодействий (на случай более чем трех факторов):

а) взаимодействие разбивается на группы слагаемых с чередующимися знаками “+”, “−”, ...;

б) первая группа (со знаком “+”) содержит всего одно слагаемое, равное регрессии, усредненной по индексам, соответствующим факторам, не участвующим в рассматриваемом взаимодействии (например,  $\eta_{ij.}$ );

в) вторая группа слагаемых (со знаком “−”) содержит усреднения регрессии из пункта а) по всем сочетаниям одного индекса;

г) третья группа (со знаком “+”) усреднения регрессии из пункта а) по двум индексам;

д) и т.д.

Предположим теперь, что при каждом сочетании уровней (в каждой ячейке) было произведено по  $M$  наблюдений и пусть  $y_{ijkm}$  есть  $m$ -ое наблюдение на  $i$ -ом уровне фактора А,  $j$ -ом уровне фактора В и  $k$ -ом уровне фактора С. Таким образом, всего получено  $Mp_A p_B p_C$  наблюдений, по которым требуется оценить регрессию  $\eta_{ijk}$  с  $p = p_A p_B p_C$  параметрами. ОМНК для параметров регрессии находятся стандартным способом с помощью подстановки в уравнения (III.8–III.11) оценки  $\eta_{ijk}^* = y_{ijk.}$ . Например, ОМНК

$$\hat{\alpha}_{ij}^{AB} = y_{ij..} - y_{i...} - y_{.j..} + y_{....} .$$

Суммы квадратов ошибок при проверке гипотез о главных эффектах и взаимодействиях легко вычислить, воспользовавшись следующей мнемонической формулой

$$SS = \sum_{ijkm} (\hat{\alpha})^2 .$$



Так, например,

$$SS_{AC} = \sum_{ijkm} (\hat{\alpha}_{ik}^{AC})^2 = Mp_B \sum_{ik} (\hat{\alpha}_{ik}^{AC})^2.$$

Число степеней свободы суммы квадратов главных эффектов, как обычно, на единицу меньше числа уровней фактора (например,  $p_A - 1$  для фактора А), а число степеней свободы суммы квадратов взаимодействий равно произведению степеней свободы факторов, участвующих в этом взаимодействии ( $(p_B - 1)(p_C - 1)$  для  $SS_{BC}$ ).

Сумма квадратов при числе наблюдений в ячейке  $M > 1$  равна

$$SS_e = \sum_{ijkm} (y_{ijkm} - y_{\dots})^2 = \sum_{ijkm} y_{ijkm}^2 - ny_{\dots}^2$$

и её число степеней свободы равно  $\nu_e = n - p = p_A p_B p_C (M - 1)$ .

В случае  $M = 1$  ДА возможен только при определенных предположениях относительно значений параметров регрессии. Сокращение числа параметров начинают, обычно, с предположения, что равны нулю все старшие взаимодействия:  $\alpha_{ijk}^{ABC} = 0, \forall i, j, k$ . Сумма квадратов ошибок в этом случае равна сумме квадратов самого старшего взаимодействия ( $-SS_{ABC}$  при 3-факторном ДА). Можно отбросить несколько групп взаимодействий. Тогда  $SS_e$  будет равна сумме  $SS$  всех отброшенных взаимодействий, а число степеней свободы  $SS_e$  – сумме соответствующих степеней свободы.

### Разбиение полной суммы квадратов

Иногда ДА строят на формальном разбиении полной суммы квадратов, которое для 3-факторной модели принимает вид

$$\begin{aligned} SS_{\Pi} &= \sum_{ijkm} y_{ijkm}^2 - ny_{\dots}^2 = \\ &= SS_A + SS_B + SS_C + SS_{AB} + SS_{AC} + SS_{BC} + SS_{ABC} + SS_e. \end{aligned}$$

К этому разложению применяют следующую теорему Cochran'a.

**Теорема III.8.** Пусть  $Q$  есть квадратичная форма полного ранга  $n$  от  $n$  независимых нормальных  $\mathcal{N}(0, 1)$  случайных величин, причем имеет место разложение  $Q$  на  $s$  квадратичных форм

$$Q = Q_1 + \dots + Q_s.$$

Обозначим через  $n_j$  ранг квадратичной формы  $Q_j$ . Тогда случайные величины  $Q_1, \dots, Q_s$  независимы и имеют  $\chi^2$ -квадрат распределение с  $n_1, \dots, n_s$  степенями свободы, соответственно, в том и только в том случае, если  $n = n_1 + \dots + n_s$ .

Предположим, что нам необходимо проверить гипотезу относительно взаимодействий факторов А и В. Проверим выполнимость условия теоремы Cochran'a при разложении полной суммы квадратов. Очевидно, сумма степеней свободы правой части (III.12) равна  $(p_A p_B p_C M - 1)$  – степени свободы  $SS_{\Pi}$ . Следовательно, отношение  $\mathcal{F} = \overline{SS}_{AB} / \overline{SS}_e$  имеет распределение Фишера с  $(p_A - 1)(p_B - 1)$  и  $p_A p_B p_C (M - 1)$  степенями свободы, и гипотеза  $\mathbf{H}_{AB} : \alpha_{ij}^{AB} = 0$  об отсутствии взаимодействий между А и В отвергается, если значение статистики  $\mathcal{F}$  достаточно велико.

**Пример.** Провести анализ влияния на содержание влаги в пищевом продукте четырех факторов:

- фактор А – сорт соли – 3 уровня;
- фактор В – количество соли – 3 уровня;
- фактор С – количество кислоты – 2 уровня;
- фактор D – пищевые примеси – 2 уровня.

По результатам нижеприведенных измерений необходимо:

а) составить таблицу ДД. Полезно рассмотреть все  $\overline{SS}$ , даже если часть из них будет объединена с  $SS_e$ : относительная величина  $\overline{SS}$  дает дополнительную информацию о влиянии факторов;

б) объединить  $SS$  всех трехфакторных и четырехфакторных взаимодействий с  $SS$  ошибок. Выявить значимость влияния оставшихся взаимодействий и главных эффектов на отклик.

Уровни А	Уровни В	Уровни С			
		1		2	
		Уровни D		Уровни D	
		1	2	1	2
1	1	8	5	8	4
	2	17	11	13	10
	3	22	16	20	15
2	1	7	3	10	5
	2	26	17	24	19
	3	34	32	34	29
3	1	10	5	9	4
	2	24	14	24	16
	3	39	33	36	34

Обработку этих данных произведем воспользовавшись возможностями пакета «STATGRAFICS». Данные занесем в файл «VLAGA» с переменной Y, содержащей количество влаги, и переменными SolSort, SolKol, Kislota и Primesi, содержащими номера уровней соответствующих факторов. Для примера приведем строку редактора данных  $\mathcal{S}_G$ ,

в которую внесено количество влаги в пищевом продукте при использовании второго сорта соли, первого уровня количества соли, первого же уровня кислоты и второго уровня примесей:

Row	Y	SolSort	SolKol	Kislota	Primesi
13	3	2	1	1	2

Для проведения четырехфакторного ДА следует в разделе «Analysis of Variance» главного меню  $\mathcal{S}\mathcal{G}$  выбрать подраздел «Multifactor Analysis of Variance». Поясним некоторые графы, которые необходимо заполнить перед запуском программы на выполнение.

1) В графу «Data» заносится имя переменной, содержащей наблюдаемый в эксперименте отклик – «VLAGA.Y», в графу «Factors» управляемые факторы – «VLAGA.SolSort», «VLAGA.SolKol», «VLAGA.Kislota», «VLAGA.Primesi».

2) Графа «Covariates» в данном случае остается пустой, поскольку должна содержать отсутствующие в нашем эксперименте непрерывные составляющие регрессии.

3) В графе «Interactions» (взаимодействия) необходимо убрать или оставить значки \* в зависимости от того, предполагается учет соответствующего двухфакторного взаимодействия или нет. Взаимодействия более высокого порядка в пакете  $\mathcal{S}\mathcal{G}$  автоматически игнорируются. Для данной задачи все значки могут быть оставлены без изменения.

Выбрав метод Шеффе построения доверительных интервалов, произведем вычисления (клавиша **F6**). На экран будет выведена следующая таблица ДА:

Source of var.	Sum of Sq.	d.f.	Mean sq.	F-ratio	Sig. level	%
MAIN EFFECTS	3633.83	6	605.63	304.93	.00	90.1
SolSort	495.05	2	247.52	124.62	.00	12.2
SolKol	2905.38	2	1452.69	731.42	.00	72.0
Kislota	3.36	1	3.36	1.69	.21	0.0
Primesi	230.02	1	230.02	115.81	.00	5.7
INTERACTIONS	364.69	13	28.053419	14.12	.00	9.0
SolSort SolKol	333.11	4	83.27	41.93	.00	8.2
SolSort Kislota	3.72	2	1.86	.93	.41	0.0
SolKol Kislota	6.05	2	3.02	1.52	.24	0.1
SolSort Primesi	4.05	2	2.02	1.02	.38	0.1
SolKol Primesi	14.38	2	7.19	3.62	.05	0.3
Kislota Primesi	3.36	1	3.36	1.69	.21	0.0
RESIDUAL	31.78	16	1.99			0.8
TOTAL (CORR.)	4030.31	35				100

Анализ предпоследнего столбца этой таблицы (Sig.Level – Уровень значимости) показывает, что все факторы, за исключением третьего, значимо влияют на содержание влаги. Из взаимодействий значимое отличие от нуля имеют только сочетания факторов АВ (сорт соли – количество соли) и ВD(количество соли – примеси). Исключив незначимые факторы и взаимодействия, получаем следующее уравнение регрессии

$$E y_{ijkl} = \eta_{ijkl} = \mu + \alpha_i^A + \alpha_j^B + \alpha_l^D + \alpha_{ij}^{AB} + \alpha_{jl}^{BD}.$$

Последний столбец таблицы (в пакете  $\mathcal{SG}$  не вычисляется) содержит долю в процентах суммы квадратов того или иного фактора в полной сумме квадратов. Этот показатель наглядно демонстрирует степень влияния фактора на отклик. Отсюда можно сделать вывод, что основное влияние на отклик оказывает количество соли.

Из дополнительных возможностей, вызываемых по клавише **F5** рассмотрим только опцию «Means Table», по которой выдается таблица средних значений с доверительными интервалами. Приведем первую страницу этой таблицы для регрессии, содержащей все факторы.

Level	Count	Average	Std. Err (internal)	Std.Err (pooled s)	95% Conf. for mean	
SolSort						
1	12	12.416	1.658	.406	11.554	13.279
2	12	19.833	3.267	.406	18.970	20.695
3	12	20.666	3.654	.406	19.804	21.529
SolKol						
1	12	6.500	.712	.406	5.637	7.362
2	12	17.916	1.588	.406	17.054	18.779
3	12	28.500	2.343	.406	27.637	29.362
Kislota						
1	18	17.944	2.652	.332	17.240	18.648
2	18	17.333	2.474	.332	16.628	18.037
Primesi						
1	18	20.167	2.478	.332	19.462	20.871
2	18	15.111	2.504	.332	14.406	15.815
Total	36	17.638	.2348	.2348	17.140	18.136

Из этой таблицы также видно, что главные эффекты уровней третьего фактора не отличаются между собой (их доверительные интервалы накладываются друг на друга). Значимое влияние первого фактора обеспечивается только эффектом первого уровня.

К сожалению, из этих таблиц нельзя установить значимость влияния взаимодействий высокого порядка. Для этого придется проводить дополнительные вычисления.

## Г л а в а I V

### МОДЕЛИ СО СЛУЧАЙНЫМИ ФАКТОРАМИ

Для моделей со случайными факторами (модель II ДА) не существует общей теории, аналогичной вышеразвитой теории для моделей с постоянными факторами. Здесь мы рассмотрим два случая ДА со случайными компонентами – однофакторный и двухфакторный ДА.

#### § 1. Однофакторный дисперсионный анализ

Представим следующий производственный эксперимент. На контроль поступают образцы изделий, произведенные по одной и той же технологии (например, лопатки турбины, изготовленные в одной и той же печи). Над каждым образцом проводится по несколько контрольных замеров. Результат  $k$ -ого измерения  $i$ -ого образца можно представить в виде

$$y_{ik} = \theta_i + \varepsilon_{ik},$$

где  $\theta_i$  – истинное значение измеряемой характеристики качества, случайно изменяющееся от образца к образцу,  $\varepsilon_{ik}$  – случайная ошибка измерения. Если определить эффект  $i$ -ого образца как разность между  $\theta_i$  и теоретическим средним  $\mu = \mathbf{E} \theta_i$  случайной величины  $\theta_i$ :  $\alpha_i = \theta_i - \mu$ , то наша модель может быть переписана в стандартном для однофакторного ДА виде  $y_{ik} = \mu + \alpha_i + \varepsilon_{ik}$  с единственным отличием, что эффекты случайны. Мы увидим далее, что таблица однофакторного ДА в модели II совпадает с таблицей однофакторного ДА в модели I, однако обоснования для этих таблиц разнятся.

Предположим, как обычно, что эффекты и ошибки распределены нормально и независимы. Кроме того, предположим, что всего было исследовано  $I$  образцов и на каждом образце было сделано по  $K$  замеров. В этом случае основные предположения принимают вид

$$\Omega : \begin{cases} y_{ik} = \mu + \alpha_i + \varepsilon_{ik}, \quad i = \overline{1, I}, k = \overline{1, K}, \\ \{\alpha_i\}, \{\varepsilon_{ik}\} \text{ — независимы в совокупности,} \\ \alpha_i \sim \mathcal{N}(0, \sigma_A^2), \quad \varepsilon_{ik} \sim \mathcal{N}(0, \sigma_e^2). \end{cases}$$

**З а м е ч а н и е** **i.** Другое существенное отличие от модели I заключается в том, что измерения на одном образце зависимы. Действительно, при различных индексах  $j$  и  $k$  ковариация  $\text{Cov}(y_{ij}, y_{ik}) = \mathbf{E}(y_{ij} - \mu)(y_{ik} - \mu) = \mathbf{E}(\alpha_i + \varepsilon_{ij})(\alpha_i + \varepsilon_{ik}) = \mathbf{E} \alpha_i^2 = \sigma_A^2$ , так как остальные математические ожидания равняются нулю в силу независимости рассматриваемых элементов и равенства нулю их средних значений.

**З а м е ч а н и е** **ii.** Так же другое значение, отличное от дисперсии ошибки, принимает дисперсия наблюдений

$$\mathbf{D}(y_{ik}) = \mathbf{D}(\alpha_i + \varepsilon_{ik}) = \mathbf{D} \alpha_i + \mathbf{D} \varepsilon_{ik} = \sigma_A^2 + \sigma_e^2.$$

**О п р е д е л е н и е** **1.** Коэффициент корреляции между наблюдениями на одном образце

$$\rho = \frac{\sigma_A^2}{\sigma_A^2 + \sigma_e^2}$$

называется *внутриклассовой корреляцией*.

Сумма квадратов ошибок в основных предположениях для этой модели находится так же, как для модели I:

$$SS_e = \sum_{ik} (y_{ik} - y_{i.})^2.$$

Для того чтобы найти ее распределение, подставим модельные предположения из  $\Omega$  в выражение для  $SS_e$ . Получим

$$SS_e = \sum_{ik} (\varepsilon_{ik} - \varepsilon_{i.})^2. \quad (\text{IV.1})$$

Рассмотрим теперь фиктивную стандартную модель однофакторного ДД с постоянными факторами

$$\Omega' : \begin{cases} \varepsilon_{ik} = a_i + e_{ik}, \quad i = \overline{1, I}, k = \overline{1, K}, \\ \{e_{ik}\} \text{ — независимы в совокупности,} \\ e_{ik} \sim \mathcal{N}(0, \sigma_e^2). \end{cases}$$

Для этой модели  $\mathcal{G}_{\Omega'}$  совпадает с выше рассмотренной суммой квадратов (IV.1):  $\mathcal{G}_{\Omega'} = SS_e$ . Исходя из общей теории, можно утверждать, что  $\mathcal{G}_{\Omega'}/\sigma_e^2 \sim \chi_{\nu_e}^2$ , где  $\nu_e = IK - I = I(K - 1)$ . Таким образом, сумма квадратов ошибок  $SS_e \sim \sigma_e^2 \chi_{\nu_e}^2$  и

$$\overline{SS_e} = \frac{1}{I(K - 1)} \sum_{ik} (y_{ik} - y_{i.})^2$$

есть несмещенная оценка дисперсии  $\sigma_e^2$ .

Перейдем теперь к рассмотрению вопроса различения гипотез. Для этой модели факторы случайны, поэтому гипотеза об отсутствии влияния фактора на отклик может быть сформулирована только в терминах дисперсии фактора:  $\mathbf{H} : \sigma_A^2 = 0$ . Сумма квадратов ошибок, обусловленная отклонением от гипотезы, как и в модели I однофакторного  $\mathcal{D}_A$ , равна

$$SS_A = \sum_{ik} (y_{i.} - y_{..})^2 = K \sum_i (y_{i.} - y_{..})^2.$$

Для того чтобы найти её распределение, положим  $g_i = \alpha_i + \varepsilon_{i.}$ . Очевидно, все  $g_i$  независимы и имеют нормальное распределение со средними ноль и дисперсией  $\sigma_g^2 = \mathbf{D} g_i = \sigma_A^2 + \sigma_e^2/K$ . Поэтому распределение

$$SS_A = K \sum_i (g_i - g_{..})^2 \sim K \sigma_g^2 \chi_{I-1}^2 = (K\sigma_A^2 + \sigma_e^2) \chi_{I-1}^2.$$

Поскольку при гипотезе  $\mathbf{H}$  средние  $SS_e$  и  $SS_A$  совпадают и равны  $\sigma_e^2$ , то для построения  $F$ -критерия осталось показать, что статистики  $SS_e$  и  $SS_A$  независимы. Этот факт сразу вытекает из стандартного разложения полной суммы квадратов и теоремы Cochran'a (III.8). Таким образом, если верна гипотеза  $\mathbf{H}$ , то статистика

$$\mathcal{F} = \frac{\overline{SS}_A}{\overline{SS}_e} = \frac{\overline{SS}_A / \sigma_e^2}{\overline{SS}_e / \sigma_e^2} \sim F_{I-1, K-1}$$

и критическая область  $\mathcal{F} \geq F_{I-1, K-1}^\alpha$  имеет уровень  $\alpha$ .

Статистику  $\mathcal{F}$  можно использовать также для проверки гипотез более общего вида:  $\mathbf{H}_{\theta_0} : \sigma_A^2 \leq (\geq) \theta_0 \sigma_e^2$ , с заданной границей  $\theta_0$ . Критическая область для  $\mathbf{H}_{\theta_0}$  имеет вид  $\mathcal{F} \geq C$ , где критическая константа определяется из соотношения

$$\mathbf{P}\{\mathcal{F} \geq C\} \leq \alpha \quad (\text{IV.2})$$

для всех параметров, определяющих модель и удовлетворяющих условиям гипотезы. Статистика

$$\mathcal{F} = \frac{\overline{SS}_A}{\overline{SS}_e} \sim \frac{K\sigma_A^2 + \sigma_e^2}{\sigma_e^2} \frac{\chi_{I-1}^2 / (I-1)}{\chi_{I(K-1)}^2 / (I(K-1))} \sim (1 + K\theta) F_{I-1, I(K-1)}$$

с параметром  $\theta = \sigma_A^2 / \sigma_e^2$ . Очевидно, вероятность в (IV.2) монотонно возрастает с увеличением параметра  $\theta$ . Поэтому для того, чтобы выполнялось неравенство (IV.2) при всех  $\theta \leq \theta_0$ , достаточно выбрать константу  $C$  из условия  $\mathbf{P}\{F_{I-1, I(K-1)} \geq C / (1 + K\theta_0)\} = \alpha$ . Т.е.

$$C = (1 + K\theta_0) F_{I-1, I(K-1)}^\alpha.$$

## Оценки компонент дисперсии

Поскольку среднее значение хи-квадрат случайной величины равно числу ее степеней свободы, то среднее

$$\mathbf{E} \overline{SS_e} = \mathbf{E} \frac{\sigma_e^2 \chi_{\nu_e}^2}{\nu_e} = \sigma_e^2$$

и, аналогично,  $\mathbf{E} \overline{SS_A} = \sigma_e^2 + K\sigma_A^2$ . Поэтому можно выбрать в качестве несмещенных оценок для компонент дисперсии  $\sigma_e^2$  и  $\sigma_A^2$  статистики

$$\hat{\sigma}_e^2 = \overline{SS_e}, \quad \hat{\sigma}_A^2 = \frac{1}{K}(\overline{SS_A} - \overline{SS_e}).$$

Ясно, что с положительной вероятностью оценка  $\hat{\sigma}_A^2$  может принимать отрицательные значения, в то время как сам параметр всегда неотрицателен. Иногда эту оценку видоизменяют, полагая ее равной некоторой положительной величине, когда оценка отрицательна. Такая модификация делает оценку смещенной, и изучение ее свойств значительно усложняется. Мы остановимся здесь только на исследовании состоятельности неизменной оценки  $\hat{\sigma}_A^2$ .

Напомним, что дисперсия хи-квадрат случайной величины с  $m$  степенями свободы равна  $2m$ . Поэтому

$$\mathbf{D} \hat{\sigma}_e^2 = \mathbf{D} \frac{\sigma_e^2 \chi_{\nu_e}^2}{\nu_e} = \frac{2\sigma_e^4}{\nu_e} = \frac{2\sigma_e^4}{I(K-1)}$$

и, в силу независимости  $\overline{SS_A}$  и  $\overline{SS_e}$ ,

$$\mathbf{D} \hat{\sigma}_A^2 = \frac{1}{K^2}(\mathbf{D} \overline{SS_A} + \mathbf{D} \overline{SS_e}) = \frac{2}{I-1}(\sigma_A^2 + \frac{1}{K}\sigma_e^2)^2 + \frac{2}{I(K-1)}\frac{\sigma_e^4}{K^2}.$$

Таким образом, дисперсия оценки  $\hat{\sigma}_e^2$  стремится к нулю (а сама оценка будет состоятельной в среднеквадратическом) при  $I \rightarrow \infty$  или  $K \rightarrow \infty$ , даже если одна из величин  $K, I$  остается конечной. По другому себя ведет дисперсия  $\hat{\sigma}_A^2$ . Очевидно, оценка  $\hat{\sigma}_A^2$  будет состоятельной только в случае, когда число групп (образцов)  $I \rightarrow \infty$ .

## § 2. Полная классификация по двум признакам

Продолжая изучение производственного эксперимента, рассмотренного в начале предыдущего параграфа, предположим, что результат измерения образца зависит кроме всего прочего от степени настройки измерительного прибора. Можно предположить, что ошибка настройки прибора оказывает также случайное влияние на отклик. В этом



случае мы приходим к двухфакторной модели  $\mathcal{D}_A$  со случайными эффектами. В рамках нормальных допущений основные предположения этой модели можно записать в виде

$$\Omega : \begin{cases} y_{ik} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \varepsilon_{ijk}, & i = \overline{1, I}, j = \overline{1, J}, k = \overline{1, K}, \\ \{\alpha_i\}, \{\beta_j\}, \{\gamma_{ij}\}, \{\varepsilon_{ik}\} & \text{— независимы в совокупности,} \\ \alpha_i \sim \mathcal{N}(0, \sigma_A^2), \beta_j \sim \mathcal{N}(0, \sigma_B^2), \gamma_{ij} \sim \mathcal{N}(0, \sigma_{AB}^2), \\ \varepsilon_{ik} \sim \mathcal{N}(0, \sigma_e^2). \end{cases} \quad (\text{IV.3})$$

Поскольку вид основных предположений этой модели ничем не отличается от вида основных предположений в модели I  $\mathcal{D}_A$ , то все  $SS$  здесь вычисляются по тем же формулам, что и для модели I. Различия появляются при установлении их распределений. Рассмотрим, например,  $SS_A = JK \sum (y_{i..} - y_{...})^2$ . Подставляя в это соотношение модельные предположения IV.3 и вводя обозначение  $g_i = \alpha_i + \gamma_{i.} + \varepsilon_{i..}$ , получаем

$$SS_A = JK \sum_i (\alpha_i + \gamma_{i.} + \varepsilon_{i..} - \alpha_{.} - \gamma_{..} - \varepsilon_{...})^2 = JK \sum_i (g_i - g_{.})^2.$$

Очевидно, все  $g_i$  имеют нормальное распределение со средним нуль и дисперсией  $\mathbf{D} \alpha_i + \mathbf{D} \gamma_{i.} + \mathbf{D} \varepsilon_{i..} = \sigma_A^2 + \sigma_{AB}^2/J + \sigma_e^2/JK$ . Следовательно,

$$SS_A \sim (JK\sigma_A^2 + K\sigma_{AB}^2 + \sigma_e^2)\chi_{I-1}^2.$$

Аналогично получаются распределения всех остальных  $SS$ :

$$\begin{aligned} SS_B &\sim (IK\sigma_B^2 + K\sigma_{AB}^2 + \sigma_e^2)\chi_{J-1}^2, \\ SS_{AB} &\sim (K\sigma_{AB}^2 + \sigma_e^2)\chi_{(I-1)(J-1)}^2, \\ SS_e &\sim \sigma_e^2\chi_{IJ(K-1)}^2. \end{aligned} \quad (\text{IV.4})$$

Построим с помощью этих распределений критерии проверки стандартных гипотез.

Если верна гипотеза  $\mathbf{H}_A : \sigma_A^2 = 0$ , отношение  $\overline{SS}_A / \overline{SS}_e$  уже не будет иметь распределение Фишера, поскольку математические ожидания рассматриваемых  $SS$  не равны. Очевидно, желаемое распределение Фишера здесь будет иметь отношение  $\overline{SS}_A / \overline{SS}_{AB}$ . Составляя аналогичные отношения для критериев проверки гипотез  $\mathbf{H}_B : \sigma_B^2 = 0$  и  $\mathbf{H}_{AB} : \sigma_{AB}^2 = 0$ , получаем следующую таблицу  $\mathcal{D}_A$ .

Таблица двухфакторного  $\mathcal{D}_A$  модель II

Источник разброса	$SS$	Ст.свободы	$\overline{SS}$	$\mathcal{F}$	$\alpha$
Фактор А (строки)	$SS_A$	$\nu_A = I - 1$	$\overline{SS}_A$	$\frac{\overline{SS}_A}{\overline{SS}_{AB}}$	$1 - F_{\nu_A, \nu_e}(\mathcal{F})$
Фактор В (столбцы)	$SS_B$	$\nu_B = J - 1$	$\overline{SS}_B$	$\frac{\overline{SS}_B}{\overline{SS}_{AB}}$	$1 - F_{\nu_B, \nu_e}(\mathcal{F})$
Взаимодействия	$SS_{AB}$	$\nu_{AB} = (I - 1)(J - 1)$	$\overline{SS}_{AB}$	$\frac{\overline{SS}_{AB}}{\overline{SS}_e}$	$1 - F_{\nu_{AB}, \nu_e}(\mathcal{F})$
Ошибки	$SS_e$	$\nu_e = IJ(K - 1)$	$\overline{SS}_e$	–	–
Полная сумма квадратов	$SS_{\Pi}$	$KIJ - 1$	–	–	–

Если  $K = 1$ , то  $\mathcal{D}_A$  может быть проведен только при дополнительных предположениях относительно уравнения регрессии. Стандартное предположение об отсутствии взаимодействий между факторами приводит здесь к равенству нулю дисперсии  $\sigma_{AB}^2$ . В этом случае в таблице  $\mathcal{D}_A$  строка «ошибки» совпадает со строкой «взаимодействия» и  $SS_{AB} = SS_e$ .

Несмещенные оценки компонент дисперсии легко получаются из представленных выше распределений для  $SS$  :

$$\hat{\sigma}_A^2 = \frac{1}{IK}(\overline{SS}_A - \overline{SS}_{AB}), \quad \hat{\sigma}_B^2 = \frac{1}{JK}(\overline{SS}_B - \overline{SS}_{AB}),$$

$$\hat{\sigma}_{AB}^2 = \frac{1}{K}(\overline{SS}_{AB} - \overline{SS}_e).$$

Так же, как и в модели II однофакторного  $\mathcal{D}_A$ , состоятельность этих оценок имеет место, только если  $I \rightarrow \infty$  (для  $\hat{\sigma}_A^2$ ),  $J \rightarrow \infty$  (для  $\hat{\sigma}_B^2$ ) и  $J, I \rightarrow \infty$  (для  $\hat{\sigma}_{AB}^2$ ).

Заметим, что если  $K = 1$ , то компонента  $\sigma_{AB}^2$  отдельно не оценивается. В этом случае, однако,  $\mathcal{D}_A$  проводится обычно в предположении, что  $\sigma_{AB}^2 = 0$ .

### Полная классификация по трем признакам

Трехфакторный  $\mathcal{D}_A$  в модели II (так же, как четырех, пяти и выше) уже существенно отличается от двухфакторного  $\mathcal{D}_A$ . Это различие обусловлено тем обстоятельством, что с ростом числа факторов увеличивается количество слагаемых, входящих в среднее значение для  $\overline{SS}$  (IV.4). Например, для  $I$  групп фактора  $A$ ,  $J$  групп фактора  $B$ ,  $K$  групп фактора  $C$  и  $M$  наблюдений при каждом сочетании факторов

средние значения ошибок равны

$$\mathbf{E} \overline{SS}_A = \sigma_e^2 + JKM\sigma_A^2 + KM\sigma_{AB}^2 + JM\sigma_{AC}^2 + M\sigma_{ABC}^2 ,$$

$$\mathbf{E} \overline{SS}_{AC} = \sigma_e^2 + JM\sigma_{AC}^2 + M\sigma_{ABC}^2 .$$

Поэтому построение критерия проверки гипотезы  $\mathbf{H}_A : \sigma_A^2 = 0$  возможно только при дополнительных предположениях относительно параметров модели. Например, если положить  $\sigma_{AB}^2 = 0$ , то такой критерий может быть основан на отношении  $\overline{SS}_A / \overline{SS}_{AC}$ , которое будет иметь в этом случае распределение Фишера. (Если предположить, что  $\sigma_{AC}^2 = 0$ , тогда распределение Фишера будет иметь статистика  $\overline{SS}_A / \overline{SS}_{AB}$ .) Однако, если мы не можем предположить равенство нулю  $\sigma_{AB}^2$  или  $\sigma_{AC}^2$ , тогда построение точного критерия становится невозможным. В этом случае применяют приближенный F-критерий, формулировку которого мы здесь опускаем.

# Г л а в а V

## ФАКТОРНЫЕ ПЛАНЫ

Перейдем теперь к рассмотрению способов построения матриц плана  $\mathbf{X}$ , используемых при проведении факторных экспериментов. Разработка этих способов связана с необходимостью построения экономичных планов с достаточно простой процедурой вычислений.

### § 1. Дробные реплики полного факторного плана

Рассмотрим схему регрессионного эксперимента, в которой результаты измерений зависят от значений  $m (> 1)$  факторов  $x_1, \dots, x_m$ . Фактор  $x_i$  может принимать конечное число значений (- уровней)  $s_i (\geq 2)$ , причем выбор сочетаний уровней для различных факторов в каждом из проводимых экспериментов полностью зависит от исследователя. Здесь мы ограничимся рассмотрением случая, когда все  $s_i = 2$ . Такие *двухуровневые планы* имеют широкое распространение на практике в силу ряда обстоятельств. Во-первых, весьма типична ситуация, когда при проведении эксперимента рассматриваемые факторы могут либо присутствовать, либо отсутствовать, т.е. имеют только два значения. Во-вторых, если какой-либо фактор принимает значения на  $s$  уровнях, то можно вместо этого одного фактора рассмотреть  $s$  новых двухуровневых факторов, соответствующих уровням исходного фактора. В-третьих, как будет показано ниже, для некоторых типов регрессий наилучший план соответствует выбору из всех возможных уровней фактора только двух крайних.

В моделях с двухуровневыми факторами обычно входные переменные  $x_i$  принимают только два значения: 1, если  $i$ -ый фактор участвует в эксперименте, и  $-1$ , если фактор отсутствует. Функция регрессии

$$\begin{aligned} \eta(x_1, \dots, x_m) = & \beta_0 + \sum_{i=1}^m \beta_i x_i + \sum_{i<j} \beta_{ij} x_i x_j + \sum_{i<j<l} \beta_{ijl} x_i x_j x_l + \dots \\ & + \sum_{i_1 < i_2 < \dots < i_m} \beta_{i_1 \dots i_m} x_{i_1} \dots x_{i_m}. \end{aligned}$$

Здесь, как обычно, параметр  $\beta_0$  называется генеральным средним, параметры  $\beta_1 (i = \overline{1, m})$  главными эффектами факторов,  $\beta_{ij}$  – эффектами взаимодействий первого порядка или эффектами двухфакторных взаимодействий и, аналогично,  $\beta_{i_1 \dots i_k}$  – эффектами взаимодействий  $(k - 1)$ -ого порядка или эффектами  $k$ -факторных взаимодействий. Общее число неизвестных параметров, входящих в уравнение регрессии, равно, очевидно,  $p = 2^m$ . Поэтому для нахождения всех оценок необходимо провести по крайней мере  $n = 2^m$  наблюдений (если предполагается проведение дисперсионного анализа относительно параметров модели, то число наблюдений должно быть больше  $2^m$ ). Понятно, что уже при  $m > 7$  такое требование не часто осуществимо на практике. С целью сокращения числа неизвестных параметров выдвигаются некоторые предположения относительно этих параметров. Обычно начинают со старших взаимодействий, которые отбрасываются, т.е. считаются равными нулю. Чаще всего используется два вида функций регрессии:

$$\eta(x_1, \dots, x_m) = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m \quad (\text{V.1})$$

и

$$\eta(x_1, \dots, x_m) = \beta_0 + \sum_{i=1}^m \beta_i x_i + \sum_{i < j} \beta_{ij} x_i x_j.$$

**О п р е д е л е н и е 1.** Факторный план, в соответствии с которым проводится по одному наблюдению для каждой комбинации уровней факторов, называется *полным* и обозначается  $2^m$ .

Обозначение полного факторного плана связано с тем, что при этом плане необходимо провести ровно  $2^m$  наблюдений. Для модели регрессии вида (V.1), которая зависит всего от  $(m + 1)$ -ого параметра, такое число наблюдений избыточно. Однако также понятно, что большое число наблюдений позволяет оценить более богатую функцию регрессии. Эти соображения лежат в основе методики построения так называемых дробных реплик полного факторного плана.

Идею этого метода проиллюстрируем на примере трехфакторной регрессии. Сначала построим полный  $2^2$  план для первых двух факторов, а уровни третьего фактора будем выбирать, исходя из соотношения  $x_3 = x_1 x_2$  (варианты –  $x_3 = -x_1 x_2, x_3 = x_1$  и т.д.). В результате получим план, задаваемый следующей матрицей

$$\begin{array}{c|ccc} x_0 & x_1 & x_2 & x_3 \\ \hline 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{array} \quad (\text{V.2})$$

где первый столбец содержит значения аддитивного фактора (- коэффициент при  $\beta_0$ ), который присутствует во всех экспериментах. Заметим, кроме того, что второй и третий столбцы образованы путем перебора всех  $2^2$  возможных сочетаний цифр  $-1$  и  $1$ .

Проведя формальный дисперсионный анализ по данным измерений с матрицей плана (V.2), можно найти оценки всех параметров линейной модели (V.1). Проанализируем теперь план (V.2) применительно к полной модели со всеми взаимодействиями. Матрица плана для полной модели выглядит следующим образом:

$$\begin{array}{c|c|c|c|c|c|c|c}
 x_0 & x_1 & x_2 & x_3 & x_1x_2 & x_1x_3 & x_2x_3 & x_1x_2x_3 \\
 \hline
 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\
 1 & -1 & 1 & -1 & -1 & 1 & -1 & 1 \\
 1 & 1 & -1 & -1 & -1 & -1 & 1 & 1 \\
 1 & -1 & -1 & 1 & 1 & -1 & -1 & 1
 \end{array}$$

В этой матрице пятый столбец совпадает с четвертым, шестой – с третьим, седьмой – со вторым и восьмой – с первым. Таким образом, при анализе данных, основанных на наблюдениях с представленной матрицей плана, нельзя различить  $x_0$  и  $x_1x_2x_3$ ,  $x_1$  и  $x_2x_3$ ,  $x_2$  и  $x_1x_3$ ,  $x_3$  и  $x_1x_2$ . Следовательно, если нет дополнительных предположений о параметрах полной модели, то могут быть найдены оценки не всех восьми неизвестных коэффициентов, а только четырех смешанных

$$\beta_0 + \beta_{123}, \beta_1 + \beta_{23}, \beta_2 + \beta_{12}, \beta_3 + \beta_{12}.$$

Рассмотренный план называется *полуреplikой* или *планом*  $2^{3-1}$ . Для модели с общим числом факторов  $m$  дробные реплики строятся аналогичным образом. Сначала для выбранных (например, первых)  $m - k$  факторов строится полный  $2^{m-k}$  план, причем  $k$  выбирается так, чтобы  $2^{m-k} \geq m + 1$ . Затем с помощью некоторых *генерирующих соотношений* выбираются уровни оставшихся  $k$  факторов.

Анализ построенного плана применительно к полной модели проводится с помощью *определяющих соотношений*, получаемых путем перемножения генерирующих соотношений так, чтобы в левой части равенства находилась единица, а в правой – какие-либо произведения факторов. Так для дробной реплики  $2^{m-k}$  первые  $k$  определяющих соотношений получаются из генерирующих соотношений путем перемножения последних на их левые части и последующей замены величины  $x_i^2$  на  $1$ . Другие определяющие соотношения получаются перемножением ранее полученных и выделением из них новых. Знание всех определяющих соотношений позволяет найти всю систему совместных оценок без изучения матрицы плана. Для того чтобы определить, с

какими взаимодействиями смешано данное, нужно на него умножить обе части всех определяющих соотношений.

Рассмотрим пример с  $m = 6$  факторами. Построим дробную реплику  $2^{6-3}$  исходя из полного  $2^3$  факторного плана для первых трех факторов. Уровни четвертого, пятого и шестого факторов зададим с помощью генерирующих соотношений

$$x_4 = x_1x_2, \quad x_5 = x_1x_3, \quad x_6 = x_2x_3.$$

Первые три определяющих соотношения получаются умножением генерирующих соотношений на их левые части:

$$1 = x_1x_2x_4, \quad 1 = x_1x_3x_5, \quad 1 = x_2x_3x_6.$$

Попарное перемножение этих соотношений дает соотношения

$$1 = x_2x_3x_4x_5, \quad 1 = x_1x_3x_4x_6, \quad 1 = x_1x_2x_5x_6.$$

Перемножая по три равенства первые, получаем

$$1 = x_4x_5x_6.$$

Для того чтобы теперь, исходя из этих семи определяющих соотношений, определить, с какими эффектами взаимодействий смешан данный эффект, нужно умножить обе части всех определяющих соотношений на этот эффект. Например, умножая все определяющие соотношения на  $x_1$ , получаем

$$\begin{aligned} x_1 &= x_2x_4 = x_3x_5 = x_3x_4x_6 = x_2x_5x_6 = \\ &= x_1x_2x_3x_6 = x_1x_4x_5x_6 = x_1x_2x_3x_4x_5. \end{aligned}$$

Таким образом, главный эффект первого фактора  $\beta_1$  смешан с двухфакторными взаимодействиями  $\beta_{24}$ ,  $\beta_{35}$ , трехфакторными взаимодействиями  $\beta_{346}$ ,  $\beta_{256}$ , четырехфакторными взаимодействиями  $\beta_{1236}$ ,  $\beta_{1456}$ , и, наконец, с пятифакторным взаимодействием  $\beta_{12345}$ . Символически это записывается как

$$\beta_1^* \rightarrow \beta_1 + \beta_{24} + \beta_{35} + \beta_{346} + \beta_{256} + \beta_{1236} + \beta_{1456} + \beta_{12345} \quad (\text{V.3})$$

и означает, что оценка главного эффекта первого фактора, построенная по дробной реплике  $2^{6-3}$  полного плана для модели регрессии без взаимодействий, является оценкой суммы параметров из (V.3), если на самом деле верна модель регрессии с взаимодействиями.

Удобство использования дробных реплик обусловлено простотой формул для ОМНК. Не вдаваясь в подробности доказательства, заметим, что информационная матрица такого плана равна  $\mathcal{S} = n\mathbf{I}$  с

общим числом наблюдений  $n = 2^{m-k}$ . Оценки эффектов равны взвешенным средним

$$\beta_j = \frac{1}{n} \sum_{i=1}^n x_{ji} y_i,$$

где  $(y_1, \dots, y_n)'$  – вектор наблюдений,  $(x_{j1}, \dots, x_{jn})$  – строка матрицы плана  $\mathbf{X}$ , соответствующая выбранному фактору  $\beta_j$ .

## § 2. Латинские планы

**Определение 2.** Латинским квадратом порядка  $M$  называется такое расположение символов в виде квадратной таблицы  $M \cdot M$ , при котором каждый символ появляется по одному разу в каждой строке и в каждом столбце.

Существование латинских квадратов и один из способов их построения демонстрирует следующий пример

A	B	C	D
D	A	B	C
C	D	A	B
B	C	D	A

Латинские квадраты предназначены для проведения трехфакторного эксперимента, в котором каждый из факторов разбит на  $M$  уровней. Общее число измерений  $n = M^2$ , при этом сочетание уровней факторов в каждом эксперименте соответствует сочетанию номера строки (уровень первого фактора) с номером столбца (уровень второго фактора) и символом (уровень третьего фактора).

В сравнении с полным трехфакторным экспериментом латинские планы требуют проведения значительно меньшего числа наблюдений ( $M^2$  против  $M^3$ ) при сохранении удобных формул для вычисления оценок и сумм квадратов. Естественно, за такое сокращение приходится платить уменьшением количества параметров модели. Латинские планы используются лишь для дисперсионного анализа моделей без взаимодействий

$$y_{ijk} = \mu + \alpha_i^A + \alpha_j^B + \alpha_k^C + \varepsilon_{ijk}.$$

В этом случае оценки  $\mu$  и  $\alpha_i^A$  равны

$$\hat{\mu} = y_{\dots} = \frac{1}{M^2} \sum_{(i,j,k) \in \Lambda} y_{ijk},$$

$$\hat{\alpha}_i^A = y_{i..} - y_{\dots} = \frac{1}{M} \sum_{(j,k) \in \Lambda_i} y_{ijk} - y_{\dots},$$



где суммирование производится по множеству  $\Lambda$  – всех  $M^2$  упорядоченных троек  $(i, j, k)$ , и множеству  $\Lambda_i$  – всех пар  $(j, k)$ , для которых  $(i, j, k) \in \Lambda$ . Оценки главных эффектов двух других факторов получаются аналогично.

Выпишем теперь суммы квадратов ошибок. Остаточная сумма квадратов

$$SS_e = \sum_{(i,j,k) \in \Lambda} (y_{ijk} - \hat{\mu} - \hat{\alpha}_i^A - \hat{\alpha}_i^B - \hat{\alpha}_i^C)^2.$$

Сумма квадратов ошибок, обусловленная влиянием фактора  $A$ , равна

$$SS_A = M \sum_i y_{i..}^2 - M^2 y_{...}^2.$$

Как видим, формулы вполне стандартные, с единственным отличием, что суммирование при усреднениях производится по суженным множествам индексов.

**Замечание i.** Выбор того или иного латинского квадрата привносит элемент субъективности в схему проведения факторного эксперимента. Чтобы этого избежать, рекомендуют строить латинский план на основе квадрата, выбранного случайным образом из множества всех возможных латинских квадратов.

### Латинские прямоугольники

Если исследуемые факторы имеют разное число уровней, то при планировании факторного эксперимента применяются *латинские прямоугольники*, которые можно определить как подмножество строк (или столбцов) латинского квадрата. Например,

E	D	G	B	C	A	F
D	C	F	A	B	G	E
A	G	C	E	F	D	B
F	E	A	C	D	B	G

Приведенный прямоугольник можно использовать для проведения трехфакторного эксперимента с факторами, разбитыми на 7, 4 и 7 уровней. Кстати, этот  $(M \cdot k)$ -прямоугольник с  $M = 7, k = 4$ , обладает тем свойством, что каждая пара символов встречается в нем ровно в  $\lambda = k(k - 1)/(M - 1) = 2$  столбцах. Такие прямоугольники называются *квадратами Юдена*.

### Греко-латинские планы

Два латинских квадрата одинакового порядка называются ортогональными, если при наложении одного на другой каждая пара символов встречается только один раз. Например, два латинских квадрата

A	B	C	$\alpha$	$\beta$	$\gamma$
C	A	B	$\beta$	$\gamma$	$\alpha$
B	C	A	$\gamma$	$\alpha$	$\beta$

при наложении образуют квадрат

$A\alpha$	$B\beta$	$C\gamma$
$C\beta$	$A\gamma$	$B\alpha$
$B\gamma$	$C\alpha$	$A\beta$

в котором каждая пара латинских и греческих букв встречается по одному разу.

Квадраты, получаемые после наложения друг на друга двух ортогональных латинских квадратов, называются *греко-латинскими*. Известно, что ортогональные латинские квадраты существуют для всех  $M \neq 2, \neq 6$ .

Эти квадраты используются для построения планов более чем трехфакторного ДД. Так при проведении  $l$ -факторного ДД необходимо построить  $l - 2$  попарно ортогональных латинских квадрата. При этом номер строки в каждом квадрате будет указывать на уровень первого фактора, номер столбца на уровень второго фактора, символы квадратов, стоящие на пересечении выбранной строки и выбранного столбца, на уровни третьего, четвертого и т.д. факторов. Следует заметить, что количество попарно ортогональных латинских квадрата не может превышать  $M - 1$ .

### § 3. Блочные схемы

В этом разделе на примере двухфакторного ДД мы приведем общую схему построения планов, охватывающую уже рассмотренные дробные реплики и латинские планы.

Пусть 1-ый фактор имеет  $I$  уровней, 2-ой  $J$  уровней. Поставим в соответствие уровням 1-ого фактора элементы  $a_1, \dots, a_I$ . Тогда любой план проведения экспериментов может быть описан набором множеств

$$\{B_1, \dots, B_J\},$$

состоящих из элементов  $a_1, \dots, a_I$  и указывающих на сочетания уровней факторов в различных экспериментах. Например, если первый фактор имеет 3 уровня, а второй 2 уровня, то план, соответствующий схеме с  $B_1 = (a_1, a_2, a_3)$  и  $B_2 = (a_1, a_3)$ , подразумевает проведение пяти экспериментов со следующими сочетаниями уровней первого и второго факторов –  $(1; 1), (2; 1), (3; 1), (1; 2), (3; 2)$ , соответственно.

**О п р е д е л е н и е 3.** *Двумерной блок-схемой* называется структура, порожденная множеством элементов  $A = \{a_1, \dots, a_I\}$  и

множеством блоков  $\mathbf{B} = \{B_1, \dots, B_J\}$ .

Введем некоторые обозначения:

$k_j$  - число элементов в блоке  $B_j$ ;

$r_i$  - число блоков, содержащих элемент  $a_i$ ;

$\lambda_{lm}$  - число блоков, содержащих пару элементов  $(a_l, a_m)$ .

**О п р е д е л е н и е 4.** Блок-схема называется

*правильной*, если все блоки имеют одинаковый размер:  $k_j = k, \forall j$ ;

*полной*, если она правильная и  $k = I$ , то есть каждый блок содержит все элементы;

*равноповторной*, если каждый элемент встречается в равном числе блоков  $r_i = r$ ;

*симметричной*, если она правильная, равноповторная и  $I = J$  (то есть оба фактора имеют одинаковое число уровней);

*сбалансированной*, если она правильная, равноповторная и любая пара элементов принадлежит одному и тому же числу блоков:  $\lambda_{lm} = \lambda$ .

Основной интерес представляют неполные сбалансированные блок-схемы, так называемые ВІВ-планы. Эти планы обладают рядом оптимальных свойств и простотой проведения ДД. Дисперсии оценок главных эффектов факторов, построенных по наблюдениям, основанным на ВІВ-планах, равны между собой.

В следующем примере приводится симметричный ВІВ-план с параметрами  $I = J = 4, k = r = 3, \lambda = 2$ .

элементы	$B_1$	$B_2$	$B_3$	$B_4$
1	+		+	+
2	+	+		+
3	+	+	+	
4		+	+	+

**З а м е ч а н и е ii.** Не для всех сочетаний параметров существуют ВІВ-планы. Можно доказать, что параметры ВІВ-плана должны удовлетворять следующим соотношениям:

а)  $kJ = rI$ ;

б)  $r(k - 1) = \lambda(I - 1)$ ;

в)  $I \leq J$ ;

г)  $k \leq r$ ;

д) если  $I = J$  - четное число, то  $k - \lambda$  есть точный квадрат.

**З а м е ч а н и е iii.** Понятие двумерной блок-схемы обобщается на многомерный случай для построения планов многофакторного эксперимента. Примером трехмерной блок-схемы является латинский квадрат, а четырехмерной блок-схемы - греко-латинский квадрат.

## Г л а в а VI

### ОПТИМАЛЬНЫЕ ПЛАНЫ РЕГРЕССИОННЫХ ЭКСПЕРИМЕНТОВ

В этой главе мы рассмотрим задачу построения планов, оптимизирующих оценки параметров линейных регрессионных моделей. В предлагаемой ниже теории результат статистического эксперимента  $y$  зависит от некоторого  $k$ -мерного вектора входных переменных  $x \in \mathcal{R}^k$ :

$$y = \sum_{j=1}^p \beta_j f_j(x) + \varepsilon.$$

Множество значений  $\mathcal{X}$ , которые может принимать вектор  $x$ , называется *пространством планирования* и предполагается ограниченным и замкнутым. Функции  $f_1, \dots, f_p$  суть известные непрерывные функции на  $\mathcal{X}$ , называемые базисными функциями. Задача состоит в том, чтобы по результатам измерений  $y_1, \dots, y_n$  в  $n$  различных точках  $x_1, \dots, x_n$  оценить неизвестные параметры  $\beta_1, \dots, \beta_p$ .

Примеры. 1) Ранее рассматривавшееся стандартное уравнение регрессии  $y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \varepsilon$  для задач дисперсионного анализа многофакторного эксперимента также может быть записано в представленном виде, если положить вектор входных переменных  $x = (x_1, \dots, x_p)$ , а функции  $f_i(x) = x_i$ .

2) Часто возникает задача выявления зависимости отклика от некоторого «непрерывного» аргумента (например, продолжительность жизни микроорганизмов в зависимости от концентрации раствора в среде их обитания). Среди возможных форм функциональной зависимости чаще всего выбирают полиномиальную или тригонометрическую. Это обусловлено тем, что все практически полезные функции связи могут быть разложены в ряд Тейлора или в ряд Фурье. Отрезок такого ряда и выбирается в качестве приближенного описания реально существующего явления.

3) Рассмотрим задачу определения массы трех предметов на весах с двумя чашками. Если истинные веса этих предметов равны  $\beta_1, \beta_2, \beta_3$ , и отклонение стрелки весов в правую сторону интерпретируется как

вес со знаком плюс, а в левую сторону как вес со знаком минус, то показания весов можно представить в виде

$$y = \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \varepsilon,$$

где  $\varepsilon$  – ошибка измерения, обусловленная неучтенными факторами, а переменная  $x_i$  принимает значение  $-1$ , если при взвешивании  $i$ -ый предмет находился на левой чашке весов,  $+1$ , если на правой, и  $0$ , если предмет отсутствовал в данном эксперименте по взвешиванию.

## § 1. Планы эксперимента и их информационная матрица

Запишем уравнение регрессии в матричной форме

$$\eta(x) = \sum_{j=1}^p \beta_j f_j(x) = \vec{\beta}' \vec{f}(x),$$

где  $\vec{\beta} = (\beta_1, \dots, \beta_p)'$  и  $\vec{f}(x) = (f_1(x), \dots, f_p(x))'$ . Тогда вектор результатов статистического эксперимента представим в виде

$$\vec{y} = \mathbf{X}'\beta + \varepsilon,$$

где матрица плана

$$\mathbf{X} = \begin{pmatrix} f_1(x_1) & \dots & f_1(x_n) \\ \vdots & \dots & \vdots \\ f_p(x_1) & \dots & f_p(x_n) \end{pmatrix} = \left( \vec{f}(x_1), \dots, \vec{f}(x_n) \right).$$

Как было показано ранее, качество ОМНК зависит от матрицы  $\mathbf{X}$  :

$$\text{Cov}(\vec{\beta}^*) = \sigma^2(\mathcal{S}')^{-1} = \sigma^2(\mathbf{X}\mathbf{X}')^{-1}, \quad (\text{VI.1})$$

то есть, в конечном счете, от выбора точек  $x_1, \dots, x_n$ .

**О п р е д е л е н и е 1.** Совокупность точек  $x_1, \dots, x_n$  называется *планом эксперимента*.

Введем некоторые обозначения для планов эксперимента, задаваемых различными способами.

План эксперимента, в котором среди заданных точек  $x_1, \dots, x_n$  имеется только  $N$  различных  $x_1, \dots, x_N$ , называется *дискретным* и обозначается

$$\xi = \left\{ \begin{array}{l} x_1, \dots, x_N \\ p_1, \dots, p_N \end{array} \right\}, \quad (\text{VI.2})$$

где  $p_i = r_i/n$  и  $r_i$  есть число повторений  $i$ -ой точки. Очевидно,  $\sum p_i = 1$ . Для этого плана информационная матрица

$$\mathcal{S}(\xi) = \mathbf{X} \mathbf{X}' = \left( \vec{f}(x_1), \dots, \vec{f}(x_n) \right) \begin{pmatrix} \vec{f}'(x_1) \\ \vdots \\ \vec{f}'(x_n) \end{pmatrix} = \sum_{i=1}^n \vec{f}(x_i) \vec{f}'(x_i).$$

С учетом повторяющихся точек  $x_i$  получаем

$$\mathcal{S}(\xi) = \sum_{i=1}^N r_i \vec{f}(x_i) \vec{f}'(x_i) = n \sum_{i=1}^N p_i \vec{f}(x_i) \vec{f}'(x_i). \quad (\text{VI.3})$$

Поскольку информационная матрица выступает в роли обобщенной меры качества плана, то более удобно рассматривать вместо самой матрицы  $\mathcal{S}$  нормированную информационную матрицу  $\mathcal{S}/n$ . В этом случае сравнение планов автоматически производится при одинаковом числе наблюдений. В дальнейшем всегда под информационной матрицей будет подразумеваться нормированная информационная матрица.

Обобщением плана (VI.2) является план, в котором точки  $x_1, \dots, x_n$  выбираются случайным образом из совокупности точек  $x_1, \dots, x_N$  в соответствии с распределением  $\{p_1, \dots, p_N\}$ . Иными словами, для определения точки  $x$ , в которой будет производиться наблюдение  $y$ , проводится дополнительный случайный эксперимент, с вероятностью  $p_i$  принимающий значение  $x = x_i$ . Такой план в литературе называется «непрерывным» или «приближенным».

Обозначение непрерывного плана совпадает с (VI.2). Информационная матрица непрерывного плана вычисляется аналогично матрице детерминированного плана по формуле (VI.3). Если ввести случайную величину  $X$  с распределением, задаваемым таблицей (VI.2), то информационная матрица может быть записана через среднее значение  $X$ :

$$\mathcal{S}(\xi) = \mathbf{E} \vec{f}(X) \vec{f}'(X). \quad (\text{VI.4})$$

**З а м е ч а н и е** **i**. Подчеркнем существенное различие между дискретными и непрерывными планами. Для дискретного плана вида (VI.2) с  $p_i = r_i/n$  общее количество наблюдений ограничено числом  $n$  или кратным ему —  $2n, 3n, \dots$ . Однако по тому же самому плану, трактуемому как непрерывный, общее число наблюдений может принимать любое значение.

Дальнейшее обобщение способа задания плана эксперимента состоит в том, что совокупность точек  $x_1, \dots, x_n$  не фиксируется, а выбор точки  $x$  для проведения наблюдения производится случайно из всего

пространства планирования  $\mathcal{X}$  в соответствии с некоторым вероятностным распределением (вероятностной мерой)  $\xi$  на  $\mathcal{X}$ . Информационная матрица этого плана вычисляется аналогично (VI.4).

Как уже отмечалось, в качестве характеристики сравнения планов естественно выбрать матрицу, обратную информационной

$$\mathcal{D}(\xi) = (\mathcal{S}(\xi))^{-1}$$

(см. (VI.1)). Эту матрицу будем называть *дисперсионной*. Заметим сразу, что дисперсионная матрица определена только для невырожденных информационных матриц, то есть для планов полного ранга  $p$ . Следующее утверждение показывает, что такие планы обязаны иметь не менее  $p$  различных входных точек  $x_1, \dots, x_N$ .

**Лемма VI.1.** *Если для плана  $\xi$  вида (VI.2) число точек  $N < p$ , то этот план вырожденный.*

*Доказательство.* Так как ранг суммы матриц не больше суммы рангов, то в силу (VI.3) ранг информационной матрицы

$$\text{rang } \mathcal{S}(\xi) \leq \sum_{i=1}^N \text{rang } \vec{f}(x_i) \vec{f}'(x_i).$$

Столбцы матрицы  $\vec{f}(x) \vec{f}'(x)$  равны  $(f_l(x)f_1(x), \dots, f_l(x)f_p(x))'$ ,  $l = \overline{1, p}$ . Поэтому, если все функции  $f_l(x) = 0$ , то ранг этой матрицы равен 0. Если же  $f_l(x) \neq 0$  хотя бы для одного  $l$ , например,  $f_1(x) \neq 0$ , то все столбцы матрицы  $\vec{f}(x) \vec{f}'(x)$  будут пропорциональны первому. Следовательно, ранг матрицы  $\vec{f}(x) \vec{f}'(x)$  не больше 1, а ранг информационной матрицы  $\text{rang } \mathcal{S}(\xi) \leq N < p$ , то есть информационная матрица вырождена.  $\otimes$

## Примеры построения и сравнения планов

**П Р И М Е Р 1.** Линейная регрессия на отрезке:

$$y = \beta_0 + \beta_1 x + \varepsilon, \quad x \in [-1; 1].$$

Для задачи определения функциональной зависимости одной переменной от другой область изменения управляемой переменной обычно представляет собой ограниченный замкнутый интервал  $\mathcal{R}^1$ . Любой такой интервал может быть линейно преобразован в интервал  $[-1; 1]$ .

Для этой задачи вектор-функция  $\vec{f}(x) = (1, x)'$  и произведение

$$\vec{f}(x) \vec{f}'(x) = \begin{bmatrix} 1 & x \\ x & x^2 \end{bmatrix}.$$

Следовательно, информационная матрица любого плана равна

$$\mathcal{S}(\xi) = \begin{bmatrix} 1 & \mathbf{E}X \\ \mathbf{E}X & \mathbf{E}X^2 \end{bmatrix}.$$

Рассмотрим сначала план

$$\xi_1 = \left\{ \begin{array}{ccc} -1 & 0 & 1 \\ 1/3 & 1/3 & 1/3 \end{array} \right\}.$$

Для этого плана  $\mathbf{E}X = 0$ ,  $\mathbf{E}X^2 = 2/3$ , следовательно,

$$\mathcal{S}(\xi_1) = \begin{bmatrix} 1 & 0 \\ 0 & 2/3 \end{bmatrix} \quad \text{и} \quad \mathcal{D}(\xi_1) = \begin{bmatrix} 1 & 0 \\ 0 & 3/2 \end{bmatrix}.$$

Второй план сосредоточим в точках  $-1$  и  $1$ :

$$\xi_2 = \left\{ \begin{array}{cc} -1 & 1 \\ 1/2 & 1/2 \end{array} \right\}.$$

Для этого плана  $\mathbf{E}X = 0$ ,  $\mathbf{E}X^2 = 1$ , поэтому

$$\mathcal{S}(\xi_2) = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \text{и} \quad \mathcal{D}(\xi_2) = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \leq \mathcal{D}(\xi_1).$$

Таким образом, можно утверждать, что если по обоим планам проводится одинаковое число наблюдений (например, 6, 12 и т.д., если они дискретные, и произвольное число, если непрерывные), то план  $\xi_2$  предпочтительнее плана  $\xi_1$ .

**П Р И М Е Р 2.** Квадратическая регрессия на отрезке  $[-1; 1]$ .

Рассмотрим уравнение регрессии вида  $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \varepsilon$  и связанные с ним два плана

$$\xi_1 = \left\{ \begin{array}{ccc} -1 & 0 & 1 \\ 1/3 & 1/3 & 1/3 \end{array} \right\}, \quad \xi_2 = \left\{ \begin{array}{ccc} -1 & 0 & 1 \\ 1/4 & 2/4 & 1/4 \end{array} \right\}.$$

Для этой регрессии вектор-функция  $\vec{f}(x) = (1, x, x^2)'$  и информационная матрица любого плана равна

$$\mathcal{S}(\xi) = \mathbf{E} \vec{f}(X) \vec{f}'(X) = \begin{bmatrix} 1 & \mathbf{E}X & \mathbf{E}X^2 \\ \mathbf{E}X & \mathbf{E}X^2 & \mathbf{E}X^3 \\ \mathbf{E}X^2 & \mathbf{E}X^3 & \mathbf{E}X^4 \end{bmatrix}.$$

Для плана  $\xi_1$  средние  $\mathbf{E}X = 0$ ,  $\mathbf{E}X^2 = 2/3$ ,  $\mathbf{E}X^3 = 0$ ,  $\mathbf{E}X^4 = 2/3$ , следовательно, информационная и дисперсионная матрицы равны

$$\mathcal{S}(\xi_1) = \begin{bmatrix} 1 & 0 & 2/3 \\ 0 & 2/3 & 0 \\ 2/3 & 0 & 2/3 \end{bmatrix} \quad \text{и} \quad \mathcal{D}(\xi_1) = \begin{bmatrix} 3 & 0 & -3 \\ 0 & 3/2 & 0 \\ -3 & 0 & 9/2 \end{bmatrix}.$$



Аналогично, для плана  $\xi_2$   $\mathbf{E} X = 0$ ,  $\mathbf{E} X^2 = 1/2$ ,  $\mathbf{E} X^3 = 0$ ,  $\mathbf{E} X^4 = 1/2$ , и

$$\mathcal{S}(\xi_2) = \begin{bmatrix} 1 & 0 & 1/2 \\ 0 & 1/2 & 0 \\ 1/2 & 0 & 1/2 \end{bmatrix} \quad \text{и} \quad \mathcal{D}(\xi_2) = \begin{bmatrix} 2 & 0 & -2 \\ 0 & 2 & 0 \\ -2 & 0 & 4 \end{bmatrix}.$$

Разность матриц

$$\mathcal{D}(\xi_1) - \mathcal{D}(\xi_2) = \begin{bmatrix} 1 & 0 & -1 \\ 0 & -1/2 & 0 \\ -1 & 0 & 1/2 \end{bmatrix}$$

не будет, очевидно, ни положительной, ни отрицательной. Иными словами, планы  $\xi_1$  и  $\xi_2$  не сравнимы.

**П Р И М Е Р 3.** Взвешивание трех предметов.

Стандартная методика состоит в поочередном взвешивании каждого из предметов на одной из чашек весов. Мы сейчас покажем, что такой «естественный» план дает очень плохие оценки искомым весов.

Уравнение регрессии, описывающее данный эксперимент, задается вектор-функцией  $\vec{f}(x) = (x_1, x_2, x_3)'$  и, следовательно, произведение  $\vec{f} \vec{f}'$  представляет собой матрицу, состоящую из всевозможных произведений  $x_l x_m$ . Для «естественного» плана, сосредоточенного в трех точках  $(1, 0, 0)$ ,  $(0, 1, 0)$ ,  $(0, 0, 1)$ , отличны от нуля будут только диагональные элементы:  $x_1 x_1 = 1$  для точки  $(1, 0, 0)$ ,  $x_2 x_2 = 1$  для точки  $(0, 1, 0)$  и  $x_3 x_3 = 1$  для точки  $(0, 0, 1)$ . Поэтому информационная матрица этого плана

$$\mathcal{S}(\xi_1) = \frac{1}{3} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \frac{1}{3} \mathbf{I} \quad \text{и} \quad \mathcal{D}(\xi_1) = 3\mathbf{I}.$$

Рассмотрим теперь другой план  $\xi_2$ , сосредоточенный в восьми различных точках вида  $(\pm 1, \pm 1, \pm 1)$  с равными весами. Таким образом, по этому плану при общем числе наблюдений  $n = 8$  будет проведено по одному взвешиванию сразу всех предметов при различных их сочетаниях на чашках весов. Для этого плана средние значения  $\mathbf{E} X_l X_m = 0$  при различных  $l, m$ , так как произведение  $x_l x_m$  в половине случаев принимает значение 1 и в половине случаев значение  $-1$ . Очевидно, среднее значение  $\mathbf{E} X_l X_l = 1$ . Следовательно, для плана  $\xi_2$

$$\mathcal{S}(\xi_2) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \mathbf{I} \quad \text{и} \quad \mathcal{D}(\xi_2) = \mathbf{I},$$

то есть этот план в 3 раза эффективней «естественного» плана.

## § 2. Критерии оптимальности

Как видно из примера 2, не каждая пара планов будет иметь сравнимые дисперсионные матрицы. Следовательно, для определения понятия оптимального плана необходимо введение некоторой числовой характеристики, зависящей от дисперсионной матрицы. Общий подход здесь состоит в рассмотрении некоторого функционала  $\Phi$ , заданного на множестве возможных дисперсионных матриц, и определении оптимального плана через соотношение

$$\xi^* = \arg \min_{\xi} \Phi[\mathcal{D}(\xi)].$$

Различные функционалы  $\Phi$  определяют различные *критерии оптимальности*.

Выбор  $\Phi$  зависит от целей проводимого эксперимента. Так, если по результатам статистического эксперимента необходимо оценить все параметры регрессии, которые представляют самостоятельный интерес (например, вес предметов), то естественно стараться уменьшить возможный разброс всего вектора получаемых оценок. Числовой характеристикой разброса векторной случайной величины является обобщенная дисперсия – определитель ковариационной матрицы.

**О п р е д е л е н и е 2.** План

$$\xi^* = \arg \min_{\xi} \|\mathcal{D}(\xi)\|$$

называется *D-оптимальным* (от английского «determinant»).

**О п р е д е л е н и е 3.** План

$$\xi^* = \arg \min_{\xi} \|A' \mathcal{D}(\xi) A\|$$

называется *обобщенно D-оптимальным*.

Здесь матрица  $A$  размера  $p \cdot s$  ( $s \leq p$ ) и полного ранга  $s$  служит для установления «предпочтения» одних характеристик вектора неизвестных параметров над другими. Например, если эксперимент управляется  $p$  факторами, а необходимо оценить влияние только первых  $s$  из них, то следует выбрать матрицу  $A' = (\mathbf{I}, 0)$ , где  $\mathbf{I}$  – единичная матрица  $s \cdot s$ . В этом случае произведение  $A' \mathcal{D} A$  равно главному минору порядка  $s$  матрицы  $\mathcal{D}$  и в качестве оптимизируемой характеристики выступает разброс первых  $s$  компонент вектора  $\vec{\beta}^*$ .

Другой характеристикой разброса является сумма дисперсий оценок всех параметров регрессии – след дисперсионной матрицы, что приводит к следующему критерию оптимальности.

**О п р е д е л е н и е 4.** План

$$\xi^* = \arg \min_{\xi} \text{tr } \mathcal{D}(\xi)$$

называется *A-оптимальным* (от английского «average» – среднее).

Следующий критерий оптимальности связан с задачей прогноза возможного значения отклика  $y$  при различных значениях входной переменной  $x$  (например, прогноз ожидаемого дохода от вложенных средств). Наилучший прогноз равен, очевидно, значению функции регрессии при выбранном значении входной переменной. По теореме Гаусса-Маркова оценка с минимальной дисперсией для функции регрессии получается при подстановке в нее ОМНК параметров регрессии:  $\eta^*(x) = \vec{\beta}^{*'} \vec{f}(x)$ . Дисперсия этой оценки равна

$$\mathbf{D} \vec{\beta}^{*'} \vec{f}(x) = \vec{f}'(x) \text{Cov}(\vec{\beta}^*) \vec{f}(x) = \frac{\sigma^2}{n} \vec{f}'(x) \mathcal{D}(\xi) \vec{f}(x). \quad (\text{VI.5})$$

**О п р е д е л е н и е 5.** Функция  $d(x; \xi) = \vec{f}'(x) \mathcal{D}(\xi) \vec{f}(x)$  называется дисперсией оценки поверхности отклика.

Приведем одно важное свойство этой функции.

**Теорема VI.1.** Пусть  $X$  - случайная величина, определяемая планом  $\xi$  с невырожденной информационной матрицей. Тогда справедливы соотношения

$$\mathbf{E} d(X; \xi) = p, \quad (\text{VI.6})$$

$$\max_{x \in \mathcal{X}} d(x; \xi) \geq p, \quad (\text{VI.7})$$

причем, если во втором соотношении достигается знак равенства, то план  $\xi$  сосредоточен в точках максимума функции  $d(x; \xi)$ .

*Доказательство.* Заметим, что  $\vec{f}'(x) \mathcal{D}(\xi) \vec{f}(x)$  есть скалярная величина, поэтому она равняется своему следу:

$$\vec{f}'(x) \mathcal{D}(\xi) \vec{f}(x) = \text{tr}\{ \vec{f}'(x) \mathcal{D}(\xi) \vec{f}(x) \} = \text{tr}\{ \mathcal{D}(\xi) \vec{f}(x) \vec{f}'(x) \}$$

(по свойству следа, внутри которого можно переставлять любые пары матриц, если такая перестановка алгебраически допустима). Следовательно, справедлива цепочка равенств

$$\begin{aligned} \mathbf{E} d(x; \xi) &= \mathbf{E} \text{tr} \left\{ \mathcal{D}(\xi) \vec{f}(X) \vec{f}'(X) \right\} = \\ &= \text{tr} \left\{ \mathcal{D}(\xi) \mathbf{E} \vec{f}(X) \vec{f}'(X) \right\} = \text{tr}\{ \mathcal{D}(\xi) \mathcal{S}(\xi) \} = \text{tr}\{ \mathbf{I} \} = p. \end{aligned}$$

Второе утверждение теоремы вытекает из первого.  $\otimes$

Если исследователя интересует прогноз в конкретной точке  $x = x_0$ , то естественно постараться уменьшить дисперсию прогноза именно в этой точке. В связи с этим вводится следующий критерий оптимальности.

О п р е д е л е н и е 6. План

$$\xi^* = \arg \min_{\xi} d(x_0; \xi)$$

называется оптимальным для экстраполяции в точку  $x_0$ .

В том случае, когда планируется использование уравнения регрессии для прогноза во всех возможных точках  $x \in \mathcal{X}$ , предлагается несколько критериев оптимальности. Один из них основан на идее минимаксности.

О п р е д е л е н и е 7. План

$$\xi^* = \arg \min_{\xi} [\max_{x \in \mathcal{X}} d(x; \xi)]$$

называется G-оптимальным.

Другой критерий использует в качестве оптимизируемой характеристики взвешенное среднее значение дисперсии поверхности отклика.

О п р е д е л е н и е 8. План

$$\xi^* = \arg \min_{\xi} \int_{\mathcal{Z}} d(x; \xi) \omega(x) dx$$

называется Q-оптимальным. Здесь  $\omega(x)$  – неотрицательная функция на множестве  $\mathcal{Z}$ , задающая относительную важность точек  $z \in \mathcal{Z}$ . Множество  $\mathcal{Z}$  состоит из точек, для которых предполагается проведение прогноза отклика.

Последние два критерия, а также критерий A-оптимальности являются частными случаями критерия L-оптимальности.

О п р е д е л е н и е 9. План

$$\xi^* = \arg \min_{\xi} \text{tr}[L \mathcal{D}(\xi)]$$

называется L-оптимальным. Здесь  $L$  – фиксированная, неотрицательно определенная матрица  $p \cdot p$ .

Возможны, конечно, и другие определения критериев оптимальности, однако, как говорили классики, «нельзя объять необъятное».

### § 3. Теоремы эквивалентности

Не существует универсального способа построения оптимальных критериев, однако имеет место ряд утверждений, позволяющих проверить «подозрительные» планы на их возможную оптимальность. Эти утверждения называются теоремами эквивалентности.

Первая теорема принадлежит Киферу и Вольфовицу.

**Теорема VI.2.** *В классе всех непрерывных планов следующие утверждения эквивалентны:*

- а) план  $\xi^*$   $D$ -оптимален;
- б) план  $\xi^*$   $G$ -оптимален;
- в)  $\max_{x \in \mathcal{X}} d(x; \xi^*) = p$ .

*Информационные матрицы всех планов, удовлетворяющих одному из трех указанных утверждений, совпадают между собой. Все эти планы сосредоточены в точках  $x_i$ , для которых  $d(x_i; \xi^*) = p$ .*

*Доказательство.* Рассмотрим два плана  $\xi^*$  и  $\xi$  и положим план  $\xi_a = (1 - a)\xi^* + a\xi$  с некоторым  $a \in (0; 1)$ . Справедлива

**Лемма VI.2.** *Обозначим через  $X$  случайную величину, определяемую планом  $\xi$ . Тогда, если план  $\xi^*$  невырожден, то*

$$\left. \frac{\partial \ln \|\mathcal{S}(\xi_a)\|}{\partial a} \right|_{a=0+} = \mathbf{E} d(X; \xi^*) - p. \quad (\text{VI.8})$$

Сначала воспользуемся этой леммой для доказательства теоремы, а доказательство леммы оставим напоследок.

Пусть  $\xi^*$  —  $D$ -оптимальный план. Тогда  $\|\mathcal{S}(\xi_a)\| \leq \|\mathcal{S}(\xi^*)\|$ , то есть в точке  $a = 0$  достигается максимум функции  $\|\mathcal{S}(\xi_a)\|$ , и поэтому

$$\left. \frac{\partial \ln \|\mathcal{S}(\xi_a)\|}{\partial a} \right|_{a=0+} \leq 0.$$

Выберем в качестве  $\xi$  план вида  $\xi = \left\{ \begin{matrix} x \\ 1 \end{matrix} \right\}$ , сосредоточенный в произвольной точке  $x$  пространства планирования  $\mathcal{X}$ . Для этого плана  $\mathbf{E} d(X; \xi) = d(x; \xi)$  и, следовательно,

$$\left. \frac{\partial \ln \|\mathcal{S}(\xi_a)\|}{\partial a} \right|_{a=0+} = d(x; \xi^*) - p \leq 0.$$

Отсюда, с учетом неравенства (VI.7), получаем, что  $\max d(x; \xi^*) = p$  и план  $\xi^*$   $G$ -оптимален. Таким образом, доказано, что из  $D$ -оптимальности плана  $\xi^*$  следует его  $G$ -оптимальность и выполняется равенство в) теоремы. То, что план  $\xi^*$  сосредоточен в точках максимума функции дисперсии отклика, следует из теоремы VI.1.

Предположим теперь, что план  $\xi^*$   $G$ -оптимален, но не является  $D$ -оптимальным. Так как, в силу только что доказанного свойства, для  $D$ -оптимального плана выполняется свойство в), то для  $G$ -оптимального плана  $\xi^*$  также имеет место равенство

$$\max_{x \in \mathcal{X}} d(x; \xi^*) = p. \quad (\text{VI.9})$$

Выберем в определении плана  $\xi_a$  в качестве плана  $\xi$  какой-либо  $D$ -оптимальный план. Позже мы докажем, что функционал  $-\ln \|\mathcal{S}(\xi)\|$  строго выпукл, то есть справедливо неравенство

$$\ln \|\mathcal{S}(\xi_a)\| \geq (1 - a) \ln \|\mathcal{S}(\xi^*)\| + a \ln \|\mathcal{S}(\xi)\|.$$

Отсюда следует, что дифференциальное отношение

$$\frac{\ln \|\mathcal{S}(\xi_a)\| - \ln \|\mathcal{S}(\xi^*)\|}{a} \geq \ln \|\mathcal{S}(\xi)\| - \ln \|\mathcal{S}(\xi^*)\| > 0.$$

Устремляя  $a$  к нулю, получаем (с учетом равенства VI.8 леммы VI.2) строгое неравенство  $\mathbf{E} d(X; \xi^*) > p$ , которое противоречит VI.9. Следовательно,  $G$ -оптимальный план будет и  $D$ -оптимальным.

Покажем, что информационные матрицы всех  $D$ -оптимальных планов совпадают. Предположим от противного, что найдутся два  $D$ -оптимальных плана  $\xi_1$  и  $\xi_2$  с разными информационными матрицами. Образует из них новый план посредством выпуклой комбинации:  $\xi = (1 - a)\xi_1 + a\xi_2$ . Для этого плана в силу строгой выпуклости функционала  $-\ln \|\mathcal{S}(\xi)\|$  логарифм определителя дисперсионной матрицы

$$\begin{aligned} \ln \|\mathcal{D}(\xi)\| &= -\ln \|\mathcal{S}(\xi)\| < \\ &< -(1 - a) \ln \|\mathcal{S}(\xi_1)\| - a \ln \|\mathcal{S}(\xi_2)\| = \\ &= (1 - a) \ln \|\mathcal{D}(\xi_1)\| + a \ln \|\mathcal{D}(\xi_2)\| = \min_{\xi} \ln \|\mathcal{D}(\xi)\|. \end{aligned}$$

То есть определитель дисперсионной матрицы плана  $\xi$  меньше наименьшего возможного значения, что, естественно, невозможно.

Для доказательства теоремы осталось показать строгую выпуклость выше рассмотренного функционала и доказать лемму VI.2. Известно, что для плотности многомерного нормального распределения

$\mathcal{N}_p(\vec{0}, A^{-1})$  с ковариационной матрицей  $A^{-1}$

$$\frac{\|A\|^{1/2}}{\pi^{p/2}} \int_{\mathcal{R}^p} \exp\{-\vec{x}' A \vec{x}\} d\vec{x} = 1.$$

Следовательно, для любых двух различных положительно определенных матриц  $A$  и  $B$  имеем

$$\begin{aligned} \frac{\pi^{p/2}}{\|(1-a)A + aB\|^{1/2}} &= \int_{\mathcal{R}^p} \exp\{-\vec{x}'((1-a)A + aB)\vec{x}\} d\vec{x} = \\ &= \int_{\mathcal{R}^p} (\exp\{-\vec{x}' A \vec{x}\})^{1-a} (\exp\{-\vec{x}' B \vec{x}\})^a d\vec{x} < \\ &< \left( \int_{\mathcal{R}^p} \exp\{-\vec{x}' A \vec{x}\} d\vec{x} \right)^{1-a} \times \left( \int_{\mathcal{R}^p} \exp\{-\vec{x}' B \vec{x}\} d\vec{x} \right)^a = \\ &= \frac{\pi^{p(1-a)/2} \pi^{pa/2}}{(\|A\|)^{(1-a)/2} (\|B\|)^{a/2}} \end{aligned}$$

в силу известного неравенства Гёльдера. Отсюда, произведя сокращения и прологарифмировав обе части неравенства, получаем требуемую выпуклость.

Докажем теперь лемму VI.2. Обозначим через  $s_{ij}(\xi)$  элементы информационной матрицы  $\mathcal{S}(\xi)$ , а через  $s^{ij}(\xi)$  элементы матрицы, обратной к  $\mathcal{S}(\xi)$ . Тогда, как известно,

$$\begin{aligned} \|\mathcal{S}(\xi)\| &= \sum_{j=1}^p (-1)^{i+j} \|\mathcal{S}_{ij}(\xi)\| s_{ij}(\xi), \\ s^{ij}(\xi) &= (-1)^{i+j} \frac{\|\mathcal{S}_{ij}(\xi)\|}{\|\mathcal{S}(\xi)\|}, \quad i, j = 1, \dots, p, \end{aligned}$$

где  $\mathcal{S}_{ij}(\xi)$  – определитель матрицы, полученной из матрицы  $\mathcal{S}(\xi)$  вычеркиванием  $i$ -ой строки и  $j$ -ого столбца. Итак, производная

$$\frac{\partial \|\mathcal{S}(\xi)\|}{\partial s_{ij}} = (-1)^{i+j} \|\mathcal{S}_{ij}(\xi)\|$$

и

$$\frac{1}{\|\mathcal{S}(\xi)\|} \frac{\partial \|\mathcal{S}(\xi)\|}{\partial s_{ij}} = s^{ij}(\xi).$$

Кроме того, так как  $s_{ij}(\xi_a) = (1-a)s_{ij}(\xi^*) + as_{ij}(\xi)$ , то

$$\frac{\partial s_{ij}(\xi_a)}{\partial a} = s_{ij}(\xi) - s_{ij}(\xi^*).$$

Воспользовавшись теперь теоремой о дифференцировании сложной функции, получаем

$$\begin{aligned}
\left. \frac{\partial \ln \|\mathcal{S}(\xi_a)\|}{\partial a} \right|_{a=0} &= \frac{1}{\|\mathcal{S}(\xi_a)\|} \frac{\partial \|\mathcal{S}(\xi_a)\|}{\partial a} = \\
&= \frac{1}{\|\mathcal{S}(\xi_a)\|} \sum_{ij} \frac{\partial \|\mathcal{S}(\xi_a)\|}{\partial s_{ij}} \frac{\partial s_{ij}(\xi_a)}{\partial a} = \sum_{ij} s^{ij}(\xi^*) [s_{ij}(\xi) - s_{ij}(\xi^*)] = \\
&= \text{tr } \mathcal{S}^{-1}(\xi^*) [\mathcal{S}(\xi) - \mathcal{S}(\xi^*)] = \text{tr } \mathcal{S}^{-1}(\xi^*) \mathcal{S}(\xi) - \text{tr } \mathbf{I}_p = \\
&= \text{tr } \mathcal{S}^{-1}(\xi^*) \mathbf{E} \vec{f}(X) \vec{f}'(X) - p = \mathbf{E} \vec{f}'(X) \mathcal{S}^{-1}(\xi^*) \vec{f}(X) - p = \\
&= \mathbf{E} d(X; \xi^*) - p.
\end{aligned}$$

Лемма доказана.  $\otimes$

### Примеры применения теоремы Кифера-Вольфовица

**П Р И М Е Р 1.** Для случая линейной регрессии на отрезке  $[-1; 1]$   $D$ -оптимальным, а соответственно и  $G$ -оптимальным является план

$$\xi^* = \left\{ \begin{array}{cc} -1 & 1 \\ 1/2 & 1/2 \end{array} \right\}.$$

Действительно, как мы уже видели, для этого плана дисперсионная матрица равна единичной  $\mathcal{D}(\xi^*) = \mathbf{I}$ . Поэтому дисперсия поверхности отклика равна  $d(x; \xi^*) = \vec{f}'(x) \mathbf{I} \vec{f}(x) = (1, x)(1, x)' = 1 + x^2$ .

Максимум этой функции равен 2 (= числу неизвестных параметров) и достигается в точках  $-1, 1$ . Следовательно, по теореме Кифера-Вольфовица план  $\xi^*$   $D$ - и  $G$ -оптимален.

По плану  $\xi^*$  можно проводить любое число наблюдений, только если он интерпретируется как непрерывный. Для дискретного плана  $\xi^*$  число наблюдений может быть только четным. Покажем, что для дискретных планов с нечетным числом наблюдений, кратным 3, критерии  $D$ -оптимальности и  $G$ -оптимальности не совпадают. Рассмотрим план общего вида

$$\xi = \left\{ \begin{array}{ccc} x_1 & x_2 & x_3 \\ 1/3 & 1/3 & 1/3 \end{array} \right\},$$

где по крайней мере две точки  $x_1, x_2, x_3$  различны (это необходимо, чтобы план был невырожденным). Информационная матрица этого плана равна

$$\mathcal{S}(\xi) = \begin{bmatrix} 1 & \mathbf{E} x \\ \mathbf{E} x & \mathbf{E} x^2 \end{bmatrix} = \frac{1}{3} \begin{bmatrix} 3 & x_1 + x_2 + x_3 \\ x_1 + x_2 + x_3 & x_1^2 + x_2^2 + x_3^2 \end{bmatrix},$$



а ее определитель  $|\mathcal{S}(\xi)| = \frac{2}{9}(x_1^2 + x_2^2 + x_3^2 - x_1 - x_2 - x_3)$ .

Если здесь существует  $D$ -оптимальный план, тогда на этом плане достигается максимум определителя информационной матрицы, то есть все частные производные  $\mathcal{S}(\xi)$  по переменным  $x_i$  равны нулю, либо этот максимум достигается на крайних точках пространства планирования:  $x_i = \pm 1$ . Система уравнений для частных производных имеет вид

$$\begin{aligned} 2x_1 - x_2 - x_3 &= 0, \\ -x_1 + 2x_2 - x_3 &= 0, \\ -x_1 - x_2 + 2x_3 &= 0. \end{aligned}$$

Эта система вырождена и имеет бесконечно много решений  $x_1 = x_2 = x_3 = c$ , где  $c$  – произвольная константа. Для любого из этих решений план  $\xi$  вырожден, так как сосредоточен всего в одной точке при двух неизвестных параметрах. На крайних точках  $\pm 1$  существует всего два невырожденных плана

$$\xi_1 = \left\{ \begin{array}{cc} -1 & 1 \\ 1/3 & 2/3 \end{array} \right\} \quad \text{и} \quad \xi_2 = \left\{ \begin{array}{cc} -1 & 1 \\ 2/3 & 1/3 \end{array} \right\},$$

которые и будут  $D$ -оптимальными. Информационная и дисперсионная матрицы для плана  $\xi_1$  (также как и для плана  $\xi_2$ ) равны

$$\mathcal{S}(\xi_1) = \begin{bmatrix} 1 & 1/3 \\ 1/3 & 1 \end{bmatrix} \quad \text{и} \quad \mathcal{D}(\xi_1) = \frac{3}{8} \begin{bmatrix} 3 & -1 \\ -1 & 3 \end{bmatrix}.$$

Легко видеть, что дисперсия поверхности отклика  $d(x; \xi_1) = \frac{3}{8}(3 - 2x + 3x^2) = 1 + \frac{1}{8}(3x - 1)^2$ . Максимум ее достигается при  $x = -1$  и равен  $\max d(x; \xi_1) = 3$ , то есть планы  $\xi_1$  и  $\xi_2$   $D$ -оптимальны.

Рассмотрим дискретный план  $\xi_3$ , сосредоточенный с равными весами в трех точках  $-1, 0, 1$ . Для этого плана дисперсия отклика  $d(x; \xi_3) = 1 + 3x^2/2$  (докажите!), а ее максимум равен, очевидно,  $5/2 < 3$ . То есть план  $\xi_3$  по критерию  $G$ -оптимальности лучше  $D$ -оптимальных планов  $\xi_1$  и  $\xi_2$ .

**П Р И М Е Р 3.** Взвешивание трех предметов.

Для второго из рассматриваемых в этом примере планов дисперсионная матрица равна единичной  $\mathcal{D}(\xi_2) = \mathbf{I}$ . Следовательно, дисперсия поверхности отклика равна  $d(x; \xi_2) = (x_1, x_2, x_3)\mathbf{I}(x_1, x_2, x_3)' = x_1^2 + x_2^2 + x_3^2$ . Максимум дисперсии отклика равен, очевидно, 3 (– числу неизвестных параметров) и достигается на точках вида  $(\pm 1, \pm 1, \pm 1)$ . По теореме Кифера-Вольфовица план  $\xi_2$   $D$ -оптимален.

## Теоремы эквивалентности для $L$ -оптимальных планов

Имеет место следующая

**Теорема VI.3.** *Если существует  $L$ -оптимальный план, тогда следующие утверждения эквивалентны:*

- а) план  $\xi^*$   $L$ -оптимален;
- б) максимум функции

$$\varphi(x; \xi^*) = \vec{f}'(x) \mathcal{D}(\xi) L \mathcal{D}(\xi) \vec{f}(x) \quad (\text{VI.10})$$

равен  $\max_{x \in \mathcal{X}} \varphi(x; \xi^*) = \text{tr}[L \mathcal{D}(\xi^*)]$ .

План  $\xi^*$  сосредоточен в точках максимума функции  $\varphi(x; \xi^*)$ .

**Замечание ii.** Легко видеть, что для критерия  $A$ -оптимальности матрица  $L = \mathbf{I}$ , а функция (VI.10) равна

$$\varphi(x; \xi) = \vec{f}'(x) \mathcal{D}^2(\xi) \vec{f}(x).$$

Для критерия  $Q$ -оптимальности  $L = \int_Z \vec{f}(x) \vec{f}'(x) \omega(x) dx$  и функция

$$\varphi(x; \xi) = \int_Z (\vec{f}'(x) \mathcal{D}(\xi) \vec{f}(z))^2 \omega(z) dz.$$

Для критерия оптимальности при экстраполяции в точку функция

$$\varphi(x; \xi) = (\vec{f}'(x) \mathcal{D}(\xi) \vec{f}(x_0))^2.$$

В качестве примера покажем, что для квадратической регрессии  $A$ -оптимальным является план

$$\xi_2 = \left\{ \begin{array}{ccc} -1 & 0 & 1 \\ 1/4 & 2/4 & 1/4 \end{array} \right\}.$$

Как мы уже видели, дисперсионная матрица этого плана равна

$$\mathcal{D}(\xi_2) = \begin{bmatrix} 2 & 0 & -2 \\ 0 & 2 & 0 \\ -2 & 0 & 4 \end{bmatrix} \quad \text{и} \quad \mathcal{D}^2(\xi_2) = \begin{bmatrix} 8 & 0 & -12 \\ 0 & 4 & 0 \\ -12 & 0 & 20 \end{bmatrix}.$$

Следовательно, функция (VI.10) равна

$$\varphi(x; \xi) = (1, x, x^2) \mathcal{D}^2 \begin{pmatrix} 1 \\ x \\ x^2 \end{pmatrix} = 8 - 20x^2(1 - x^2).$$

Очевидно, что  $\max \varphi(x; \xi_2) = 8 = \text{tr} \mathcal{D}(\xi_2)$  и достигается в точках  $-1, 0, 1$ , откуда следует  $A$ -оптимальность плана  $\xi_2$ .

## § 4. $D$ -оптимальные планы для полиномиальной и тригонометрической регрессии

### Полиномиальная регрессия на отрезке $[-1; 1]$

В следующей теореме дается способ построения  $D$ -оптимальных планов для полиномиальной регрессии вида  $\eta(x) = \beta_0 + \beta_1 x + \dots + \beta_m x^m$ ,  $x \in [-1; 1]$ . Предварительно введем одно вспомогательное

**Определение 10.** Полиномом Лежандра степени  $n$  называется полином  $P_n(x)$ , определяемый соотношением

$$P_n(x) = \frac{1}{n!2^n} \frac{d^n}{dx^n} (x^2 - 1)^n.$$

Для первых трех степеней полином Лежандра равен

$$P_0(x) = 1, \quad P_1(x) = x, \quad P_2(x) = \frac{1}{2}(3x^2 - 1), \quad P_3(x) = \frac{1}{2}(5x^3 - 3x).$$

**Теорема VI.4.** Для полиномиальной регрессии степени  $m$  на отрезке  $[-1; 1]$  существует единственный непрерывный  $D$ -оптимальный план. Этот план сосредоточен с равными весами  $1/(m+1)$  в  $(m+1)$ -ой точке, которые являются решениями уравнения

$$(1 - x^2) \frac{d}{dx} P_m(x) = 0,$$

где  $P_m(x)$  – полином Лежандра степени  $m$ .

**П Р И М Е Р 2.** Для квадратичной регрессии степень  $m = 2$ , производная  $P_2'(x) = 3x$ , а уравнение  $(1 - x^2)3x = 0$  имеет три решения  $x = -1, x = 0, x = 1$ . Итак, для квадратичной регрессии непрерывный  $D$ -оптимальный план сосредоточен с весами  $1/3$  в точках  $-1, 0, 1$ .

### Тригонометрическая регрессия на отрезке $[0; 2\pi]$

Другой важный частный случай уравнения регрессии возникает при аппроксимации отклика отрезком ряда Фурье, т.е.

$$\eta(x; \beta) = \beta_1 + \sum_{j=1}^k [\beta_{2j} \cos jx + \beta_{2j+1} \sin jx] \quad (\text{VI.11})$$

с  $p = 2k + 1$  неизвестными параметрами и пространством планирования  $\mathcal{X} = [0, 2\pi]$ . В следующей теореме приводится вид  $D$ -оптимальных планов для тригонометрической регрессии; в ней, в частности, фигурирует непрерывный план  $\xi^* = \frac{1}{2\pi} dx$ , выбирающий точки в соответствии с равномерным распределением на отрезке  $[0, 2\pi]$ .

**Теорема VI.5.** *Если уравнение регрессии имеет вид (VI.11), тогда равномерный непрерывный план  $\xi^* = (1/2\pi)dx$  является  $D$ -оптимальным непрерывным планом. План  $\xi_N$ , сосредоточенный в  $N \geq 2k + 1$  точках*

$$x_i = \frac{i-1}{N} 2\pi, \quad i = \overline{1, N},$$

*также будет  $D$ -оптимальным непрерывным планом.*

## § 5. Оптимальные планы первого порядка

Предположим теперь, что отклик зависит от  $m$  переменных (факторов) и уравнение регрессии имеет вид

$$\eta(\vec{x}) = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m \quad (\text{VI.12})$$

с входной точкой  $\vec{x} = (x_1, \dots, x_m)'$ . Совокупность входных точек образует стандартную матрицу плана  $\mathbf{X}'$ , строки которой равны  $(1, x_{1i}, \dots, x_{mi}), i = \overline{1, n}$ . Важную роль для задач с подобной регрессией играют ортогональные планы, для которых информационная, а, следовательно, и дисперсионная матрицы диагональны и невырождены. Для ортогональных планов выполнены соотношения

$$\begin{aligned} \sum_{i=1}^n x_{ji} x_{li} &= 0, \quad j, l = 1, \dots, m, \quad l \neq j, \\ \sum_{i=1}^n x_{ji} &= 0, \quad i = 1, \dots, m. \end{aligned}$$

Пространство планирования, как обычно, предполагается ограниченным замкнутым подмножеством пространства  $\mathcal{R}^m$ . Здесь удобно будет записать это пространство в виде шаров

$$\sum_{i=1}^n x_{ji}^2 \leq C_j^2, \quad j = 1, \dots, m, \quad (\text{VI.13})$$

где  $C_j$  заданы.

**Теорема VI.6.** *Пусть (VI.12) — модель регрессии с дисперсией ошибки  $\sigma^2$ . Тогда для дискретных планов с пространством планирования вида (VI.13) дисперсии ОМНК всех неизвестных параметров*

$$\mathbf{D}(\beta_j^*) \geq \frac{\sigma^2}{C_j^2}, \quad j = 0, \dots, m, \quad (\text{VI.14})$$

с  $C_0^2 = n$ . Причем, если план ортогональный и в (VI.13) выполняется равенство, тогда и в (VI.14) достигается знак равенства, то есть этот план является  $A$ -оптимальным в классе дискретных планов с пространством планирования (VI.13).

**П Р И М Е Р 3.** Для задачи взвешивания трех предметов  $x_j$  принимает три значения  $-1, 0, 1$ , поэтому точки планирования удовлетворяют соотношениям (VI.13) с  $C_j^2 = n$ . Следовательно,  $A$ -оптимальным будет любой ортогональный план, при котором в каждом взвешивании участвуют все три предмета. В частности, таковым будет план, сосредоточенный в 8-ми точках вида  $(\pm 1, \pm 1, \pm 1)$ . Также  $A$ -оптимальным будет план, построенный как дробная реплика  $2^{3-1}$  полного факторного плана и требующий 4 взвешивания с расположением предметов  $(-1, -1, 1), (1, -1, -1), (-1, 1, -1), (1, 1, 1)$ .

Построение  $D$ -оптимальных планов для регрессии (VI.12) возможно, если предположить, что пространство планирования совпадает с  $m$ -мерным кубом  $\mathcal{X} = [-1; 1]^m$ , т.е. все  $x_j$  принимают значения из интервала  $[-1; 1]$ . Имеет место следующая

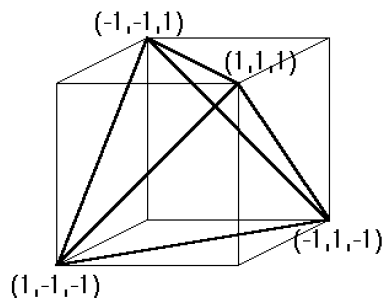
**Теорема VI.7.** Пусть уравнение регрессии задается (VI.12) с пространством планирования  $\mathcal{X} = [-1; 1]^m$ . Тогда

1) полный факторный план, сосредоточенный во всех  $2^m$  вершинах куба  $\mathcal{X}$ ,  $D$ -оптимален;

2)  $D$ -оптимальным является также план, точки сосредоточения которого лежат в вершинах куба  $\mathcal{X}$  и образуют правильный  $m$ -мерный  $(m + 1)$ -угольник (симплекс-план).

### З а м е ч а н и е iii.

На плоскости при размерности  $m = 2$  симплекс-планы не существуют. При  $m = 3$  любой правильный тетраэдр, например,  $(1, 1, 1), (1, -1, -1), (-1, 1, -1)$  и  $(-1, -1, 1)$  образует 3-мерный симплекс.



Симплекс-планы, если они существуют, также  $A$ -оптимальны.

# Литература

- [1] Шеффе Г. Дисперсионный анализ. – М.: Наука, 1980. – 512 с.
- [2] Ермаков С.М., Жиглявский А.А. Математическая теория оптимального эксперимента. – М.: Наука, 1987. – 320 с.
- [3] Тюрин Ю.Н., Макаров А.А. Анализ данных на компьютере. – М.:ИНФРА-М, Финансы и статистика, 1995. – 384 с.
- [4] Себер Дж. Линейный регрессионный анализ. – М.: Мир, 1980. – 456 с.