

Дисперсионный анализ

Н.П. Алексеева,

СПбГУ, математико-механический факультет

2020 г.

Линейная модель с фиксированными эффектами

Предполагаем, что у нас есть r нормально распределенных популяций $\mathcal{N}(\mu_1, \sigma), \dots, \mathcal{N}(\mu_r, \sigma)$, где через μ_i обозначены генеральные средние внутри каждой популяции.

Переменная x_{ij} означает j -е наблюдение в i -й подпопуляции. Это предположение можно записать в виде:

$$x_{ij} = \mu_i + e_{ij}, \quad i = 1, \dots, r, \quad j = 1, \dots, n_i, \quad (1)$$

где "ошибки" $e_{ij} \sim \mathcal{N}(0, \sigma)$. Перепишем модель однофакторного дисперсионного анализа в виде

$$x_{ij} = \mu + \alpha_i + e_{ij}, \quad i = 1, \dots, r, \quad j = 1, \dots, n_i, \quad (2)$$

где μ генеральное среднее объединенной популяции, α_i главные эффекты подпопуляций.

Модель с фиксированными эффектами

Наилучшими оценками параметров модели

$$x_{ij} = \mu + \alpha_i + e_{ij}, \quad i = 1, \dots, r, \quad j = 1, \dots, n_i,$$

в смысле метода наименьших квадратов являются

$$\hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^r \sum_{j=1}^{n_i} x_{ij}, \quad \hat{\alpha}_i = \bar{x}_{i\cdot} - \bar{x} = \frac{1}{n_i} \sum_{j=1}^{n_i} x_{ij} - \bar{x}.$$

Ошибкой в этой модели является выражение

$$Q_2 = \sum_{i=1}^r \sum_{j=1}^{n_i} (x_{ij} - \hat{\mu} - \hat{\alpha}_i)^2 = \sum_{i=1}^r \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_{i\cdot})^2.$$

Ее число степеней свободы равно $\nu_2 = n - r$. Принятие гипотезы $H_0 : \alpha_1 = \dots = \alpha_r = 0$ означает справедливость эквивалентной гипотезы $H_0 : \mu_1 = \dots = \mu_r = \mu$.

Проверка нулевой гипотезы $H_0 : \alpha_1 = \dots = \alpha_r = 0$

При $H_0 : \alpha_1 = \dots = \alpha_r = 0$ имеет место усеченная модель вида

$$x_{ij} = \mu + e_{ij}, \quad i = 1, \dots, r; \quad j = 1, \dots, n_i. \quad (3)$$

Наилучшей оценкой параметра μ в смысле метода наименьших квадратов является $\hat{\mu} = \bar{x}$. Ошибкой \widetilde{Q}_2 с числом степеней свободы $\widetilde{\nu}_2 = n - 1$ является выражение

$$\widetilde{Q}_2 = \sum_{i=1}^r \sum_{j=1}^{n_i} (x_{ij} - \hat{\mu})^2 = \sum_{i=1}^r \sum_{j=1}^{n_i} (x_{ij} - \bar{x})^2 = Q.$$

Статистика критерия правдоподобия имеет вид:

$$F = \frac{(\widetilde{Q}_2 - Q_2)/(\widetilde{\nu}_2 - \nu_2)}{Q_2/\nu_2} = \frac{Q_1/(r-1)}{Q_2/(n-r)}. \quad (4)$$

Значимость определяется при помощи доверительного уровня вероятности

$$\alpha_* = P\{F > F_*\},$$

где F_* наблюдаемое значение статистики F . При $\alpha_* < \alpha = 0.05$ гипотеза $H_0 : \alpha_1 = \dots = \alpha_r = 0$ отвергается.

Матричная модель дисперсионного анализа

Пусть $Y = (y_{11}, \dots, x_{1n_1}, y_{21}, \dots, x_{2n_2}, \dots, y_{r1}, \dots, x_{rn_r})'$ – вектор наблюдений, $\sum_{i=1}^r n_i = n$, матрица плана размерности n на r имеет вид:

$$X = \begin{bmatrix} 1 & 1 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 1 & 1 & 0 & \dots & 0 \\ \hline 1 & 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 1 & 0 & 1 & \dots & 0 \\ \hline \dots & \dots & \dots & \dots & \dots \\ \hline 1 & -1 & -1 & \dots & -1 \\ \dots & \dots & \dots & \dots & \dots \\ 1 & -1 & -1 & \dots & -1 \end{bmatrix}, \quad \beta = \begin{bmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \dots \\ \alpha_{r-1} \end{bmatrix}, \quad Y = \begin{bmatrix} y_{11} \\ \vdots \\ y_{1n_1} \\ y_{21} \\ \vdots \\ y_{2n_2} \\ \vdots \\ y_{r1} \\ \dots \\ y_{rn_r} \end{bmatrix}$$

Матричный вид модели однофакторного дисперсионного анализа

$$y_{ij} = \mu + \alpha_i + e_{ij}, \quad i = 1, \dots, r, \quad j = 1, \dots, n_i, \quad \text{где } \sum_{i=1}^r \alpha_i = 0$$

$$Y = X\beta + \epsilon, \quad \text{где } \epsilon \sim N(\mathbb{O}, \sigma^2 \mathbb{I})$$

Оценки параметров по методу наименьших квадратов имеют вид:

$$\hat{\beta} = (X'X)^{-1}X'Y.$$

Для проверки гипотезы $H_0 : \alpha_1 = \alpha_2 = \dots = \alpha_r = 0$ используют модель с ограничением на параметры при помощи матрицы H размерности r (число параметров) на $s = r - 1$

$$H'\beta = \theta_0, \quad \text{где} \quad H' = \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix} \quad \text{и} \quad \theta_0 = \begin{bmatrix} 0 \\ 0 \\ \dots \\ 0 \end{bmatrix}$$

Нужно показать, что $Z = H'\hat{\beta} \sim \mathcal{N}_s(H'\beta, \sigma^2 D)$, где $D = H'(X'X)^{-1}H$ и $R_0^2 = (Y - X\hat{\beta})'(Y - X\hat{\beta}) \sim \sigma^2 \chi^2(n - r)$ распределены независимо. Кроме того $(Z - \theta_0)'D^{-1}(Z - \theta_0) = R_1^2 - R_0^2 \sim \sigma^2 \chi^2(s)$, где $R_1^2 = (Y - X\beta^*)'(Y - X\beta^*)$, β^* оценка параметров усеченной модели. Тогда при справедливости нулевой гипотезы

$$F = \frac{n - r}{s} \frac{R_1^2 - R_0^2}{R_0^2} \sim F(s, n - r).$$

Для получения оценок МНК дифференцируем квадратичную форму.

$$\mathcal{L}_1 = (Y - X\beta)^T(Y - X\beta) = \left(y_1 - \sum_{i=1}^m x_{1i}\beta_i\right)^2 + \dots + \left(y_n - \sum_{i=1}^m x_{ni}\beta_i\right)^2.$$

Система нормальных уравнений имеет вид

$$\begin{cases} -\frac{\partial \mathcal{L}_1}{\partial \beta_1} = 2 \left(x_{11} \left(y_1 - \sum_{i=1}^m x_{1i}\beta_i \right) + \dots + x_{n1} \left(y_n - \sum_{i=1}^m x_{ni}\beta_i \right) \right) = 0, \\ \dots \\ -\frac{\partial \mathcal{L}_1}{\partial \beta_m} = 2 \left(x_{1m} \left(y_1 - \sum_{i=1}^m x_{1i}\beta_i \right) + \dots + x_{nm} \left(y_n - \sum_{i=1}^m x_{ni}\beta_i \right) \right) = 0, \end{cases}$$

$$\frac{\partial \mathcal{L}_1}{\partial \beta} = 0 \iff X^T(Y - X\beta) = 0.$$

При ОМК с ограничением на параметры дифференцируется выражение

$$\begin{aligned} \mathcal{L}_2 = \lambda'(1, s)H'(s, m)\beta(m, 1) &= (\lambda_1, \dots, \lambda_s) \begin{bmatrix} h_{11} & \dots & h_{m1} \\ \dots & \dots & \dots \\ h_{1s} & \dots & h_{ms} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \dots \\ \beta_m \end{bmatrix} = \\ &= (\lambda_1, \dots, \lambda_s) \begin{bmatrix} \sum_{i=1}^m h_{i1}\beta_i \\ \dots \\ \sum_{i=1}^m h_{is}\beta_i \end{bmatrix} = \sum_{j=1}^s \lambda_j \left(\sum_{i=1}^m h_{ij}\beta_i \right). \end{aligned}$$

Система нормальных уравнений

$$\begin{cases} \frac{\partial \mathcal{L}_2}{\partial \beta_1} = \sum_{j=1}^s \lambda_j h_{1j} = 0, \\ \dots \\ \frac{\partial \mathcal{L}_2}{\partial \beta_m} = \sum_{j=1}^s \lambda_j h_{mj} = 0, \end{cases} \iff \iff \begin{bmatrix} h_{11} & \dots & h_{1s} \\ \dots & \dots & \dots \\ h_{m1} & \dots & h_{ms} \end{bmatrix} \begin{bmatrix} \lambda_1 \\ \dots \\ \lambda_s \end{bmatrix} = H\lambda = 0.$$

Определение. A^{-} - обобщенная обратная матрица, такая, что для любого вектора Y , при котором система $AX = Y$ совместна, $X = A^{-}Y$ является ее решением.

A^{-} существует $\iff AA^{-}A = A \iff A^{-}AA^{-} = A^{-}$.

Задача 1

Пусть A — симметричная матрица n на n ранга r , B — матрица m на n , BY — m линейных функций от $Y = (y_1, \dots, y_n)^T$, y_i независимые нормально распределенные случайные величины. Для независимости BY и $Y^T AY$ достаточно условия $BA = 0_{m,n}$, где $0_{m,n}$ матрица m на n из нулей.

Спектральное разложение симметричной матрицы A ранга r :

$$A = \lambda_1 P_1 P_1^T + \dots + \lambda_r P_r P_r^T, \quad \lambda_i \neq 0, \quad i = 1, \dots, r. \quad (5)$$

Пусть $BA = 0$. Умножим (5) слева на B и справа для примера на P_1 .

$$BAP_1 = \lambda_1 BP_1 P_1^T P_1 + \dots + \lambda_r BP_r P_r^T P_1 = 0_{m,1}.$$

Из ортогональности собственных векторов получаем $BP_1 = 0_{m,1}$.

Аналогично для остальных $BP_i = 0_{m,1}$, т.е. линейные функции BY не зависят от функций $P_i^T Y$ (если $a^T b = 0$, то для Y с независимыми компонентами $\mathbb{E}(a^T Y)(Y^T b) = a^T \mathbb{E} Y Y^T b = a^T b = 0$) и не зависят от $Y^T AY$, так как

$$Y^T AY = \sum \lambda_i (P_i^T Y)^2.$$

Задача 2

а) $X = X(X'X)^{-1}X'$,

б) $I - X(X'X)^{-1}X'$ идемпотентна.

а) $G = X(I - (X'X)^{-1}X'X)$,

$$\begin{aligned}\implies G'G &= (I - (X'X)^{-1}X'X)'X'X(I - (X'X)^{-1}X'X) = \\ &= (I - (X'X)^{-1}X'X)'(X'X - X'X(X'X)^{-1}X'X) = \\ &= (I - (X'X)^{-1}X'X)'(X'X - X'X) = 0, \implies G = 0.\end{aligned}$$

$$\begin{aligned}\text{б) } &(I - X(X'X)^{-1}X')(I - X(X'X)^{-1}X') = \\ &= I - X(X'X)^{-1}X' - X(X'X)^{-1}X' + \underbrace{X(X'X)^{-1}X'X(X'X)^{-1}X'}_X = \\ &= I - X(X'X)^{-1}X' .\end{aligned}$$

Задача 3

Пусть $\beta = (\beta_1, \dots, \beta_m)'$ вектор параметров, X - матрица плана ранга r , $(X'X)^-$ обобщенная обратная, $\hat{\beta} = (X'X)^-X'Y$, $\hat{\beta}$ решение уравнения $(X'X)\beta = X'Y$. Матрица H' размерности s на m ранга s . Показать, что распределены независимо $Z = H'\hat{\beta}$ и

$$R_0^2 = (Y - X\hat{\beta})'(Y - X\hat{\beta}).$$

Введем матрицу C такую, что $H = X'C$. Тогда

$$\begin{aligned} H'\hat{\beta} &= \overbrace{C'X(X'X)^-X'}^B Y = BY. \\ R_0^2 &= Y'Y - \hat{\beta}'X'Y - Y'X\hat{\beta} + \hat{\beta}' \underbrace{X'X\hat{\beta}}_{X'Y} = Y'Y - Y'X\hat{\beta}, \end{aligned} \quad (6)$$

$$R_0^2 = Y'Y - Y'X(X'X)^-X'Y = Y' \underbrace{(I - X(X'X)^-X')}_A Y = Y'AY.$$

$$\begin{aligned} \text{Таким образом } BA &= C'X(X'X)^-X'(I - X(X'X)^-X') = \\ &= C' \underbrace{(X(X'X)^-X' - X(X'X)^-X'X(X'X)^-X')}_X = 0. \end{aligned}$$

Теорема Фишера-Кочрена

Пусть $Y = (y_1, \dots, y_n)$ с независимыми компонентами, $y_i \sim \mathcal{N}(0, 1)$, квадратичные формы Q_1, \dots, Q_k рангов n_1, \dots, n_k , такие что

$$Y^T Y = Q_1 + \dots + Q_k.$$

Тогда $n = n_1 + \dots + n_k \implies Q_i \sim \chi^2(n_i)$ и Q_1, \dots, Q_k независимы.

Необходимость очевидна. Достаточность. Пусть $n = n_1 + \dots + n_k$ и $Q_i = Y^T A_i Y$, где матрица A_i ранга n_i . Тогда существует матрица B_i размерности n_i на n , при помощи которой квадратичная форма приводится к диагональному виду, то есть $Q_i = (B_i Y)^T \Delta_i (B_i Y)$, где Δ_i диагональная матрица с элементами ± 1 . Соберем k матриц B_i в одну матрицу B размерности n на n .

$$\begin{aligned} [B_1^T(n, n_1) | \dots | B_k^T(n, n_k)] \cdot \begin{bmatrix} \Delta_1 & \dots & 0 \\ \vdots & \dots & \vdots \\ 0 & \dots & \Delta_k \end{bmatrix} \cdot \begin{bmatrix} B_1(n_1, n) \\ \vdots \\ B_k(n_k, n) \end{bmatrix} &= \sum_i B_i^T \Delta_i B_i \\ \implies Y^T Y &= \sum_i Q_i = \sum_i (B_i Y)^T \Delta_i (B_i Y) = \\ &= Y^T \left(\sum_i B_i^T \Delta_i B_i \right) Y = Y^T B^T \Delta B Y, \implies B^T \Delta B = I. \end{aligned}$$

Ранг матрицы B должен быть равен n , так как с одной стороны, он не может быть больше n , а с другой ранг произведения матриц не превышает минимального ранга сомножителей,

$$n = \text{rank}(I) = \text{rank}((B^T \Delta)B) \leq \min(\text{rank}(B^T \Delta), \text{rank}(B)).$$

Тогда $\Delta = (B^T)^{-1}B^{-1}$ положительно определенная матрица, следовательно, $\Delta = I$ и матрица B является ортогональной. Тогда компоненты вектора $X = BY$ независимы и нормальны, кроме того $Y^T Y = X^T X$.

$$\begin{aligned} Q_1 &= x_1^2 + \dots + x_{n_1}^2, \\ Q_2 &= x_{n_1+1}^2 + \dots + x_{n_1+n_1}^2, \\ Q_3 &= x_{n_1+n_2+1}^2 + \dots + x_{n_1+n_2+n_3}^2, \\ &\dots \end{aligned}$$

Осюда все Q_i независимы и имеют распределение $\chi^2(n_i)$ соответственно. Таким образом достаточность $n = n_1 + \dots + n_k$ установлена.

Для того чтобы квадратичная форма $Y^T A Y$ имела распределение χ^2 , необходимо и достаточно условие идемпотентности матрицы $A = A^2$, причем $df = \text{rank}(A) = \text{Tr}A$.

Достаточность. Пусть $A = A^2$. Применим теорему Фишера-Кочрена к квадратичной форме ранга n

$$Y^T Y = Y^T A Y + Y^T (I - A) Y$$

и воспользуемся неравенством Сильвестра для $A = A^2$ и $C = I - A$,

$$\begin{aligned} \text{rank}(A) + \text{rank}(C) - \text{rank}(AC) &\leq n \\ \text{rank}(A) + \text{rank}(I - A) - \underbrace{\text{rank}(A(I - A))}_{=0} &\leq n \end{aligned}$$

$$n = \text{rank}(I) = \text{rank}(A + I - A) \leq \text{rank}(A) + \text{rank}(I - A),$$

Таким образом,

$$n = \text{rank}(A) + \text{rank}(I - A),$$

и по теореме Фишера-Кочрена $Y^T A Y \sim \chi^2(\text{rank}(A))$.

Для доказательства необходимости (пусть квадратичная форма $Y^T A Y$ имеет распределение χ^2) заметим, что существует ортогональная матрица C такая, что при преобразовании $Y = CX$

$$Y^T A Y = X^T C^T A C X = \lambda_1 x_1^2 + \dots + \lambda_m x_m^2,$$
$$Y^T Y = X^T X = x_1^2 + \dots + x_n^2,$$

где $\lambda_1, \dots, \lambda_m$ ненулевые собственные числа матрицы A , $m = \text{rank}(A)$. Так как $Y^T A Y \sim \chi^2(p)$, то характеристическая функция суммы $\lambda_1 x_1^2 + \dots + \lambda_m x_m^2$ имеет вид

$$[(1 - 2i\lambda_1 t) \dots (1 - 2i\lambda_m t)]^{-\frac{1}{2}} = (1 - 2it)^{-\frac{p}{2}}.$$

Сравнение х.ф. дает $p = m$, все $\lambda_i = 1$, $C^T A C$ - диагональная матрица с нулями или единицами на главной диагонали.

$$C^T A C = C^T A C C^T A C = C^T A^2 C, \implies A^2 = A.$$

Число степеней свободы равно $p = m = \text{rank}(A)$, то есть $\text{Tr}A$, так как $A^2 = A$.

Задача 4

- а) $Z = H'\hat{\beta} \sim \mathcal{N}_s(H'\beta, \sigma^2 D)$, где $D = H'(X'X)^{-1}H$ и
б) $R_0^2 \sim \sigma^2 \chi^2(n-r)$, где r – ранг матрицы плана X .

$$\begin{aligned} & E(\hat{\beta} - E\hat{\beta})(\hat{\beta} - E\hat{\beta})' = \\ &= E((X'X)^{-1}X'(Y - EY))((X'X)^{-1}X'(Y - EY))' = \\ &= (X'X)^{-1}X'\sigma^2 I((X'X)^{-1}X')' = \\ &= \sigma^2(X'X)^{-1}X'X(X'X)^{-1} = \sigma^2(X'X)^{-1}. \text{ Отсюда } D = H'(X'X)^{-1}H. \end{aligned}$$

$$\begin{aligned} & \text{б) } R_0^2 = Y'(I - X(X'X)^{-1}X')Y, \\ & \text{матрица } I - X(X'X)^{-1}X' \text{ идемпотентна, } \implies R_0^2 \sim \sigma^2 \chi^2(n-r), \\ & \text{где } n-r = \text{Tr}I - \text{Tr}X(X'X)^{-1}X' = \\ &= n - \text{Tr}(X'X)^{-1}X'X = n - \text{rang}(X'X) = n - \text{rang}(X) = n - r. \end{aligned}$$

Задача 5

Минимизация выражения $(Y - X\beta)'(Y - X\beta)$ по β при условии $H'\beta = \theta_0$ приводит к минимизации $\mathcal{L} = (Y - X\beta)'(Y - X\beta) + \lambda(H'\beta - \theta_0)$ и к системе уравнений:

$$\begin{cases} X'X\beta + H\lambda = X'Y, \\ H'\beta = \theta_0 \end{cases} \quad (7)$$

λ^* и β^* – ее решение. $Z = H'\hat{\beta} \sim \mathcal{N}_s(H'\beta, \sigma^2 D)$, где $D = H'(X'X)^{-1}H$, $R_0^2 = (Y - X\hat{\beta})'(Y - X\hat{\beta})$.

$$R_1^2 = (Y - X\beta^*)'(Y - X\beta^*) = (Z - \theta_0)'D^{-1}(Z - \theta_0) + R_0^2.$$

Заметим, что

$$\begin{aligned} R_1^2 &= (Y - X\beta^* + X\hat{\beta} - X\hat{\beta})'(Y - X\beta^* + X\hat{\beta} - X\hat{\beta}) = \\ &= (Y - X\hat{\beta})'(Y - X\hat{\beta}) + (\hat{\beta} - \beta^*)'X'X(\hat{\beta} - \beta^*), \\ \text{так как } &(\hat{\beta} - \beta^*)'X'(Y - X\hat{\beta}) = 0 \text{ по определению } \hat{\beta}. \end{aligned}$$

Достаточно показать, что

$$(\hat{\beta} - \beta^*)'X'X(\hat{\beta} - \beta^*) = (Z - \theta_0)'D^{-1}(Z - \theta_0).$$

$$(\hat{\beta} - \beta^*)'X'X(\hat{\beta} - \beta^*)' = (Z - \theta_0)'D^{-1}(Z - \theta_0)$$

Пусть $H = X'XC$, тогда

$$D = H'(X'X)^{-1}H = C'X'X(X'X)^{-1}X'XC = C'X'XC.$$

Вычитаем из $X'X\hat{\beta} = X'Y$ уравнение $X'X\beta^* + H\lambda^* = X'Y$.

$$X'X(\hat{\beta} - \beta^*) = H\lambda^* = X'XC\lambda^*, \quad \hat{\beta} - \beta^* = (X'X)^{-1}X'XC\lambda^*,$$

$$(\hat{\beta} - \beta^*)'X'X(\hat{\beta} - \beta^*) = \underbrace{\lambda^{*'}C'X'X(X'X)^{-1}X'XC\lambda^*}_{(\hat{\beta} - \beta^*)'}$$

$$= \lambda^{*'}C'X'XC\lambda^* = \lambda^{*'}D\lambda^*$$

$$Z - \theta_0 = H'\hat{\beta} - H'\beta^* = H'(\hat{\beta} - \beta^*) = C'X'X(X'X)^{-1}X'XC\lambda^* =$$

$$= C'X'XC\lambda^* = D\lambda^*$$

$$\implies \lambda^* = D^{-1}(Z - \theta_0),$$

$$\implies \lambda^{*'}D\lambda^* =$$

$$= (Z - \theta_0)'D^{-1}DD^{-1}(Z - \theta_0) = (Z - \theta_0)'D^{-1}(Z - \theta_0).$$