

Вопросы к экзамену
по курсу «Прикладной многомерный статистический анализ»
(Билеты)

1. Основные задачи многомерного статистического анализа.
2. Гильбертово пространство случайных величин. Задача о наилучшей линейной оценке.
3. Корреляционный и регрессионный анализ.
4. Коэффициенты корреляции.
5. Простая линейная регрессия. Метод наименьших квадратов.
Свойства оценок.
6. Множественная линейная регрессия. МНК. Свойства оценок.
7. Т-критерий значимости влияния фактора.
8. Проверка линейных гипотез. F-критерий.
9. Проверка адекватности модели. Коэффициент детерминации.
10. Равенство уравнений регрессии. Тест Чоу.
11. Фиктивные переменные.
12. Модель линейной регрессии с гетероскедастичностью.
13. Модель линейной регрессии с автокорреляцией в ошибках.
Критерий Дарбина-Уотсона.
14. Бинарная модель дискретного выбора.
15. Многомерная модель дискретного выбора в случае
упорядоченных альтернатив.
16. Многомерная модель дискретного выбора в случае
упорядоченных альтернатив.
17. Однофакторный дисперсионный анализ.
18. Двухфакторный дисперсионный анализ.
19. Дискриминантный анализ: постановка задачи и ее решение в
случае известных параметров.
20. Решение задачи дискриминантного анализа в случае
неизвестных параметров.

21. Кластерный анализ: постановка задачи и основные понятия.
22. Кластерный анализ: схема последовательного построения факторов.

Дополнительные вопросы к экзамену
по курсу «Прикладной многомерный статистический анализ»
(Знать наизусть!)

1. Основные задачи многомерного статистического анализа:
 - корреляционный анализ
 - регрессионный анализ
 - снижение размерности
 - дисперсионный анализ
 - дискриминантный анализ
 - кластерный анализ
2. Гильбертово пространство случайных величин.
3. Что такое наилучшая линейная оценка.
4. Лемма о перпендикуляре.
5. Простой коэффициент корреляции и что он измеряет.
6. Множественный коэффициент корреляции и что он измеряет.
7. Частный коэффициент корреляции и что он измеряет.
8. Множественная линейная регрессия: модель и основные ограничения.
9. Описание МНК для оценки параметров.
10. Явный вид оценок параметров по МНК.
11. Общая схема проверки гипотезы о параметре.
12. Для чего используется Т-критерий.
13. Основное различие Т-критерия и F-критерия в задаче проверки значимости влияния фактора.
14. Адекватность модели. Постановка задачи.
15. Коэффициент детерминации и что он измеряет.
16. Задача о равенстве двух регрессий.
17. Что такое модель с гетероскедастичностью в ошибках?
18. Что проверяет тест Дарбина-Уотсона?
19. Описание модели бинарного выбора?
20. Что такое Пробит-модель?
21. Что такое Логит-модель?
22. В чем различие моделей упорядоченного и неупорядоченного выбора?
23. Основная задача в однофакторном дисперсионном анализе.
24. Основная задача в двухфакторном дисперсионном анализе.

25. Основная задача дискриминантного анализа.
26. Кластерный анализ: постановка задачи.
27. Кластерный анализ: последовательное построение факторов.

Прикладной многомерный статистический анализ (ПМСА)

04.09

Хоклов Юрий Степанович yskhokhlov@yandex.ru

титулование: "Многомерный статистич. анализ"

Экз. магистерский (Lec 2a)

Введение.

$\vec{X} = (x_1, \dots, x_p)^T$ — случайный вектор с неизвест. расположением
р разнотипных характеристиках.

Векторы: $\vec{X}_1, \dots, \vec{X}_N$.

$\vec{X}_j = (x_{1j}, \dots, x_{pj})$

$X = (X_{kj})$ — матрица наблюдений.

Матрица средних наблюдений:

$$\bar{X}_k = \frac{1}{N} \sum_{j=1}^N X_{kj}$$

$$\vec{x}_j = (x_{1j}, \dots, x_{pj})^T = (x_{1j} - \bar{x}_1, \dots, x_{pj} - \bar{x}_p)$$

Тема 1. Основные задачи многомерного стат. анализа

В курсе предполагается, что величина — многомерная нормальная.

- 1) Есть ли зависимость между различными характеристиками \Rightarrow корреляционный анализ
- 2) Найти взаимную функциональную связь одиних характеристиках от других.
 \Rightarrow функциональный,
регрессионный анализ
- 3) Задачи описание различия
 \Rightarrow анализ главных компонент
 \Rightarrow дисперсионный анализ
- 4) Есть ли вхождение в среднем на нех. признаках заданных набора характеристиках?
 \Rightarrow дисперсионный анализ
- 5) Помехи приведены для описание наблюдений к одному из классификации,]
[классификации]
 \Rightarrow дисперсионный анализ
- 6) Есть перво-по вторые признаки \Rightarrow классификационный анализ

Теорема 2. Гауссово представление случайных величин.

Нужно решить задачу о наименьшем квадратном приближении.

§1. Гауссово приближение случайных величин

] ($\mathcal{L}, \mathbb{E}, \mathbb{P}$) - фиксированное вероятностное пр-во.

Одномерный резерв L_2 пр-ва с. в. ξ : $\mathbb{E}(|\xi|^2) < \infty$ (1)

Проверим, что L_2 - линейное пр-во. (конечно)

Для $\forall \xi_1, \xi_2 \in L_2$ очевидно что $(\xi_1, \xi_2) := \mathbb{E}(\xi_1 \cdot \xi_2)$ (2)

Оно \exists и конечно. Теперь можно проверить, что выполняется следующее сл-во:

$$1) (\xi, \xi) \geq 0 \text{ и } (\xi, \xi) = 0 \iff \xi = 0 \text{ н.ч.}$$

$$2) (\xi_1, \xi_2) \in L_2 : (\xi_1, \xi_2) = (\xi_2, \xi_1)$$

$$3) \forall \xi_1, \xi_2, \xi_3 \in L_2 \quad \forall c_1, c_2 \in \mathbb{R}^1$$

$$(c_1 \xi_1 + c_2 \xi_2, \xi_3) = c_1 (\xi_1, \xi_3) + c_2 (\xi_2, \xi_3)$$

$\Rightarrow (\xi_1, \xi_2)$ -свойство произведения в L_2 .

Можно очевидным образом $\forall \xi \in L_2$:

$$\|\xi\| = \sqrt{(\xi, \xi)} = (\mathbb{E}(|\xi|^2))^{1/2} \quad (3)$$

Норма - длина элемента в линейном пр-ве

Норма \rightarrow длина \rightarrow сх-ма

Оп. Требование с. в. $\{\xi_n\}$ сходится в среднемквадр.

$$\text{к. с. в. } \xi_0 \iff \|\xi_n - \xi_0\|^2 = \mathbb{E}(|\xi_n - \xi_0|^2) \rightarrow 0 \text{ при } n \rightarrow \infty \quad (4)$$

L_2 -норма определяет сходимость в среднемквадр

$\Rightarrow L_2$ -недеформативное пр-во.

$\forall \xi_1, \xi_2 \in L_2$ очевидно что ξ_1 и ξ_2 не могут быть

$$\cos(\varphi) = \frac{(\xi_1, \xi_2)}{\|\xi_1\| \cdot \|\xi_2\|} \quad \text{см. рисунок.} \quad (5)$$

Оп. ξ_1 и $\xi_2 \in L_2$ ортогональны $\iff \xi_1 \perp \xi_2 : (\xi_1, \xi_2) = \mathbb{E}(\xi_1 \cdot \xi_2) = 0$.

Замечание!] ξ_1 и $\xi_2 \in L_2$.] $\mathbb{E}(\xi_1) = \mathbb{E}(\xi_2) = 0$.

1) $(\xi_1, \xi_2) = \mathbb{E}(\xi_1 \cdot \xi_2) = \text{cov}(\xi_1, \xi_2) = 0 \iff$ когда ξ_1 и ξ_2 не коррелируют

с. в. ξ_1 и ξ_2 ортогональны - т.е. они независимы

2) $\cos \varphi = \rho(\xi_1, \xi_2)$ - коэф. корреляции

§2. Задача нахождения минимальной оценки для с.б.

$\hat{z} \in L_2$ - наименьшее линейное подпространство
 $h \in L_2$ (не обозн. $z \in L$)

Оп. с.б. \hat{z} наз-ся наименьшим линейным приближением зе
 z в уп-бе L , если:

- 1) $\hat{z} \in L$
- 2) $\forall z \in L \quad \|h - \hat{z}\| \leq \|h - z\|$

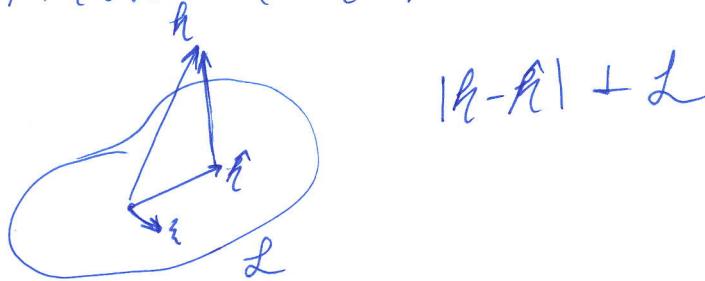
(7)

лемма о перпендикуляре

\hat{z} есть наименьшее линейное приближение z в $L \Leftrightarrow$

- 1) $\hat{z} \in L$
- 2) $\forall z \in L \quad (h - \hat{z}, z) = 0$

(8)



$L = z_0 + L_0$, где L_0 - нрн. подпр-во.

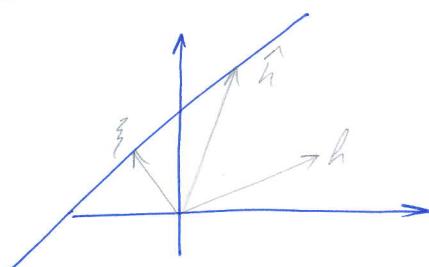
L -линейность

Определение симметричности ЭЛ-на доказано не м.

лемма о перпендикуляре

h есть наименьшее лин. приближение \hat{h} в L (линейность) \Leftrightarrow

- 1) $\hat{h} \in L$
- 2) $(h - \hat{h}, z - \hat{h}) = 0 \quad \forall z \in L$



Задача 2 с.б. { и h. Недоказано наименее пакутие
представление h в виде линейной комбинации гр-унм от ξ .

В наимене мн. упр-ва h равнотрив с.б. альгебраческого вида:

$$\alpha + \beta \cdot \xi = \alpha \cdot 1 + \beta \cdot \xi = \alpha \cdot \xi_1 + \beta \cdot \xi_2, \quad \text{так как } \xi_1 = 1 \\ \xi_2 = \xi.$$

\hat{h} -наименее мн. представление h в L, т.е. $\hat{h} = \hat{\alpha} + \hat{\beta} \cdot \xi$

Приемами линий о перенесении:

$$(h - \hat{h}, \xi_1) = \mathbb{E}[(h - \hat{h}) \cdot \xi_1] = \mathbb{E}[h - \hat{\alpha} - \hat{\beta} \cdot \xi] = \mathbb{E}(h) - \hat{\alpha} - \hat{\beta} \cdot \mathbb{E}(\xi) = 0 \quad (9)$$

$$(h - \hat{h}, \xi_2) = \mathbb{E}[(h - \hat{h}) \cdot \xi_2] = \mathbb{E}[(h - \hat{\alpha} - \hat{\beta} \cdot \xi) \cdot \xi] = \mathbb{E}[h \cdot \xi] - \hat{\alpha} \cdot \mathbb{E}(\xi) - \hat{\beta} \cdot \mathbb{E}(\xi) = 0 \quad (10)$$

$$\Rightarrow \hat{\beta} = \frac{\text{cov}(\xi, h)}{D(\xi)} = \rho(\xi, h) \cdot \frac{\sigma_h}{\sigma_\xi}$$

$$\hat{\alpha} = \mathbb{E}(h) - \rho(\xi, h) \frac{\sigma_h}{\sigma_\xi} \mathbb{E}(\xi)$$

$$\Rightarrow \hat{h} - \mathbb{E}\xi = \rho(\xi, h) \cdot \frac{\sigma_h}{\sigma_\xi} [\xi - \mathbb{E}\xi].$$

Теорема 3. Многомерное нормальное распределение.

Случайный вектор $\xi = (\xi_1, \dots, \xi_n)^T$ имеет многомерное нормальное распределение в \mathbb{R}^n , если он имеет многомерное азимутальное распределение вида:

$$\begin{aligned} P_\xi(x_1, \dots, x_n) &= P_\xi(\mathbf{x}) = (2\pi)^{-\frac{n}{2}} \det(A) \cdot \exp\left\{-\frac{1}{2}(A(x-m), (x-m))\right\} = \\ &= (2\pi)^{-\frac{n}{2}} \det(A) \exp\left\{-\frac{1}{2} \sum_{i,j=1}^n a_{ij}(x_i - m_i)(x_j - m_j)\right\} \end{aligned}$$

$$m = (m_1, \dots, m_n)^T \in \mathbb{R}^n$$

$A = (a_{ij})$ — симметрическая положительного определения матрица

$$\Sigma = A^{-1} = (\sigma_{pq})$$

Определение: m и Σ — параметры многомерного норм. распред.

Если $t = (t_1, \dots, t_n) \in \mathbb{R}^n$, то

$$\varphi_t(t) := \mathbb{E}[e^{i(t, \xi)}] = \mathbb{E}[e^{i(t_1 \xi_1 + \dots + t_n \xi_n)}] \quad (1)$$

$$(1) \Rightarrow \varphi_t = \exp[i(t, m) - \frac{1}{2}(\Sigma \cdot t, t)] \quad (2)$$

$$\Rightarrow \mathbb{E}(\xi_k) = m_k$$

$$\text{cov}(\xi_p, \xi_q) = \sigma_{pq}.$$

Если Σ не положена, то многомерное норм. распред. может называться некоррелированным.

Если распред. векториз., то \exists подразделение на меньшую разницу, в которой М.К.Р. не выражено.

Если $\Sigma = I$ (единичная матрица), а $m = \vec{0}$, то распределение наз-ся сингарифмом.

Теорема 3.1.

Если $C: \mathbb{R}^m \rightarrow \mathbb{R}^n$ — линейное (на всей вып.), то, если $\xi \sim N(m, \Sigma)$, то $C \cdot \xi \sim N(Cm, C \cdot \Sigma \cdot C^T)$

Следствие. ξ имеет сингарифм. многомерное норм. распред. в \mathbb{R}^n , C — ортогональная проекция ξ на \mathbb{R}^m , $m < n \Rightarrow$ вектор $h = C \cdot \xi$ имеет сингарифм. норм. распред. в \mathbb{R}^m .

Теорема 3.2

$\xi = (\xi_1, \dots, \xi_n)$ имеет норм. распред. с параметрами m, Σ :
 $\xi \sim N(m, \Sigma)$. Разбейте ξ на 2 независимых ξ' и ξ'' разницами n_1 и n_2 , $n_1 + n_2 = n$.

$$\text{Тогда } m^T = (m', m'')^T \text{ и } \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

1) $\xi' \sim N(m', \Sigma_{11})$

2) ξ' и ξ'' независимы $\Leftrightarrow \Sigma_{12} = \Sigma_{21}^T$ состоят из нулей.

Теорема 3.3.

Если $\xi \sim N(m, \Sigma)$, то условное распред. ξ'' при условии, что $\xi' = X'$ является многомерной норм. с вектором средних $m'' + \Sigma_{21} \cdot \Sigma_{11}^{-1}(X' - m')$ и ковариационной матрицей

$$\Sigma_{22,1} = \Sigma_{22} - \Sigma_{21} \cdot \Sigma_{11}^{-1} \cdot \Sigma_{12}$$

Теорема 3.4.

$\xi \sim N(m, \Sigma)$ — небивариант. Тогда с.в. $(\xi - m)^T \cdot \Sigma^{-1} (\xi - m)$ имеет χ^2 -распределение с n степенями свободы.

Таблица 4. Понятие о корреляционном и регрессионном анализе.

§1. Основные задачи корреляционного и регрессионного анализа

Задачи 1-го типа: есть ли связь между с.в. и показателем она смысла. Обычно применяют коэффиц. корреляции.
 \Rightarrow корреляционный анализ.

Пусть есть зависимость, что связь есть и она линейная смыла.

Из курса ТВ известно, что наилучшее линейное применение в средней классификации с.в. Y с помощью с.в. X_1, \dots, X_n занимается с помощью устойчивого мат. ожидания = q -типа регрессии.

$$\hat{Y} = g(X_1, \dots, X_n) + \varepsilon$$

$$g(X_1, \dots, X_n) = \mathbb{E}[Y | X_1 = x_1, \dots, X_n = x_n]$$

\Rightarrow Регрессионный анализ

- 1) Выбор класса моделей (будет изучать только линейные)
- 2) Нахождение оценок параметров
- 3) Проверка номенклатуры параметров моделей
- 4) Проверка адекватности моделей
- 5) Проверка выполнения основных ограничений

§2. Корреляционный коэффициент

1. Простой корреляционный коэффициент

Опред. Простой корреляционный коэффициент простых независимых с.в. ξ_1 и ξ_2 называют ρ -корр. коэф. от ξ_1 и ξ_2 :

$$\rho = \rho(\xi_1, \xi_2) = \frac{\text{cov}(\xi_1, \xi_2)}{\sqrt{D(\xi_1) \cdot D(\xi_2)}} \quad D \neq 0$$

Теорема 4.1.

- 1) $|\rho(\xi_1, \xi_2)| \leq 1$
- 2) $\rho(\xi_1, \xi_2) = 0 \Leftrightarrow \xi_1$ и ξ_2 - независимы, т.е. $\text{cov}(\xi_1, \xi_2) = 0$.
- 3) Если $h_1 = a_1 \xi_1 + b_1 \quad a_1, a_2 \neq 0$
 $h_2 = a_2 \cdot \xi_2 + b_2$ с помощью же знако
 $\text{то } \rho(h_1, h_2) = \rho(\xi_1, \xi_2) \cdot \text{sign}(a_1, a_2)$
- 4) если $\rho(\xi_1, \xi_2) = \pm 1$, то $\exists a > 0$ и $b \in \mathbb{R}^1$: $\xi_2 = a \cdot \xi_1 + b$
 $\text{Если } \rho(\xi_1, \xi_2) = -1$, то $\exists a < 0$ и $b \in \mathbb{R}^1$: $\xi_2 = a \cdot \xi_1 + b$
ничего другого быть не может

Если $\rho > 0$, то \hat{Y}_1 и \hat{Y}_2 положит. коррелированы
(если одна растянута, то другая растянута в среднем)

Если $\rho < 0$, то \hat{Y}_1 и \hat{Y}_2 отрицательно коррелированы

$0 < \rho^2 \leq 0,25$ — слабо корр.

$0,25 < \rho^2 \leq 0,5$ — средне корр.

$\rho^2 > 0,5$ — сильно корр.

Таким образом для с.в. Y и X . Рассмотрим наименее квадратичное приближение \hat{Y} через X .

$$D(\hat{Y}) = \rho^2(X, Y) \cdot D(Y).$$

Если $e = Y - \hat{Y}$ — ошибка приближения, то

$$D(e) = 1 - \rho^2(X, Y) \cdot D(Y)$$

$\rho^2(X, Y)$ называем, какую горю членение с.в. Y можно обяснить линейной зависимостью X .

Задача. Есть ли связь X и Y ?

18.09

Равнозначно проверяется гипотеза

$$H_0: \rho = 0$$

(3)

$$H_1: \rho \neq 0$$

Предположим, что пара (X, Y) имеет двумерное нормальное распределение.

Таким образом имеем повторную двумерную выборку $(X_1, Y_1), \dots (X_n, Y_n)$ из двумерного норм. распределения и проверяем гипотезу (3).

Рассмотрим выборочный коррел. коэффиц. корреляции R :

$$R = \frac{S_{xy}}{S_x \cdot S_y} \quad (4)$$

$$S_{xy} = \frac{1}{N} \sum_j (X_j - \bar{X})(Y_j - \bar{Y})$$

$$S_x = \frac{1}{N} \sum_j (X_j - \bar{X})^2, \quad \bar{X} = \frac{1}{N} \sum X_j$$

$$S_y = \frac{1}{N} \sum_j (Y_j - \bar{Y})^2, \quad \bar{Y} = \frac{1}{N} \sum Y_j$$

$$\text{Рассмотрим с.в. } T_{N-2} = \frac{R}{\sqrt{1-R^2}} \sqrt{N-2} \quad (5)$$

Теорема

Если верна H_0 , то с.в. T_{N-2} имеет распределение

Симметрична $T_{N-2} \sim S_t$ с $(N-2)$ степенями свободы.

Для заданного уровня значимости $\alpha: 0 < \alpha < 1$ и подчиняется распределению Стьюдента находить критическое значение $t_{n-2}(\alpha) > 0$: $P(|T_{n-2}| > t_{n-2}(\alpha) | H_0) = \alpha$

(6)



1) Если реальное наблюдаемое значение T_{n-2}^* статистики T_{n-2} :
 $|T_{n-2}^*| > t_{n-2}(\alpha) \Rightarrow$ отвергаем $H_0 \Rightarrow H_1$ верна. (некорректное значение X на Y)

Если $|T_{n-2}^*| \leq t_{n-2}(\alpha) \Rightarrow$ не можем отвергнуть H_0 . (значение X на Y не однозначно).

В пакетах прикладных программ по статистике имеются программы для односторонней проверки (на какой же, задаваемый, результат).
 Тогда T_{n-2}^* — реально наблюдаемое значение статистики T_{n-2} . Решением верно ли это, что оно с.в. больше t_{n-2}^* :

$$P(|T_{n-2}| > |T_{n-2}^*| | H_0) \quad (7)$$

Дениситометрический уровень значимости: p-value.

Чем меньше p-value, тем более это свидетельство в пользу H_1 .

2. Многовariateльный корреляционный коэффициент

Изучение изменение с.в. Y в зависимости от нескольких факторов X_1, \dots, X_p , где $p \geq 2$. Рассмотрим наименее критическое приближение

$$\hat{Y} = \alpha + \beta_1 X_1 + \dots + \beta_p X_p \quad (8)$$

через X_1, \dots, X_p .

Определим многовariateльный корреляционный коэффициент по-сле что

$$R_{Y, X_1, \dots, X_p} := P(Y, \hat{Y}) \quad (9)$$

$e = Y - \hat{Y}$ — ошибка приближения.

Ошибка оценивается по формулам $X \Rightarrow e \sim \hat{Y}$ — некорр.

$$Y = \hat{Y} + e \Rightarrow D(Y) = D(\hat{Y}) + D(e)$$

$$\Rightarrow D(\hat{Y}) = R_{Y, X_1, \dots, X_p}^2 \cdot D(Y) \quad (10)$$

$$D(e) = (1 - R_{Y, X_1, \dots, X_p}^2) D(Y) \quad (11)$$

Квадрат многовariateльного корреляционного коэффициента называется, какую долю суммы квадратов X ($D(Y)$)

у大局е обозначим совместное минимальное значение
коэффициента джаккноров.

Максима.

1) $|R_{Y, x_1, \dots, x_p}| \leq 1$

2) $|R_{Y, x_1, \dots, x_p}| = \max$ по модулю коэффициенту корр. корр. между
 Y и $C_0 + C_1 x_1 + \dots + C_p x_p$

3) Если $R_{Y, x_1, \dots, x_p} = 0$, тогда наименьшее значение приближение
 $\hat{Y} = E(Y)$.

Если $R_{Y, x_1, \dots, x_p} = 1$ ($n - 1$), то Y и \hat{Y} связанны нормой
поверхности.

3°. Частичный коэффициент корреляции.

Изучаемое влияние на с.в. Y наводит с.в. x_1, \dots, x_p , $p > 2$.
Внедрен неизвестное k : $1 \leq k \leq p$.

Задача: каково "чистое" влияние x_k на Y ?

Обозначим через C наводку всех остальных джаккноров
(остальных с.в.).

Наименьшее значение минимальное приближение
 \hat{Y}_C с.в. Y через C . Аналогично наименьшее значение X_k^C — наименьшее
минимальное приближение X_k через C .

Пусть $Z_Y = Y - \hat{Y}_C$

$Z_{X_k} = X_k - X_k^C$

Оп. Частичный коэффициент корреляции с.в. Y и джаккнора
 X_k находит из

$R_{Y X_k \cdot C} := \rho(Z_Y, Z_{X_k})$

(12)

Максима

1) $|R_{Y X_k \cdot C}| \leq 1$

2) $1 - R_{Y, x_1, \dots, x_k}^2 = (1 - R_{Y, x_1, \dots, x_{k-1}}^2)(1 - R_{Y, x_k, x_1, \dots, x_{k-1}}^2)$

3) Если X_C и X_k некорр. со всеми остальными джаккнорами
из C , то $R_{k|C} = \rho_{k|C}$

Чтобы найти частичный коэффициент корреляции, нужно
делить необратимый ранее множеством (дисперсией) удалено
единство добавленной ещё одного джаккнора.

$1 - \rho_{Y, X}^2$ — это необратимый множеством

Замечание. м. к. не оценивается никаким показателем не выше 2-го порядка и так. комбинации с. в., то и показательный, и генеральный коэффициент корр. можно принять, зная простые коэффициенты корр.

§3. Пример

По некоторым рабочим 37-ти однородных предпринимательских единиц производственного ресурса определение следующие показатели их работы:

$Y = X_0$ - среднепечатное количество копий на один рабочий час;

X_1 - среднепечатное число производимых граверов

X_2 - среднепечатное число единиц труда

Вычислим парные (простые) коэффициенты корр.:

$$R_{01} = 0,105$$

$$R_{02} = -0,024$$

$$R_{12} = -0,996$$

Полученное значение also производят оптимизацию работы в данной области.

Вычислим генеральный коэффициент корреляции:

$$R_{01 \cdot 2} = 0,907$$

$$R_{02 \cdot 1} = -0,906$$

Этот результаты дают с тем, что X_1 и X_2 сильно связаны между собой и оптимизировано корр., поэтому их совместное влияние находит свое выражение в показателе производительности в отдельности.

Таблица 5. Многомодельная линейная модель регрессии

§1. Классическая линейная модель.

У нас есть одна независимая переменная Y и несколько объясняющих факторов или предикторов X_1, \dots, X_m .

Для этого имеются N наблюдений и одновременных измерений этих величин и предполагается, что они связаны следующим образом: $\forall Y_j$ это называется функцией:

$$Y_j = g(X_{j1}, \dots, X_{jm}) + \varepsilon_j, \quad j=1, N \quad (1)$$

Основные ограничения:

1) Модель линейна по параметрам, т.е.

$$Y_j = \alpha + \beta_1 X_{j1} + \dots + \beta_m X_{jm} + \varepsilon_j \quad (2)$$

2) Равнотипные измерения для каждого, т.е. X_{jk} — одинаковые величины

3) $\mathbb{E}(\varepsilon_j) = 0 \quad \forall j$, т.е. нет систематических ошибок

4) $D(\varepsilon_j) = \sigma^2 > 0 \quad \forall j$ — одинаковая гетероскедастичность (одинаковая дисперсия)

5) ε_j и ε_k — независимы (независимы), $j \neq k$

6) ε_j имеет нормальное распределение

Давно Решёв следующее обозначение:

$$Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_N \end{bmatrix}, \quad \vec{X}_k = \begin{bmatrix} X_{1k} \\ \vdots \\ X_{Nk} \end{bmatrix}, \quad \vec{X}_0 = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}, \quad \boldsymbol{\theta} = \begin{bmatrix} \alpha \\ \beta_1 \\ \vdots \\ \beta_m \end{bmatrix} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_m \end{bmatrix}, \quad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_N \end{bmatrix},$$

$$X = (X_{jk}) = \begin{bmatrix} X_{00} & X_{11} & \dots & X_{1m} \\ \cdots & \cdots & \cdots & \cdots \\ X_{n0} & X_{n1} & \dots & X_{nm} \end{bmatrix} \Rightarrow$$

$$Y = X \cdot \boldsymbol{\theta} + \varepsilon = \theta_0 \vec{X}_0 + \theta_1 \vec{X}_1 + \dots + \theta_m \vec{X}_m + \varepsilon \quad (3)$$

Три основных ограничения:

$$\mathbb{E} Y = X \cdot \boldsymbol{\theta}$$

$$\sum_y = \sum_\varepsilon = \sigma^2 \cdot E \quad (4)$$

единичная матрица

Основные задачи:

- 1) Оценка параметров модели, т.е. $\theta_0, \dots, \theta_m$ и σ^2 .
- 2) Проверка гипотез о параметрах
- 3) Оценка значимых коэффициентов
- 4) Проверка адекватности модели
- 5) Проверка основных ограничений (как правило, это не является в силу линейной структуры, но это не всегда ходит.)

§ 2. Оценка параметров

Для оценки параметров θ приложим метод наименьших квадратов. Рассмотрим величину

$$Q(\theta) = \|Y - X\theta\|^2 = \sum_{j=1}^n [y_j - (x_{j0}\theta_0 + \dots + x_{jm}\theta_m)]^2 \rightarrow \min_{\theta} \quad (5)$$

Несколько условий экстремума: минимум $Q'(\theta) = 0$
 $\Rightarrow m+1$ уравнение.

После некоторого преобразования находим:

$$X^T X \cdot \theta = X^T Y \quad \text{-система нормальных уравнений} \quad (6)$$

Также получим $X^T X$ - невыполн. $\Rightarrow X$ имеет ранг $m+1 \Rightarrow$
 $\Rightarrow \vec{x}_0, \vec{x}_1, \dots, \vec{x}_m$ - лин. нез. (базис)

\Rightarrow решение системы

$$\hat{\theta} = (X^T X)^{-1} X^T Y = \theta + (X^T X)^{-1} X^T \varepsilon \quad (7)$$

$\hat{\theta}$ - оценка по МНК.

$\hat{Y} = X \cdot \hat{\theta}$ - вектор предсказанных значений
 $e = Y - \hat{Y}$ - вектор остатков

$$Y = \hat{Y} + e$$

$$\sigma^2 = D(e_j) = E(e_j^2)$$

$$\text{Рассмотрим: } \frac{1}{N} \sum_{j=1}^N e_j^2 \rightarrow \sigma^2, \quad N \rightarrow \infty$$

В качестве оценки σ^2 предполагаем форму:

$$\hat{\sigma}^2 = S^2 = \frac{1}{N-(m+1)} \sum_{j=1}^N e_j^2 \quad (8)$$

Чем хуже МНК? Рассмотрим на примере задачу о наилучшей линейной оценке вектора Y в линейном ур-ве L , которое порождено векторами $\vec{x}_0, \dots, \vec{x}_m$

§3. Свойства оценок параметров

Лемма 1.

Пусть Y -случ.ベктор с конечным ожиданием 2-го порядка, A -несингулярная матрица. Тогда ожидание вектора $Z = A \cdot Y$

$$\mathbb{E} Z = A \cdot \mathbb{E} Y$$

$$\Sigma_Z = A \cdot \Sigma_Y \cdot A^T$$

коэффициент. матрица

Пусть Y -вектор измерений, \hat{Y} -надежд. параметров. \hat{X} -оценка \hat{Y} с погрешностью $\hat{\varepsilon}$.

$$Y = \theta_0 \vec{X}_0 + \dots + \theta_m \vec{X}_m + \varepsilon$$

Оценка \hat{Y} параметров $\hat{\theta}$ наз-ся максимумом (по Y), если она имеет вид $\hat{\theta} = A \cdot Y$

Лемма 2. Пусть выполнены условия 1)-5) основных ограничений. Максимум оценка $\hat{\theta} = A \cdot Y$ параметров θ будет несущей. $\Leftrightarrow A \cdot X = E$

Доказательство: $\mathbb{E}(\hat{\theta}) = \mathbb{E}(AY) = \mathbb{E}(A(X\theta + \varepsilon)) = AX\theta = \theta \quad \forall \theta \Leftrightarrow A \cdot X = E$. ■

Предположение 1. Пусть выполнены условия 1)-5) основных ограничений. Тогда $\hat{\theta}$ -оценка по МНК - это максимум и несущий.

Доказательство:

$$\hat{\theta} = (X^T X)^{-1} X^T \cdot Y = A \cdot Y - \text{максимум}$$

Давно: $A \cdot X = (X^T X)^{-1} X^T \cdot X = E$ - несущий. ■

Предположение 2. Пусть выполнены 1)-5). Тогда матрица коэффициентов оценки $\hat{\theta}$ имеет вид:

$$\Sigma_{\hat{\theta}} = \sigma^2 \cdot (X^T X)^{-1} \quad (10)$$

Доказательство:

$$\hat{\theta} = (X^T X)^{-1} X^T \cdot Y = A \cdot Y$$

По условию 1: $\Sigma_{\hat{\theta}} = A \cdot \Sigma_Y A^T = (X^T X)^{-1} \cdot X^T \cdot \sigma^2 \cdot E \cdot X \cdot (X^T X)^{-1} = \sigma^2 (X^T X)^{-1}$. ■

Предположение 3. Пусть выполнены 1)-5) и диагональные элементы спадают к 0 при $n \rightarrow \infty$. Тогда оценки $\hat{\theta}$ параметров θ по МНК являются состоятельными.

(Это служит из условия несущийности и нер-ва Кохнича)

Теорема (Файса - Маркова)

Пусть выполнены условия 1)-5) основных ограничений, тогда оценка $\hat{\theta}$ по МНК наз-на ОЛ-СА оптимальной в среднеквадр. в смысле всех линейных и нелинейных оценок.

D-Б:

$$\hat{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \cdot \mathbf{X}^T \cdot \mathbf{Y} = A \cdot Y$$

$\tilde{\theta} = B \cdot Y$ - другая ин. нелиней. оценка

$$B = A + C$$

В силу несинг. $CX = 0$

$$\sum \tilde{\theta} = B \cdot \sum_y \cdot B^T = (A + C)\sigma^2 \cdot E \cdot (A + C)^T = (A \cdot A^T + A \cdot C^T + C \cdot A^T + C \cdot C^T) \cdot \sigma^2$$

$$A \cdot C^T = (C \cdot A^T)^T = (\mathbf{X} \cdot \mathbf{X}^T)^{-1} \cdot \mathbf{X}^T \cdot C^T = (\mathbf{X} \cdot \mathbf{X}^T)^{-1} \underbrace{(C \cdot X)^T}_{=0} = 0$$

$$\Rightarrow \sum \tilde{\theta} = (A \cdot A^T + C \cdot C^T) \sigma^2 = \sum \hat{\theta} + C \cdot C^T \sigma^2$$

$C \cdot C^T$ - неотриц. опред. матрица

$$\Rightarrow \sum \tilde{\theta} - \sum \hat{\theta} = (C \cdot C^T) \sigma^2 - \text{неотр. опред.}$$

\Rightarrow по гипотезам стат. неотр. мат

$$\Rightarrow D(\tilde{\theta}_k) \geq D(\hat{\theta}_k) \quad \forall k$$



Задача оценки где $\hat{\sigma}^2 = S^2 = \frac{1}{N-(m+1)} \sum_{j=1}^N e_j^2$

Предположение 4. Пусть выполнены условия 1)-6). Тогда оценка S^2 по МНК для σ^2 является нелинейной и состоятельной

§4. Построение доверительных интервалов

Всегда далее будем считать, что выполнены сб-ла 1)-6). Каждому $\hat{\theta}_k$ соответствует оценка о нормальности остатков

Теорема.

Пусть $\hat{\theta}$ -оценка по МНК наз-на ОЛ-СА, тогда справедливо след.

Д-Б:

1) $\hat{\theta}$ имеет многомерное нормальное распред. с вектором средних: $X \cdot \theta$ и матрицей ковариации: $\sum \hat{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \cdot \sigma^2$

2) С.в. $(\hat{\theta} - \theta)^T \cdot \mathbf{X}^T \cdot \mathbf{X} \cdot (\hat{\theta} - \theta) / \sigma^2$ имеет χ^2 -распределение с $(m+1)$ степенями свободы

3) С.в. $\sum_j (y_j - \hat{y}_j)^2 / \sigma^2 = \sum_j e_j^2 / \sigma^2$ имеет χ^2 -распределение с $(N-(m+1))$ степенями свободы

4) С.в. S^2 и $\hat{\theta}$ (или \hat{Y}) - независим.

Фиксируем некот. k : $0 \leq k \leq m$. Обозначим $B = (\mathbf{X}^T \mathbf{X})^{-1} = (b_{jk})$.

В силу теоремы С.в. $\hat{\theta}_k$ имеет одномерное норм. распред. со средним θ_k и дисперсией $b_{kk} \cdot \sigma^2$.

$$\Rightarrow \frac{\hat{\Omega}_k - \Omega_k}{\sqrt{S_{kk} \cdot S^2}} : \frac{\sigma}{S} = \frac{\hat{\Omega}_k - \Omega_k}{S \cdot \sqrt{S_{kk}}} = \frac{\hat{\Omega}_k - \Omega_k}{S_k} - \text{имеем расп. Согласно с } (N-(m+1)) \text{ степенным членам}$$

$S_k = S \cdot \sqrt{S_{kk}}$ - стандартная ошибка оценки пар-ра Ω_k .

По заданному доверительному уровню $\delta-1$ находим по таблицам критического константу $t_{N-(m+1)}(\delta) > 0$:

$$P(|T_k| < t_{N-(m+1)}(\delta)) = \delta$$

$$|T_k| = \left| \frac{\hat{\Omega}_k - \Omega_k}{S_k} \right| < t_{N-(m+1)}(\delta)$$

$$\Rightarrow \hat{\Omega}_k - t_{N-(m+1)}(\delta) \cdot S_k < \Omega_k < \hat{\Omega}_k + t_{N-(m+1)}(\delta) \cdot S_k \quad (\text{с вер-мн } \delta)$$

§5. Тестовая значимость одногранного.

Рассмотрим задачу: доказать H_0 (нарв): $0 \leq k \leq m$.

Проверка гипотеза $H_0: \Omega_k = \Omega_{k0}$

$$H_1: \Omega_k \neq \Omega_{k0}$$

В силу теории из §4 имеем, что с.в. $T_k = \frac{\hat{\Omega}_k - \Omega_k}{S_k}$ имеет при верной H_0 распределение Согласно с $(N-(m+1))$ степенным членом.

Для заданного уровня значимости $\alpha > 0$ находим критическое const $t_{N-(m+1)}^*(\alpha) \quad (= t_{N-(m+1)}(\delta=1-\alpha))$

$$P(|T_k| > t_{N-(m+1)}^*(\alpha) \mid H_0) = \alpha$$

Дано:

1) Если реально полученное значение T_k^* статистики T_k : $|T_k^*| > t_{N-(m+1)}^*(\alpha)$, то отвергаем H_0

2) В противоположной ситуации H_0 не отвергается.

] $\Omega_{k0} = 0$, т.е. гипотеза не отвергаем в силу на Y .

\Rightarrow Гипотеза о значимости влияния гранца X_k . (один гранец!)

\Rightarrow T-критерий Согласия. (см. выше)

Если гранец придан незначим, но его удалением из модели и не остатком заново.

§6. Проверка линейных гипотез

10. Однократные линейные гипотезы.

Член модели измерений:

$$Y = X \cdot \theta + \varepsilon \quad (1)$$

Всегда дает линейную зависимость 1)-6).

Проверка гипотезы:

$$H_0: A \cdot \theta = a$$

$$H_1: A \cdot \theta \neq a \quad (2)$$

A - заданные неизвестные параметры называются $p \times (m+1)$.

a - заданный неизвестный вектор называется p .

Предполагаем, что $p \leq m+1$.

$\text{rank}(A) = p$ (p мн/нез. ограничения)

Описание алгоритма проверки однократной линейной гипотезы:

1) Оцениваем модель (1) без гипотезы (2) и находим сумму квадратов остатков ESS_{VR}

2) Оцениваем модель (1) с гипотезой (2) и находим сумму квадратов остатков ESS_R

3) Тест верности H_0 с.в. $ESS_{VR}/10^2$ и $(ESS_R - ESS_{VR})/10^2$ независим и имеет χ^2 -распределение с $(N-(m+1))$ и p степенями свободы

4) Тест верности гипотезы H_0 с.в. $F = \frac{(ESS_R - ESS_{VR})/P}{ESS_{VR}/(N-(m+1))}$ - имеет

распределение Сnedecora - Фишера с $(p, N-(m+1))$ степенями свободы.

5) Для заданного α находим по таблицам критерия F значение $F(\alpha) = \alpha$: $P(F > F(\alpha) | H_0) = \alpha$

6) Если реальное значение F^* статистики F больше, чем $F^* > F(\alpha)$, то отвергаем H_0 . В противном случае H_0 не отвергается экспериментальным данным.

$\Rightarrow F$ -критерий. (с.в. выше)

(менее приемлемое значение для критерия)

F-критерий

2. Проверка значимости влияния одного фактора

Проверка

$H_0: \theta_k = 0$

нуль

$H_1: \theta_k \neq 0$

В случае простой линейной регрессии F- и T-критерии дают один и тот же ответ.

3. Проверка значимости влияния групповых факторов.Видели $1 \leq i_1 < i_2 < \dots < i_p \leq m$

Проверка ипотезы

$H_0: Q_{i_1} = Q_{i_2} = \dots = Q_{i_p} = 0$

нуль $\exists k:$

$H_1: Q_{ik} \neq 0.$

Если H_0 не отвергается, то все факторы Q_{ij} одновременно независимы (их значимое влияние не обнаружено).4. Проверка адекватности модели.

В нашем случае под адекватностью будем понимать, что полученный набор коэффициентов означает значение влияния исследуемого независимого.

$H_0: Q_1 = Q_2 = \dots = Q_m = 0$

нуль $\exists k: 1 \leq k \leq m$

$H_1: Q_k \neq 0.$

Можно показать, что проверка является следующим образом:

$$\sum_j (Y_j - \bar{Y})^2 = \sum_j (X_j - \hat{Y}_j)^2 + \sum_j (\hat{Y}_j - \bar{Y})^2$$

$TSS = ESS + RSS$

TSS - общее сумма квадратов - ^{значение} Y под влиянием всех
total sum of squares, или в независимых коэффициентах.ESS - основная сумма квадратов - на чисто независимом Y ,
error какое-либо отклонение независимым коэффициентамRSS - остаточная сумма квадратов - на чисто независимом Y ,
regression отклонение общему влиянию независимых коэффициентовТаким $ESS \leftrightarrow RSS$ (в независимые коэффициенты)

Справедлива следующая

теоремаЕсли верна H_0 , то

1) ESS и RSS - независимы

2) ESS/σ^2 и RSS/σ^2 имеют χ^2 -распред. с $(N-(m+1))$ и $(m+1)$ степенями дис3) С.В. $F = \frac{RSS/(m+1)}{ESS/(N-(m+1))}$ имеет распред. Снег-Фишера с $((m+1), (N-(m+1)))$ степ. дис.

Данное утверждение не относится к многомерной статистике.

Оп. Величина

$$R^2 := \frac{RSS}{TSS} = 1 - \frac{ESS}{TSS}$$
 назв. коэффициентом детерминации.

Он характеризует качество модели.

Для хороших моделей значение R^2 близко к 1.

$$\Rightarrow F = \frac{N-m-1}{m+1} \cdot \frac{R^2}{1-R^2}$$

5^o Пример.

"Влияние на 10%"

Изучаемое влияние качества пшеницы и зерна (Y) в зав-ии от количества BBT (X₁) и значимости количества зерна изображения пшеницы и BBT (X₂).

Dанные за 2 лет:

Y	X ₁	X ₂
100	100	100 < ложный уровень
106	104	99 < 10% к ложному нулю
107	106	110
120	111	126
110	111	113
116	115	103
123	120	102
133	124	103
137	126	98

Рассматриваем модель уравнений: (наи. вер. модели)

$$Y_j = \alpha_1 + \beta_1 X_{j1} + \beta_2 X_{j2} + \varepsilon_j$$

Причины к неоднозначности данных - наил. модель МНК:
(стартовая, наим. EViews)

Почему модель?

Var	Coef	SE	t-stat	p-value
X ₀	-49.34	24.06	-2.0506	0.0862
X ₁	1.364	0.143	9.5299	0.001
X ₂	0.114	0.143	0.7943	0.4573

$$S^2 = 12.92 \dots$$

$$F^* = 45.78$$

$$\Rightarrow p\text{-value} = 0.0002$$

$$R^2 = 0.9385$$

Было:

объяснение

1) Выявлено \checkmark значение бывшее значение

2) Построение графиков обнаружено наличие Y на 99 %

3) Равнозначность X_1 означает значение бывшее. \checkmark более

4) X_2 обнаружено незначимость гравитации. \Rightarrow

\Rightarrow Уравнение X_2 не является значимым.

Var	coef	SE	t-stat	p-value
X_0	-35.08	24.06	-2.25	0.059
X_1	1.345	0.137	9.8	0.000...

$R^2 = 0.92 \quad F^* = 96, \quad p\text{-value} = 0.000.$

6. Виды моделей

В нашем примере мы подбирали по формуле R^2 , но можно по нему оценить надежда (если 5 штуковых показателей останутся на 1% меньше, то они не нужны), а еще она меняется. \Rightarrow Виды неизменяющихся.

$$R^2_{\text{adj}} = 1 - (1 - R^2) \cdot \frac{N-1}{N-m-1} \quad - \text{независимый коэф. для оценки надежности.}$$

7. Графическое описание регрессии

(Когда соединяющее значение из разных количественных и качественных, несравнено, то оно из одной регрессии (однотипной модели).)

Сострано для набора измерений.

В некотором смысле след. модель измер.

$$Y_j = \beta_1' X_{j1} + \beta_2' X_{j2} + \dots + \beta_m' X_{jm} + \varepsilon_j, \quad j=1, N_1 \quad (1)$$

По библиотеке

$$Y_j = \beta_1'' X_{j1} + \beta_2'' X_{j2} + \dots + \beta_m'' X_{jm} + \varepsilon_j'', \quad j=\overline{N_1+1, N_1+N_2} \quad (2)$$

$$\sigma^2 = \sigma''^2$$

Проверка на наличие членов

$$H_0: \beta_k' = \beta_k'' \quad \forall k = 1, m$$

против

$$H_1: \exists k: \beta_k' \neq \beta_k''.$$

Факторный многочленный образ:

$$Y_j = \beta_1' \tilde{X}_{j1} + \dots + \beta_m' \tilde{X}_{jm} + \beta_1'' \tilde{X}_{j,m+1} + \dots + \beta_m'' \tilde{X}_{j,2m} + \varepsilon_j, \quad j=\overline{1, N_1+N_2} \quad (4)$$

где $\tilde{X}_{jk} = \begin{cases} x_{jk}, & 1 \leq j \leq N_1 \\ 0, & в \text{ иных случаях} \end{cases} \quad \forall k \leq m$

$$\tilde{X}_{jk} = \begin{cases} 0, & 1 \leq j \leq N_1 \\ x_{j,k-m}, & N_1+1 \leq j \leq N_1+N_2 \quad \text{и} \quad m+1 \leq k \leq 2m. \end{cases}$$

В рамках модели (4) проверяют наличие (3) с помощью F-критерия.

\Rightarrow Критерий Коу (Chow). Для него нужно можно знаять значение

N_1 - максимум количества наблюдений XDD-реализации

Пример. Учебная разница в уровне знаний музыки и математики в 2 классах.

$X_1 = \text{age} - \text{program}$

$X_2 = \text{edu} - \text{уровень образования}$

N_1 музыка, N_2 математика. $N_1 = N_2 = 75$.

$W = \text{wage} - \text{уровень заработка. (высшее) уро.}$

Балансировка ур-ия регрессии сег. 1-го:

16.10

$$W_j = \alpha + \beta_1 \cdot \text{age}_j + \beta_2 \cdot \text{edu}_j + \varepsilon_j$$

Оценки по МНК это уравнение для генеральной популяции.

$$W^{(1)} = -3.37 + 0.479 \cdot \text{age} + 3.943 \cdot \text{edu}$$

$$ESS_1 = 5672.328$$

Для генеральной популяции получим:

$$W^{(2)} = -0.20 + 0.414 \cdot \text{age} + 2.305 \cdot \text{edu}$$

$$ESS_2 = 1788.344$$

$$N_1 = N_2 = 72, m = 3 \text{ группы}$$

2 группы упр-я это как одно, но для временного отрезка ESS_{UR} = ESS₁ + ESS₂ = 7460.672

Оценки для временного отрезка, это близко это 2 группы упр-я.

Проверка оценки ур-ия по каждой генерации.

$$W = -3.06 + 4.78 \cdot \text{age} + 3.254 \cdot \text{edu}$$

$$ESS_R = 8080.443$$

Нахождение надежности знания статистикой F:

$$F^* = \frac{(ESS_R - ESS_{UR}) / m}{ESS_{UR} / (N - 2m)} = 3.342$$

150-6

По модулю расп-я Студенхорна-Рамера с (3,744) см. что оно надежно, это означает знание на уровне ≤ 0.05 .
=> Достоверно, в. формирующим 3/и музыки и математики обнаружено значимое различие.

§10. Однородная модель линии.

Числ. мт. модель

$$Y = X \cdot \theta + \varepsilon$$

(1)

X - матрица линия эксперимента.

Если обратим θ по МНК, то

$$\hat{\theta} = (X^T \cdot X)^{-1} \cdot X^T \cdot Y$$

$$X^T \cdot X \cdot \theta = X^T \cdot Y$$

(2)

Оп. Матрица линия X поддается однородной, если однородные стоят друг друга.

$\Rightarrow X^T \cdot X$ - диаг. матрица.

В этом случае (2) поддается на $n+1$ независимых ур-ий.

Потом при удалении каждого из ур-я все других коэффициентов остаются теми же самими \Rightarrow не нужно будет пересчитывать при удалении/добавлении коэффициентов.

Если $X^T \cdot X$ симм. к вспомогательной, то решение есть независимо. Это удобное свойство: коэффициенты решаются одновременно. \Rightarrow надо убрать кратные коэффициенты, чтобы её не было.

§11. Регрессионные переменные.

Когда есть категор. признаки, тогда дихотомиче (dummy) переменные.

Пример 1.

Учтываема зависимость Y в зав-ти от количественных колич. факторов X_1, \dots, X_p а еще от номиц/омкодных факторов образование.

$$Y_j = \beta_1 \cdot X_{j1} + \dots + \beta_m \cdot X_{jm} + \varepsilon_j$$

$$d_j = \begin{cases} 1, & \text{если есть фактор} \\ 0, & \text{в ином случае} \end{cases}$$

\Rightarrow расширяем регресс. модель

$$Y_j = \beta_1 \cdot X_{j1} + \dots + \beta_m \cdot X_{jm} + \delta \cdot d_j + \varepsilon_j$$

Но не всегда интересующий фактор принимает только 2 значения...

Пример 2.

Учтываема зависимость Y в форме сезона.

$$Y_j = \beta_1 \cdot X_{j1} + \dots + \beta_m \cdot X_{jp} + \varepsilon_j$$

\Rightarrow Всё же необходимо колич. переменных, а не один со значениями 0, 1, 2, 3, т.е. колич. Каждый из них один и тот же, меняя значение между сезонами будет независимое (сезон-зима $>$ весна-лето).

Нужно ли нам 4 сезона? Нет, если мы забываем оно 3-х.

Вважаємо, що у нас немає:

$$d_j^{(1)} = \begin{cases} 1, & \text{если } j \text{ має} \\ 0, & \text{без кирице} \end{cases}$$

$$d_j^{(2)} = \begin{cases} 1, & \text{если } j \text{ має} \\ 0, & \text{без кирице} \end{cases}$$

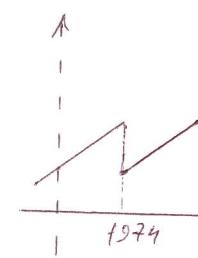
$$d_j^{(3)} = \begin{cases} 1, & \text{если } j \text{ має} \\ 0, & \text{без кирице} \end{cases}$$

$$Y_j = \beta_1 \cdot X_{j1} + \dots + \beta_p \cdot X_{jp} + \gamma_1 \cdot d_j^{(1)} + \gamma_2 \cdot d_j^{(2)} + \gamma_3 \cdot d_j^{(3)} + \varepsilon_j$$

Приклад 3

1973 - арабо-українська лінія

\Rightarrow 1974 ?:



\Rightarrow юдеї євр. можуть \Rightarrow вважаємо їх юдеями.

$$d_j = \begin{cases} 1, & t < 1974 \\ 0, & \text{інакше} \end{cases}$$

Іншім членам можна додати:

\Rightarrow юдеї євр. можна додати.

Приклад 4

Дізнатися σ залежності зі W в Y від віку age та статі S (чи $age \cdot S$).

Уточнення залежності зі W в Y від залежності від статі (фактора) та віку підтверджено.

age - чинник

S - чинник радянського

$$S_j = \begin{cases} 1, & \text{если } j \text{ має} \\ 0, & \text{без кирице} \end{cases}$$

Розглядаємо цілі числа:

$$W_j = \alpha + \beta \cdot age_j + \gamma \cdot S_j + \varepsilon_j$$

Уточнення МНК (коф. коеф. залежності)

coef	Est	SE	t	p-value
α	0.11	2.39	2.56	0.011
β	0.53	0.06	8.58	0.000...
γ	-3.73	1.35	-2.76	0.0065

$$\Rightarrow W = -0.11 + 0.53 \cdot age - 3.73 \cdot S$$

$$R^2 = 0.39$$

$$F^* = 0.000...$$

$$\boxed{\gamma = -3.73}$$

\Rightarrow зі W в Y в W залежність від статі підтверджена на 3.73 високою розподільованістю

Табл 6. Основы классической линейной модели.

§1. Основы регрессионных моделей

Модель огни залежаного переменного и m обозначающих факторов. Проверка N гипотез, которые связаны ссы. оценки:

$$Y_j = \theta_0 + \theta_1 X_{j1} + \dots + \theta_m X_{jm} + \varepsilon_j \quad (1)$$

$$Y = X \cdot \theta + \varepsilon \quad (1')$$

Основные априорные.

- 1) модель лин. по параметрам
- 2) X -неслучайная матрица
- 3) $\{\varepsilon_j\}$ - незав. (негод.)
- 4) $E(\varepsilon_j) = 0 \quad \forall j$ - нет систематических ошибок
- 5) $D(\varepsilon_j) = \sigma^2 \quad \forall j$ - гомогенность дисперсии оценки огни и оценок
- 6) $\{\varepsilon_j\}$ - несвр. норм. распред.

Если матрица матрица X невесн., то

$$\hat{\theta} = (X^T \cdot X)^{-1} \cdot X^T \cdot Y = \theta + (X^T \cdot X)^{-1} \cdot X^T \cdot \varepsilon$$

Оценка дис. σ^2 :

$$S^2 = \frac{1}{N-(m+1)} \cdot \sum_{j=1}^N e_j^2, \text{ где } \hat{Y} = X \cdot \hat{\theta}, e_j = Y_j - \hat{Y}_j.$$

Основные оценки.

- 1) $\hat{\theta}$ - мн. по Y .
- 2) $\hat{\theta}$ - не свр.
- 3) $\hat{\theta}$ - оптим. в ср. квадр. в классе мн. и несмн. оценок.
- 4) Если $\frac{1}{N} \cdot X^T \cdot X \rightarrow \Sigma_X, \quad N \rightarrow \infty$, и Σ_X - квадр. матр. порядка $(m+1)$, то $\hat{\theta}_n$ - сим.
- 5) S^2 - несмн. и сим. оценка для σ^2 .

§2. Модель со смешан. факторами

1. Описание модели.

Проверка N гипотез.

- 1) модель мн. по параметрам, т.е.

$$Y = X \cdot \theta + \varepsilon \quad (1)$$

- 2) X -свр. матрица, но $X_{j0} = 1 \quad \forall j$.

$$\text{и } E X_{jk} = 0 \quad \forall k = 1, m$$

- 3) Проверка: X и ε -независимы.

$$4) \mathbb{E}(\varepsilon_j) = \mathbb{E}(\varepsilon_j | X) = 0 \quad \forall j$$

$$5) \sum \varepsilon = \mathbb{E}(\varepsilon \cdot \varepsilon^T) = \mathbb{E}(\varepsilon \cdot \varepsilon^T | X) = \sigma^2 \cdot E.$$

6) cl. лемма ($y_j, x_{j1}, \dots x_{jm}$) имеет одно и то же многомерное норм. вероятн. распред.

Если выполнено первое условие выше означается, то в силу 35Ч имеет:

$$\frac{1}{N} \varepsilon^T \cdot \varepsilon = \frac{1}{N} \sum_{j=1}^n \varepsilon_j^2 \rightarrow \sigma^2 \quad (2)$$

$$\frac{1}{N} X^T \cdot X \rightarrow \Sigma_{xx} \quad \text{cov}(X, X) \quad (3)$$

$$\frac{1}{N} X^T \cdot \varepsilon \rightarrow 0 \quad \text{недостаточный анализ взаимной зависимости и симметрии} \quad (4)$$

при $N \rightarrow \infty$.

2. Оценки параметров и их свойства.

Применим к (1) МНК. \Rightarrow

$$\hat{\theta} = (X^T \cdot X)^{-1} \cdot X^T \cdot Y = \theta + (X^T \cdot X)^{-1} \cdot X^T \cdot \varepsilon \quad (5)$$

Одозначим $\hat{Y} = X \cdot \hat{\theta}$, $e = Y - \hat{Y} \Rightarrow$

$$\hat{\sigma}^2 = S^2 = \frac{1}{N-(m+1)} \cdot \sum_{j=1}^n e_j^2 \quad (6)$$

Пусть даны основные опр. 1)-6).

Свойства оценок.

1) $\hat{\theta}$ -инвариантна по Y

2) $\hat{\theta}$ -нелиней. оценка θ . (получается через линейное преобразование $\theta = \hat{\theta} + \varepsilon$ (лвз))

3) $\hat{\theta} = \hat{\theta}_n$ -состоит из независимых

$$\frac{\partial \hat{\theta}}{\partial \theta} = \frac{\partial \theta}{\partial \theta} + (X^T \cdot X)^{-1} \cdot X^T \cdot \varepsilon = \theta + \left(\frac{1}{N} X^T \cdot X \right)^{-1} \cdot \left(\frac{1}{N} X^T \cdot \varepsilon \right) \Rightarrow \theta + \Sigma_{xx}^{-1} \cdot 0 = \theta$$

Важный момент: для этого это, что матрица не коррелирует с остатками. Если это не так, то $\hat{\theta}$ -то состоятельность нарушается. \blacksquare

4) $\hat{\theta}$ -оценивается в среднем в виде лин. и нелиней. оценок.

Д-бо: начали не-лн. гл. функции си. функ. X , а потом упростили по X .

5) S^2 -стоит нелиней. и лин. оценка си. функ. σ^2 .

23.10.

(2)

(3)

(4)

(5)

(6)

3. Метод исклучающих переменных

Если некот. генерал. корр. с ошибками, то на него
зависимость на другое генерал. изменение не будет корр-но.
Пусть имеем некот. модель $\mathcal{Z} \in \text{матрица } N \times (m+1)$,
 $Z_{j0} \equiv 1$, при этом

$$\frac{1}{N} Z^T \cdot \varepsilon \rightarrow 0 \quad (7)$$

$$u \quad \frac{1}{N} Z^T \cdot X \xrightarrow{P} \Sigma_{zx}, \quad (8)$$

где Σ_{zx} — некорр. квадр. матрица.

Зависимости некорр. оценки

$$\tilde{\theta} = (Z^T \cdot X)^{-1} \cdot Z^T \cdot Y = \theta + (Z^T \cdot X)^{-1} \cdot Z^T \cdot \varepsilon = \theta + \left(\frac{1}{N} Z^T \cdot X \right)^{-1} \cdot \left(\frac{1}{N} Z^T \cdot \varepsilon \right) \rightarrow \\ \rightarrow \theta + \Sigma_{zx}^{-1} \cdot 0 = \theta.$$

Конечное выражение в \mathcal{Z} — исклучающие переменные.

Экспонометрическое выражение в виде матрицы Σ_{zx} называется экзогенное выражение.

§3. Особенности МНК.

В классич. ин. модели:

$$Y = X \cdot \theta + \varepsilon \quad (1)$$

при этом

$$\Sigma_{yy} = \Sigma_{\varepsilon\varepsilon} = \sigma^2 \cdot E \quad (2)$$

Пусть имеем

$$\Sigma_{\varepsilon\varepsilon} = \sigma^2 \cdot \Omega, \quad (3)$$

где Ω — целесообразная статист. положим. оцен. матрица.

Применим общий МНК. \Rightarrow

$$\hat{\theta} = (X^T \cdot X)^{-1} \cdot Y = \theta + (X^T \cdot X)^{-1} \cdot X^T \cdot \varepsilon \Rightarrow$$

1) $\hat{\theta}$ — ин. но Y .

2) $\hat{\theta}$ — несинг.

3) при них не удастся $\hat{\theta}$ -состоит.

4) сб-бо оптимальности не выполняется.

Чт. индейкой видим, что \exists некорр. матрица P :

$$\Omega = P \cdot P^T \quad (4)$$

Перейдём к новым переменным: (ჩე үзүүмдөн нө P^{-1})

$$Y_* = P^{-1} \cdot Y, \quad X_* = P^{-1} \cdot X, \quad \varepsilon_* = P^{-1} \cdot \varepsilon \Rightarrow$$

$$Y_* = X_* \cdot \theta + \varepsilon_* \quad (5)$$

Показано, что (5) — это классич. ин. модель.

$$\sum_{\varepsilon^* \varepsilon^*} = \mathbb{E}(\varepsilon_* \varepsilon_*^\top) = \mathbb{E}(P^{-1} \cdot \varepsilon \cdot \varepsilon^\top (P^{-1})^\top) = P^{-1} \cdot \mathbb{E}(\varepsilon \cdot \varepsilon^\top) \cdot (P^{-1})^\top = \\ = \sigma^2 \cdot P^{-1} \cdot \Sigma \cdot (P^{-1})^\top \stackrel{(4)}{=} \sigma^2 \cdot P^{-1} \cdot P \cdot P^\top (P^\top)^{-1} = \sigma^2 \cdot E$$

Применение к (5) общему МНК. Понятие оценки для уравнения регрессии.

$$(Y_* - X_* \cdot \theta)^\top \cdot (Y_* - X_* \cdot \theta) = (Y - X \cdot \theta)^\top \cdot \Sigma^{-1} (Y - X \cdot \theta) \rightarrow \min_{\theta} - \quad (6)$$

- σ^2 -однозначный МНК.

$$\hat{\theta}_* = (X_*^\top X_*)^{-1} \cdot X_*^\top Y_* = (X^\top \Sigma^{-1} X)^{-1} \cdot X^\top \Sigma^{-1} Y - \text{оценка для } \theta$$

Вашим гаиниң айрым:

Түбөн Σ -дағы матрица с элементами w_1, w_2, \dots, w_N .
И.е. симметрическую, но, возможно, неявно разные значения
и их знаек соотношение между их значениями.

Σ^{-1} -дағы матрица с элементами $\frac{1}{w_1}, \dots, \frac{1}{w_N}$.

$$(6) \Rightarrow \sum_{j=1}^N [Y_j - (X \theta)_j]^2 / w_j \rightarrow \min_{\theta} - \text{безвешенний МНК.} \quad (7)$$

Табл 7. Температуростатичность и аномальность

$$Y = X \cdot \theta + \varepsilon \quad (1)$$

$$\text{В квадр. ин. можем } \sum_{\varepsilon\varepsilon} = \sigma^2 \cdot E \quad (2)$$

$$\text{Мн. расч.-и аргам, когда } \sum_{\varepsilon\varepsilon} = \sigma^2 \cdot S_2, \text{ где } S_2 \text{ известна.} \quad (3)$$

Если S_2 неизвестна, то имеем $\frac{N(N+1)}{2}$ дополнительных измерений, кроме θ и σ^2 .
т. измерений буде N . (!)

2 измерения в $\sum_{\varepsilon\varepsilon}$: когда измерим разницу (на разн. не влияет величина) и квадр-и: где разн не нужна \Rightarrow б-2 измерения на разности.

§ 1. Температуростатичность. (когда ошибки несогр по средн. цен.)

1. Описание модели

Базис. модель

$$Y = X \cdot \theta + \varepsilon \quad (1)$$

$$\sum_{\varepsilon\varepsilon} = \sigma^2 \cdot S_2,$$

где S_2 -квадр. коварианца.

Пусть S_2 -неизвестна!

Применим одинаковую МНК. Получим оценки

1) ин.

2) нелинейн.

3) симметричн.

4) не описываемы!

Возьмем N доп. измер., кроме θ и σ^2 \Rightarrow останутся уравнения не линейные, но все решения входят в класс симметричн. измерений. Модель описываема выше табл. не-се моделью ин. пересечения с температуристичностью в модели измерений.

2. модель температуристичности

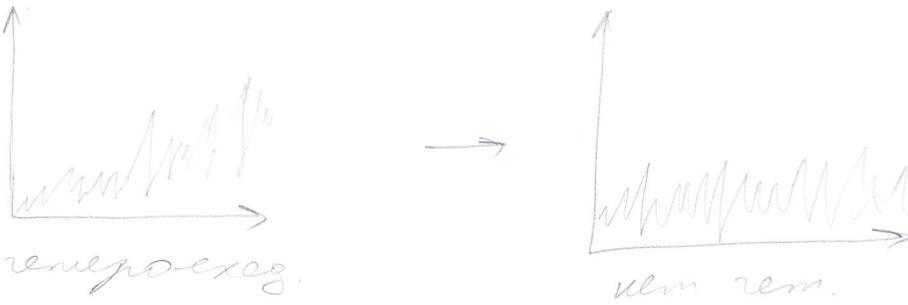
Различия между моделями, в которых ошибки описываются различными измерениями, т. е.

$$\hat{\sigma}_j^2 = E(\hat{\varepsilon}_j) = \sigma^2 \cdot |X_{je}|^2 \quad (3)$$

различны, это очевидно, т.е. применение различных МНК. Рассмотрим оценку измерений следующим образом: переходом к новым переменным:

$$Y_j^* = Y_j / |X_{je}|, X_{jk}^* = X_{jk} / |X_{je}|, \text{ а также применением одинаковой МНК. 28}$$

- Следует отметить, что модель имеет одинаковую форму, что модели имеют одинаковую форму (3): Не априорно, просто чтобы определить эту модель нет необходимости знать её форму.
- 1) Применение основной МНК к набору данных, где не учитывается нелинейность времени $|X_{j1}|$.
 - 2) Наборы остатков e_j и сплошные гладкие $|e_j|$.
 - 3) Если на графике логика предполагается линейная форма времени $|e_j|$, это в пользу (3).
 - 4) Переход от линии непрерывности, где на $|X_{j1}|$. Стока определяем (1) по МНК, наборы остатков e_j и сплошные гладкие $|e_j|$.
 - 5) Если на графике $|e_j|$ логика хаотического ходят (и.е. нет линейной тенденции), то это против в пользу модели (3).



Другие модели.

$$\hat{\sigma}_j^2 = \gamma_0 + \gamma_1 Z_{j1} + \dots + \gamma_p Z_{jp} \quad (4)$$

Предполагается следующая форма для оценки моделей:

- 1) Определяем (1) по основной МНК и наборам остатков e_j . Остаток данных себе время несёт на оценку.
- 2) Используя остатки e_j , определяем фундаментальное пересечение следующего логика:

$$e_j^2 = \hat{\gamma}_0 + \hat{\gamma}_1 Z_{j1} + \dots + \hat{\gamma}_p Z_{jp} + u_j \quad (5)$$

Далее решим

$$\hat{\gamma}_j^2 = \hat{\gamma}_0 + \hat{\gamma}_1 Z_{j1} + \dots + \hat{\gamma}_p Z_{jp} \quad (6)$$

- 3) Применение к (1) фундаментальной МНК с формулой $\hat{\sigma}_j^2$.
 - 4) Но чтобы остатки определялись гладкими оценками $\hat{\theta}$.
- Двухстадийный МНК.

Тесты на генетическая зависимость.

Нулевая гипотеза имеет вид:

$$H_0: \sigma_1^2 = \sigma_2^2 = \dots = \sigma_N^2 = \sigma^2$$

H_1 : это не так

Метод Гандрича - Куангма

$$H_1: \sigma_j^2 = \sigma^2 \cdot X_{je}^2$$

Ходим проверить гипотезу на возрастание для этого теста X_{je} по возрастанию.

Описание процедуры:

- 1) Все независимые упоминаются по возрастанию независимой $|X_{je}|$.
- 2) У тех данных удалены среднее и измерение (однако $d \sim N/2$ - назнача рекомендации).
- 3) Стартует 2 независимые регрессии по первым $(N-d)/2$ измерениям и последним $(N-d)/2$ измерениям и находят и мал, и мал величина становится e_1 и e_2 .
- 4) Но останется e_1 и ее критическое сумма называемая ESS_1 , а ESS_2 и статистика $F = \frac{ESS_2}{ESS_1}$.
- 5) Если верна H_0 , то статистика F имеет распределение Снедгоуза - Римера с $\left(\frac{N-d}{2} - p, \frac{N-d}{2} - p\right)$ степ. добр.
- 6) Далее процедура повторяется с помощью теста H_0 , формируется p -value, ...

Если обнаружена генетическая зависимость, то измерениях взаимно МНК с весами $\sigma^2 X_{je}^2$ (из нулевой H_1).

§2. Модели с аномальностью в ошибках.1. Описание модели.

Запишем математическую модель измерений:

$$Y = X \cdot \Theta + \varepsilon \quad (1)$$

Модель измерений, что ошибки измерений удовлетворяют квад. соотношению:

$$\varepsilon_j = \rho \cdot \varepsilon_{j-1} + \delta_j, \quad -\text{модель аномальности}$$

$$1) |\rho| < 1$$

$$2) \delta_j \sim N(0, \sigma^2)$$

$$(2) -\text{модель аномальности 1-го порядка}$$

(2)

(кв. наложение биах-ся трех ед. все на высшем)

Таким образом, ожидается, что набор $\{\varepsilon_j\}$ не имеет норм. распредел., т.к. же, что $E(\varepsilon_0) = 0$, $D(\varepsilon_0) = \frac{\sigma^2}{1-p^2} = \sigma_\varepsilon^2$.
 Тогда $E(\varepsilon_j) = 0$, $D(\varepsilon_j) = \sigma_\varepsilon^2 \quad \forall j$
 $Cov(\varepsilon_j, \varepsilon_{j+m}) = \sigma_\varepsilon^2 \cdot p^m$
 $\Rightarrow \{\varepsilon_j\}$ - стационарная в широком смысле сущ. корр. врем. Структурная матрица корреляции:

$$\Sigma_\varepsilon = \sigma_\varepsilon^2 \begin{pmatrix} 1 & p & p^2 & \dots & p^{N-1} \\ p & 1 & p & \dots & p^{N-2} \\ \dots & \dots & \dots & \dots & \dots \end{pmatrix}$$

Если применить основной МНК, то оценки для θ будут несущественны, но не будут оптимальными.
 Тогда p известно. Рассмотрим только модели:

$$Y_j^* = Y_j - pY_{j-1} = \theta_0(1-p) + \theta_1(X_{j1} - X_{j-1,1}p) + \dots + \theta_m(X_{jm} - X_{j-1,m}p) + \varepsilon_j \quad (4)$$

(4) - аутокоррессионная модель. иллюстрирует 1-м наблюдением

2. Оценивание модели с аутокоррессией.

Процедура Кохрена-Дракманта:

- 1) Применение основной МНК к модели (1) и нахождение остатков e .
- 2) Рассмотрим модель регрессии вида:
 $e_j = pe_{j-1} + \eta_j, \quad j \geq 2$
- \Rightarrow находим оценку p по МНК
- 3) Поставляем p вместо p в (4). Оцениваем модель (4) по основному МНК.
- 4) Находим остатки белых остатков $\tilde{e} = y - \hat{x}\hat{\theta}$.
- 5) Повторяем 2)-4), пока не будем достичь необходимой мерности остатков p .

Процедура Хиддемана-Лу

- 1) Разбиваем $[-1; 1]$ на некую сетку
- 2) Для каждого значения p проводим аутокоррессионное преобразование.
- 3) Оцениваем по МНК преобр. модель
- 4) Выбираем то значение p , где одна пакетная сумма квадратов остатков.
- 5) В зависимости найденного p строим более некую сетку и повторяют процедуру пока не достичь необходимой мерности.

Теорема для проверки, есть ли автокорреляция:

3. Критерий Дарбина-Уотсона.

"Дарбин Ватсон"

Рассмотрим линейную модель:

$$Y = X \cdot \theta + \varepsilon$$

(1)

Ошибки модели:

$$\varepsilon_j = \rho \cdot \varepsilon_{j-1} + \delta_j$$

(2)

$H_0: \rho = 0$ - нет автокорр.

$H_1: \rho > 0$.

Оценка (1) не однозначна МНК.

$\Rightarrow \hat{\theta}$,

$$\hat{Y} = X \cdot \hat{\theta},$$
$$e = Y - \hat{Y}.$$

Рассмотрим статистику DW:

$$DW = \frac{\sum_{j=2}^n (e_j - e_{j-1})^2}{\sum_{j=1}^n e_j^2} \approx 2 \cdot (1-\rho), \text{ где } \rho - \text{коэффициент корреляции}$$

Так $H_0: DW \approx 2$

$H_1: DW < 2$.

Численным методом удалось доказать, что значение этой статистики зависит от коэффициентов, входящих в модель.

Соответственно, что модель должна содержать неизвестные
 \Rightarrow сгенерировано первое правило!

Для заданных d , N , m существует такое значение:

$d_{\text{крит}}$ \uparrow \uparrow \uparrow
уров. значимости \rightarrow количество генерирований \rightarrow количество генерирований

$d_{\text{крит}} < d_r$, то $d_e \leq d(d)$ $\leq d_r$,

где $d(d)$ - квантиль критической констатанты.

$d_e \leq d_r$ - констатанта, не зависящая от генерирований, так что мы можем оценить $d(d)$.

После:

1) Если $DW < d_e$, то принимаем H_1

2) Если $DW > d_r$, то H_0 не противоречит экспериментальным данным (автокорреляцию не обнаружено)

3) $d_e \leq d(d) \leq d_r \Rightarrow$ решение нет

Типичен

$$N=100$$

$$m=1$$

$$\lambda=0,05$$

$$de = 1,65$$

$$dr = 1,69.$$

Теория 8. Модели дискретного выбора.

§1. Бинарные модели выбора

1. Описание модели.

Делаем выбор - о чём?

Y - признакое решение.

X = (X₁, ..., X_p) - факторы.

Основная задача: найти $P(Y=1|X)$. $P(Y=0|X) = 1 - P(Y=1|X)$.

Предположим, что

$$P(Y=1|X) = F(X \cdot \theta) = F(\theta_1 \cdot X_1 + \dots + \theta_p \cdot X_p),$$

где $F(z)$ - неубывающая функция со значением из $[0, 1]$.

Давно предполагалось, что $F(z)$ - функция распределения.

"То есть что? Я прошу. Может максимум? Давно я спросил
они говорят, почему это предположение?"

Предположим, что \exists некоторое непрерывное (непрерывная)
переменное Y^* , которое подчиняется следующей модели:

$$Y^* = X \cdot \theta + \varepsilon,$$

ε -с.в. с д.п. $F(z)$, которой соответствует ожидание.

Тогда $Y = \begin{cases} 1, & \text{если } Y^* \geq 0, \\ 0, & \text{если } Y^* < 0. \end{cases}$ Согласно (2)

$$\Rightarrow P(Y=1|X) = P(Y^* = X \cdot \theta + \varepsilon \geq 0 | X) = P(\varepsilon \geq -X \cdot \theta) = 1 - F(-X \cdot \theta) = F(X \cdot \theta)$$



Другой подход к обоснованию:

Такое описание нашего решения имеет наследство
ко всем u_i , причём

$$u_0 = X \cdot \theta^{(0)} + \varepsilon_0$$

$$u_1 = X \cdot \theta^{(1)} + \varepsilon_1$$

ε_0 и ε_1 - случайные независимые факторы.

Естественно принимать то решение, Y которого наследует
форму. \Rightarrow

$$P(Y=1|X) = P(u_1 - u_0 \geq 0 | X) = P(\varepsilon_1 - \varepsilon_0 \geq -X \cdot (\theta^{(1)} - \theta^{(0)}))$$

Обозначим $\theta^{(1)} - \theta^{(0)} = \theta$, $\varepsilon_1 - \varepsilon_0 = \varepsilon$ -] имеем д.п. $F(z)$, которой
соотв. ожидание. \Rightarrow получаем модель (1).

2. Примеры моделей бинарного выбора

1) Ещё $F(z) = \Phi(z)$ - д.п. стандартн. норм. закона, то имеем
Probit-модель.

Не стандартное, но аналогичное описание к этому курсу.

2) Если $F(z) = \Lambda(z)$ - ф.р. смеш. концептуального распределения, то именем Logit-модель.

$$\Lambda(z) = \frac{e^z}{1+e^z} = 1 - \frac{1}{1+e^z}$$



3. Оценка модели.

МНК здесь неприменимо, оно и есть однозначность, не знаем.

В этой задаче все оценки параметров稱 называют локальной линией наилучшего приближения:

Также имеем N независимых наблюдений

$y_i, x_{i1}, \dots, x_{ip}$.

Запишем форму найлучшего приближения:

$$L(\theta) = \prod_{j=1}^N [F(x_j; \theta)]^{y_j} [1 - F(x_j; \theta)]^{1-y_j} \quad (3)$$

$$l(\theta) = \log L(\theta) = \sum_{j=1}^N (y_j \log F(x_j; \theta) + (1-y_j) \log (1 - F(x_j; \theta))) \rightarrow \sup_{\theta} \quad (4)$$

Процессом анализа показываем, что в случае Logit-модели ф-я $l(\theta)$ является нестатистически функцией.

\Rightarrow У этой ф-и есть $\exists!$ максимум.

Если ф-я "хорошая", а эта хорошая \Rightarrow асимптотическая норм. оценки.

4. Виды моделей.

Если где модель:

- одноблочная с параметром θ_0

- автоматическая, которая имеет некомп. число параметров (гиперпараметров) $\rightarrow \theta_A$.

Наглядно это не сильно.

McFadden (McFadden) критерий качества логит-модели генетической (аналог критерия генетической):

$$R_M^2 := 1 - \frac{l(\hat{\theta}_A)}{l(\hat{\theta}_0)} \quad (5)$$

Когда $R_M^2 = 1$, то есть лучше $\hat{\theta}_A$ чем $\hat{\theta}_0$ и именем одноблочной.

Хотим чтобы предсказание лежало между 0 и 1, или лучше одноблочной.

5. Задача оценки.

Все наблюдения разделяем на 2 части раздела N_1 и N_2 :

$$N_1 + N_2 = N.$$

По разделам N_1 измер. оценивают модель и находят оценку $\hat{\theta}$.

Далее получим

$$\hat{Y}_j = \begin{cases} 1, & \text{если } F(x_j; \hat{\theta}) \geq c \\ 0, & \text{в ином случае} \end{cases}, \quad j = \overline{N_1 + 1, N_1 + N_2}$$

При этом $c = \frac{1}{2}$, если хотим, чтобы оценка оказалась на нейтральном уровне, то иначе.

Для оценки качества упаковки спроси:
Моделью предсказаний и реализаций.

		Предсказание	
		$\hat{Y}_j = 1$	$\hat{Y}_j = 0$
Показат.	$Y_j = 1$	P_{11}	P_{10}
	$Y_j = 0$	P_{01}	P_{00}
		$P_{.1}$	$P_{.0}$

Был 1, предсказан 0

Случайно не упаковка

Качество упаковки можно оценить по демпфивности предсказанных упаковок = $P_{11} + P_{00}$.

Мера Maxgaggena: $F_1 = \frac{P_{11} + P_{00} - P_{11}^2 - P_{00}^2}{1 - P_{11}^2 - P_{00}^2}$

6. Виды между днями покупки товара

На рынке есть 2 сорта консервов: Heinz и Hunts.

Составлено данные по 2798 покупкам для 300 семей.

Для каждого семейства отмечено последнее покупку ($N_2 = 300$).

На этом основе было построено распределение карамба.

Оценивают модели по 2498 измерениям. В качестве временных признаков используют следующие переменные:

1) лог отношение цен:

$$\ln p_{10} = \ln \left(\frac{\text{price Heinz}}{\text{price Hunts}} \right)$$

1-Heinz

0-Hunts

2) X_{11} - дни в неделю супермаркет Heinz наименее.

(последние недели 10 или 1, остальное)

X_{12} - дни в неделю Heinz в первом

X_{13} - дни в неделе второго консерва одновременно.

X_{21} | дни в неделю для Hunts.

X_{22}

X_{23}

Оценки Logit-моделей в пакетах:

Перем	Коэффи	...	p-value	X_{12} и X_{13} - незначим.
C	1.85		0.0000	X_{22} и X_{23} - значим.
X_{11}	0.27		0.0357	Почему? Кого сюда хотят сортировать по X_{11} \Rightarrow оно - а само X_{11} в него не попадает \Rightarrow оно не является критерий для X_{11} .
X_{12}	0.19		0.2334	необходимо раз $>$ критерий Хи-квадрат \Rightarrow более распределение монотонное и не линейное \Rightarrow критерий Хи-квадрат не применим.
X_{13}	0.25		0.3051	
X_{21}	-0.37		0.0125	
X_{22}	-0.57		0.0036	
X_{23}	-1.09		0.0001	
$Z = \hat{\rho}_{10}$	-3.27		0.0000	

Если вписали модели предсказанных / подтверждаемых значений, то получим:

коэф. предсказаний = 0.77. \Rightarrow хорошие предсказания. У экзаменов и обследование лучше 40% симптомов удачно предсказанных.

математическое выражение

§2. Многомерные модели неподтверждаемого метода.

13.11

Как раз модель от Маргаггена.

1. Описание модели \exists несколько моделей, с одинак.

Нужно необходимо выбрать одно из нескольких решений:

1, 2, ..., J.

Чтобы необходимо обозначить значения переменных:

X_1, X_2, \dots, X_m .

Нужно приложить решения к всем полезностям u_{ij} для i -го индивидуума, причём

$$u_{ij} = X_i \cdot \theta_j + \varepsilon_{ij} = \bar{u}_{ij} + \varepsilon_{ij} \quad (1)$$

$$\text{здесь } X_i = (X_{i1}, \dots, X_{im}) = (X_1, \dots, X_m) \quad \text{и} \quad \theta_j = (\theta_{j1}, \dots, \theta_{jm})$$

ε_{ij} - влияние неприм. факторов, супр. фак.

Прикладываемые по решению, у которого величина самая большая.

$$P(Y_i=j | X_i) = P(u_{ij} \geq u_{ik} \ \forall k | X_i) = P(u_{ij} = \max_k u_{ik} | X_i) \quad (2)$$

Кого как-то конкретизировано назначение X_i .

В ТБ есть описание модели в виде машины (мн)

Таким же E_{ij} - и.о. р.с.в., имеющие равноделение ожиданий:

$$P(E_{ij} < z) = e^{-e^{-z}}, \quad z \in \mathbb{R}.$$

Максимизируя наше, мы в максимуме модели

$$P(Y_i = j | X_i) = \frac{e^{x_{ij}\theta_j}}{\sum_k e^{x_{ik}\theta_k}} = \frac{e^{x_{ij}\theta_j}}{\sum_k e^{x_{ik}\theta_k}} \quad - \quad (3)$$

- многомерная Logit-модель

Почему Logit, если у нас не Logit-распределение? Потому что это не линейная модель, но неприведенное значение не линейно, как в Logit-модели.

Замечание.

Равноделение в решении: $j \neq l$ и в противном

~~$$P(Y_i = j | X_i) = \frac{e^{x_{ij}\theta_j}}{e^{x_{il}\theta_l}}$$~~
$$\frac{P(Y_i = j | X_i)}{P(Y_i = l | X_i)} = \frac{e^{x_{ij}\theta_j}}{e^{x_{il}\theta_l}} \quad - \quad (4)$$

- независимость от посторонних аномалий

Есть примеры, которые показывают, что с экспоненциальной моделью решений это независимо (как мы знаем, что на рынке наличие 3-и аномалий в Сингапуре С.Лишилл не уверен на это определение?) \Rightarrow модели не очень хороши

Хорош в то-примерах моделей вроде наезд, где ошибка залива и вероятность тоже есть в этом виде.

2. Оценивание моделей.

Пусть имеем N независимых одновременных измерений y_i, x_i .
Причина для этого называется МНП (меньшее квадратов остатков). Задача оценки параметров модели:

$$L(\theta) = \prod_{i=1}^N \prod_{j=1}^J [P(Y_i=j | X_i)]^{\mathbb{1}(Y_i=j)} \quad (5)$$

Перенесём в логарифмы:

$$l(\theta) = \sum_{i=1}^N \sum_{j=1}^J \mathbb{1}(Y_i=j) \cdot \log P(Y_i=j | X_i) \rightarrow \max_{\theta} \quad (6)$$

Данный задачу на экстремум можно решить методом, в этом духе выше выше.

Для реш. модели $l(\theta)$ скажем правило Канделл \Rightarrow ед. максимум max.

Но в этом случае эта модель не реализована (ибо в логарифме близких... \Rightarrow если не хотим вспоминать, пишем как :)

3. Бином одно из четырех сортов молока - пример

Есть 4 сорта молока на рынке: 1, 2, 3, 4.

Сейчас гамма да $N=3292$ покупок в 136 супах
Что является (гипотеза):

- один купленного молока A₁
- один из других молок (из 3-x не поддаются) A₂
- какое-либо покупка \rightarrow различные переключения:
 - d₁ - покупка молока B₁
 - d₂ - покупка в супе B₂
 - d₃ - то и другие виды B₃

Была проверена гипотеза о том, что можно предположить, что молоко купленное молоко - регул. зависимое перек. \rightarrow знания.
Каждый покупки независимы, но есть (3) противоречия между теми что и что молоко.

Составляем модельную передачу и реальных знаний, как-то правила оценки по % правильно передаваемых знаний.

В данных задаче для молока передаваемые знания = 58%.

§3. Многомерные модели прогнозирования будущего.

1. Описание модели

Вариант $1 < 2 < 3 < \dots < T$. \rightarrow оценки на будущее, ($2 < 3 < 4 < 5$)
 \nwarrow будущее

Будет простая ситуация, что в §2. Решение с помощью замененной переменной.

Также имеем некоторую зависимость (неподдающуюся) переменную y_i^* , которую называем моделью:

$$y_i^* = x_i \cdot \theta + \varepsilon_i \quad (1)$$

$x_i = (x_{i1}, \dots, x_{im})$ - числ. признаков,

$\theta = (\theta_1, \dots, \theta_m)$ - набор параметров.

Предполагаем, что ε_i - н.о.р.с.в с ф.п. $F(z) = P(\varepsilon_i < z)$.

Также будем называть порогом:

~~называем как $x_{i1} < a_1 & \dots & x_{im} < a_m$~~ $-\infty < a_0 < a_1 < \dots < a_{k-1} < a_k < \dots < a_T = +\infty$

Предполагаем, что

$$\boxed{Y_i = k | X_i} \Leftrightarrow a_{k-1} \leq y_i^* < a_k \quad (2)$$

$$\Rightarrow P(Y_i = k | X_i) = F(a_k - x_i \cdot \theta) - F(a_{k-1} - x_i \cdot \theta)$$

Для этого, мы должны заложить описание модели, надо б. языке будущего F .

Есть несколько вариантов:

- 1) Logit - модель
- 2) Probit - модель
- 3) модели для дискретных переменных

Параметры:

$$\alpha_k, \theta, \sigma$$

Часто бывает, что модель не линейна или не имеет линейных предположений. (если $\beta_0 = \text{нек-ая} * \text{const}$) \Rightarrow симметричный $\sigma = 1$.

Если в пар-ре $\beta_0 = \text{const}$, то коэффициенты β_i из $\hat{\alpha}_k$ \Rightarrow модель не линейна и имеет избыток. \Rightarrow condition $\alpha_k = 0$ или $\alpha_k = 1$.

Пусть $\alpha_k = 0$.

Также для оценки модели применяется МНГ. Здесь и в Logit - модели ответ в явном виде получать не удобно, поэтому решается мк. методом (так же в бинарной форме оценки есть.)

2. Пример.

Одна из тем тематики в определении риска здоровья.

Вопрос: можно ли выделить факторы, а не на родите?

Воп. отважен:

- 1) коммерческий не сориентирован
- 2) не сориентирован
- 3) сориентирован
- 4) полностью сориентирован

Обследование проводено среди 3705 семей с детьми из СМУ.

Y-ответ (1, 2, 3, 4).

Факторы:

- 1) возраст
- 2) физ. переменная = 1, если $g > 14$ лет сестра рождалась
- 3) физ. переменная = 1, если мама родом из юга России $g < 14$ лет
- 4) число детей, замужество матери на момент
- 5) доход семьи
- 6) число братов и сестер
- 7) prob = 1, если промышленное производство
- 8) cath = 1, если католическое производство
- 9) $south$ = 1, если из южных регионов
- 10) urb = 1, если живут в крупном городе

Оцениваются по МНГ в лог. и другой модел.

8-10 граммах ощущение неподвижности. Оказывается, что да и да
известно (и.e. ответы "известно" и "неизвестно" не имеют
различия, надо засчитывать на "не знаю").
Однако О, верши решения, супервизорами.

Теория дисперсионного анализа

Первый этап расчленения Динес.

Классическая литература: Медведь, "Дисперсионный анализ".

§1. Однодimensionalный дисперсионный анализ

Член следующего модель выражения:

$$Y_{kj} = \mu_k + \varepsilon_{kj} \quad (1)$$

k -образное различие, $k = \overline{1, m}$

j -номер измерения, $j = \overline{1, N}$

$$n = m \cdot N$$

ε_{kj} - с.б., где ненулевое среднеквадратичное:

- 1) $\varepsilon_{kj} \sim N(0, \sigma^2) \quad \forall k, j$
- 2) ε_{kj} - независимы для $\forall k, j$.

Внешним образом сдвиги среднее:

$$\mu = \frac{1}{m} \sum_{k=1}^m \mu_k$$

и тогда получим измерение:

$$z_k = \mu_k - \mu.$$

$$\Rightarrow Y_{kj} = \mu + z_k + \varepsilon_{kj} \quad (2)$$

Основная задача:

Выяснить, есть ли различия в средних для различных измерений измерений.

Формально изображено такими же

$$H_0: z_k = 0 \quad \forall k$$

нульные гипотезы

$$H_1: \exists k, \text{ где ненулевое } z_k \neq 0.$$

Пример. Y -измерение приходится в ходе 4 приемов пятью разами.

Причины МНК, т.е. расчленение следующим экспериментальным задачу:

$$Q = \sum_k \sum_j (Y_{kj} - \mu - z_k)^2 \rightarrow \min_{\mu, z_k} \quad (3)$$

$$\Rightarrow \hat{\mu} = \bar{Y}_.. = \frac{1}{n} \sum_k \sum_j Y_{kj} \quad (4)$$

$$\hat{Y}_k = \bar{Y}_{k.} = \frac{1}{N} \sum_j Y_{kj} \quad (5)$$

Справедливо следующее соотношение

$$SS_T = SS_B + SS_R \quad (6)$$

Total	Between	Residual
-------	---------	----------

$$SS_T = \sum_k \sum_j (Y_{kj} - \bar{Y}_{..})^2 \quad \text{- полная сумма квадратов}$$

$$SS_B = \sum_{k=1}^m (\bar{Y}_{k.} - \bar{Y}_{..})^2 \cdot N \quad \text{- разброс между группами}$$

$$SS_R = \sum_k \sum_j (Y_{kj} - \bar{Y}_{k.})^2 \quad \text{- разброс внутри групп}$$

Таким образом $H_0: SS_B$ должно быть мало по сравнению с SS_R .

Теорема Таким образом H_0

1) т.ч. SS_B и SS_R независимы и имеют χ^2 -распределение с $(m-1)$ и $n-m = (N-1) \cdot m$ степенями свободы

2) т.ч. $F = \frac{SS_B / (m-1)}{SS_R / (n-m)}$ имеет распределение Снедекор-Фишера

с $((m-1), (n-m))$ степенями свободы

Далее проверка гипотезы проходит по статистической схеме (сравниваем значение статистики с крит. const, ...)

Другая постановка задачи, в которой

$$SS_R \sim \chi^2 \sim N(0, \sigma^2)$$

Носят задача о проверке гипотезы об отсутствии к

$$H_0: \sigma^2 = 0 \quad \text{нуль}$$

$$H_1: \sigma^2 > 0.$$

Пример. Используя нормальную форму групп на одинаковом априори для 3-х групп: смесицейкой, спирином, эпиреном:

$$Y_1, Y_2, Y_3$$

$$m = 3$$

$$N = 9$$

Получим след. результаты для средних

$$\bar{Y}_1 = 9.78, \quad \bar{Y}_2 = 18.93, \quad \bar{Y}_3 = 13.1$$

$$\bar{Y}_{..} = 13.93$$

$$SS_B = 128.73, \quad SS_R = 496.62$$

Крит. F имеет расп. Снедекор-Фишера с $(2, 24)$ степ. свобод.

$$F^* = 3.11$$

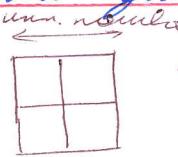
Но получено достаточно убедительное значение p-value:

$$P\text{-value}_F = 0.063$$

\Rightarrow Данные на 5% уровне не значимо различия.

§2. Двухфакторный дисперсионный анализ

Причина: $m_{k,j}$
↑
средний



rest-to средний

(Даны 2 главные переменные, две которых не находятся в взаимодействии.)

Использование схемы измерений:

$$Y_{kj} = \mu_{kj} + \varepsilon_{kj}$$

(7)

k -уровень первого фактора

j -уровень второго фактора.

Предположения:

$$1) \varepsilon_{kj} - \text{независимые}$$

$$2) \varepsilon_{kj} \sim N(0, \sigma^2) \quad \forall k, j$$

$$k = \overline{1, m}, \quad j = \overline{1, N}, \quad n = m \cdot N$$

Возможны следующие обстоятельства:

$$\mu_{..} = \frac{1}{n} \sum_{k,j} \mu_{kj}$$

(8)

$$\mu_{.j} = \frac{1}{m} \sum_k \mu_{kj}$$

(9)

$$\mu_{k.} = \frac{1}{N} \sum_j \mu_{kj}$$

(10)

$$\varepsilon_{.j} = \mu_{.j} - \mu_{..}$$

$$\rho_k = \mu_{k.} - \mu_{..}$$

$\tau_{.j}$ и ρ_k — элементы средних и среднегрупповых соотношений.

Предположение (!), что

$$\mu_{kj} = \mu_{..} + \tau_{.j} + \rho_k,$$

(11)

т.е. предполагается линейное взаимодействие между факторами уровня и среднегрупповыми средними.

Зависимый метод (7) с учетом (11).

Для оценки метода применение МНК.

$$\Rightarrow \hat{\mu}_{..} = \bar{Y}_{..} = \frac{1}{n} \sum_{k,j} Y_{kj}$$

$$\hat{\rho}_k = \frac{1}{N} \sum_j (Y_{kj} - \bar{Y}_{..})$$

$$\hat{\Sigma}_j = \frac{1}{m} \sum_k (Y_{kj} - \bar{Y}_{..})$$

с математическим ожиданием и дисперсией, а также
различия в сумме квадратов остатков (в случае симметричных
групп) равны

$$\Rightarrow Y_{kj} = \bar{Y}_{..} + (\bar{Y}_{k..} - \bar{Y}_{..}) + (\bar{Y}_{..j} - \bar{Y}_{..}) + e_{kj} \quad (12)$$

Можно показать, что это называется разложением.

Однозначно

$$S_1^2 = \sum_k \sum_j e_{kj}^2 \text{ - сумма квадратов остатков} \quad (13)$$

В генетическом анализе есть 2 нормальных случая:

Первый - есть ли различия в среднем по группам?

Решение предложен

$$H_0: \rho_1 = \dots = \rho_m = 0$$

нуль гипотеза

$$H_1: \exists k: \rho_k \neq 0.$$

\Rightarrow расчет средней величины и проверка через F-критерий.

Пример. Установить %, имеющиеся по форме глаз
в 3-х ген. группах в различии по различию
формы.

Задача: оценка, разные формы.

$$N=3, m=5.$$

Нормированные

Очерт формы	Мн-Бора	Амелия	Сор- Брандеско	Среднее по группам
<10	4.65	4.55	4.44	4.547
10-50	4.76	4.40	4.99	4.717
50-100	4.60	4.92	5.00	4.840
100-500	4.83	4.70	4.97	4.833
>500	4.99	4.46	4.93	4.793
среднее по группам	4.766	4.606	4.866	
однотипное среднее				4.746

Проверка номенклатуры

$$H_0: \rho_1 = \dots = \rho_N = 0$$

нормальность

$$H_1: \exists k: \rho_k \neq 0.$$

Наблюданное значение статистики Ранкея $F^* = 1.11$

Число степеней свободы (2,8).

Соответствует уровню значимости

$$p\text{-value} = 0.58$$

Как ураган, в ППП двухфакторный или однократный анализ?

One-way

Two-way

§3. Типы эффектов взаимодействия.

27.11

$$Y_{kjt} = \mu_{kj} + \varepsilon_{kjt} \quad (1)$$

$$k = \overline{1, m}, \quad j = \overline{1, N}, \quad t = \overline{1, T}, \quad n := m \cdot N \cdot T$$

$$\mu_{kj} = \mu_{..} + \rho_k + \tau_j + I_{kj}$$

(2)

I_{kj} - оценивает эффект взаимодействия факторов

Оцениваем модели по МНК. \Rightarrow

$$\hat{\mu}_{..} = \frac{1}{n} \sum_{k,j,t} Y_{kjt} = \bar{Y}_{..}$$

$$\hat{\rho}_k = \frac{1}{N \cdot T} \sum_{j,t} Y_{kjt} - \bar{Y}_{..} = \bar{Y}_{k.} - \bar{Y}_{..}$$

$$\hat{\tau}_j = \frac{1}{m \cdot T} \sum_{k,t} Y_{kjt} - \bar{Y}_{..} = \bar{Y}_{.j} - \bar{Y}_{..}$$

$$\hat{I}_{kj} = \frac{1}{T} \sum_t Y_{kjt} - \bar{Y}_{k.} - \bar{Y}_{.j} - \bar{Y}_{..} = \bar{Y}_{kj} - \bar{Y}_{k.} - \bar{Y}_{.j} - \bar{Y}_{..}$$

$$\Rightarrow Y_{kjt} = \bar{Y}_{..} + (\bar{Y}_{k.} - \bar{Y}_{..}) + (\bar{Y}_{.j} - \bar{Y}_{..}) + (\bar{Y}_{kj} - \bar{Y}_{k.} - \bar{Y}_{.j} - \bar{Y}_{..}) + \varepsilon_{kjt} \quad (3)$$

Можно подозревать, что это разложение на ошиб. компоненты.

Обозначим

$$S_i^2 = \sum_k \sum_j \sum_t \varepsilon_{kjt}^2 \quad - \text{нашная сумма квадратов}$$

S_i^2 после оценки параметров имеет $n-m-N$ степ. свободы. В такой модели можно судить об отсутствии взаимодействий и проводить проверку того же, как это.

Тесты о взаимодействии:

Равнозначно проверка

$$H_0: I_{kj} = 0 \quad \forall (k,j)$$

нуль

$$H_1: I_{kj} \neq 0 \text{ для некоторой пары } (k,j)$$

Если гипотеза верна, то проверка номинальной H_0 и

имеет МНК, но получается, что сумма квадр. отклон.

$$S_I^2 = T \cdot \sum_k \sum_j (\bar{Y}_{kj} - \bar{Y}_k - \bar{Y}_j - \bar{Y}_{..})^2 \quad \text{имеем } (m-1) \cdot (N-1) \text{ степ. добр. (4)}$$

Теорема

При верной H_0 статистика

$$F = \frac{S_I^2 / ((m-1)(N-1))}{S_e^2 / (n - m \cdot N)} \quad (5)$$

имеет распределение Снедекора-Фишера с $((m-1)(N-1), (n - m \cdot N))$ степ. добр.

Данная проверка H_0 проверяется по стандартной схеме.

(Рассмотрено в формальной записи ППП)

18.12 - ненормальное распределение

Табл 10. Дискриминантный анализ (наиболее значимые характеристики)

это некоторый статистич. метод, который позволяет изучать различия между 2 и более группами (классами) по определенным измерениям нескольких непрерывных (характеристик).

Примеры

- 1) исследование причины не работоспособности
- 2) анализ причин заболеваний
- 3) изучение факторов при анализе метода лечения
- 4) разделение спиритуалов по уровню знаний

Вопросы и задачи:

I. Члены семейства - можно ли по измеряющимся параметрам хар-к членам разделить на предложенное семейство

II. Клиенты банка - надо проделать анал. кол-во физических единиц от характеристик, но помочь по определенному признаку. Тогда отнести к тому или иному классу.

Всегда дается предположение, что разделение можно, и будем решать задачу II.

Решение задачи

Пусть имеем общий набор характеристик наследования хар-к (X_1, \dots, X_p) = X . X имеет многомерное нормальное распределение с вектором средних $\mu = (\mu_1, \dots, \mu_p)$ и ковариационной матрицей $\Sigma = (\sigma_{ij})$. Тогда можно 2 (или более). \Rightarrow

$$\mu^{(1)} = (\mu_{11}, \dots, \mu_{1p})$$

$$\mu^{(2)} = (\mu_{21}, \dots, \mu_{2p})$$

и одна и та же матрица ковар.

Задачу решают методом главных компонент

$$Z = z_1 \cdot X_1 + \dots + z_p \cdot X_p \quad (1)$$

Все решаемые при этом будут иметь следующий вид:

Правильное решение:

- 1) Если $Z \geq c$, то отнести объект к первому классу
- 2) Если $Z < c$, то ко второму классу

z_1, \dots, z_p, c - параметры решаемого правила.

Хотим подразумевать максимум, которого минимизируем вероятность ошибочной классификации. Основанка задачи

§2. Классификация в случае избесимых параметров.

Наша задача задать $\mu^{(1)}$, $\mu^{(2)}$ и Σ .

$W^{(1)}$ — первое сополутие

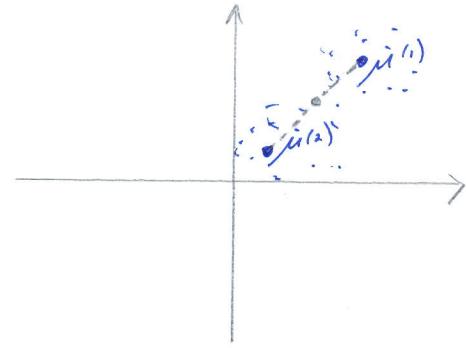
$W^{(2)}$ — второе

Если X из $W^{(1)}$, то с.в. Z имеет одномерное нормальное распределение со средним

$$m_1 = d_1 \cdot \mu_{11} + \dots + d_p \cdot \mu_{1p}$$

и дисперсией

$$\sigma_z^2 = \sum_{i,j=1}^p d_i \cdot \sigma_{ij} \cdot d_j$$



(2)

аналогично, если X из $W^{(2)}$, то с.в. Z имеет одномерное нормальное распределение со средним

$$m_2 = d_1 \cdot \mu_{21} + \dots + d_p \cdot \mu_{2p}$$

и тем же дисперсией σ_z^2 .

На интуитивном уровне понятно, что для оптимального правила решения неравенств $\mu^{(1)}$ и $\mu^{(2)}$ надо помимо гаусса, но с учётом дисперсии σ_z^2 .

Этот задачей занимается автор Махаланобис (Mahalanobis).

Он предложил расц. критерий Биномия

$$\Delta^2 = \frac{(m_1 - m_2)^2}{\sigma_z^2} \quad \text{— расстояние Махаланобиса}$$

Раньше Фишер (1936) показал, что максимум Δ^2 достигается

$$\text{на } d_1, \dots, d_p, \text{ который есть решение систему}$$

$$d_1 \cdot \sigma_{11} + \dots + d_p \cdot \sigma_{1p} = \mu_{11} - \mu_{21}, \quad \leftarrow \text{СЛАУ}$$

$$i = \overline{1, p}$$

Далее мы c=? обозначим

$P(1|2)$ — вероятность ошибки первого рода $\leftarrow W^{(1)}$, когда он из $W^{(2)}$,

$P(2|1)$ — — — $\leftarrow W^{(2)}$, когда он из $W^{(1)}$.

И мы знаем, что минимизация ожидаемого risk-функционала не даётся.

Вариационное $P(1|2) + P(2|1) \rightarrow \min_c$. Можно показать, что min

достижимо вида

$$c = \frac{m_1 + m_2}{2}$$

это экспериментальное решение.

(3)

(4)

(5)

(6)

Береже решение этой задачи в доказательной установке правил Байеса (когда задача заложена гипотезой)

Но Тогда q_1 - априорная вероятность $x \sim W^{(1)}$.
 q_2 - ... $x \sim W^{(2)}$.

Тогда X имеет априорное распределение $f_1(x)$ и $f_2(x)$ при $W^{(1)}$ и $W^{(2)}$.

Выводим по формуле Байеса априорные вероятности:

$$P(W^{(i)} | x) = \frac{q_i \cdot f_i(x)}{q_1 \cdot f_1(x) + q_2 \cdot f_2(x)} \quad (7)$$

Будем выражать эту вероятность, где конечной априорной вероятности. Получаем следующие результаты:

Правильное решение:

1) Принимаем $W^{(1)}$, если $\frac{q_1 f_1(x)}{q_2 f_2(x)} \geq 1$ (8)

2) Принимаем $W^{(2)}$, если $\frac{q_1 f_1(x)}{q_2 f_2(x)} < 1$ (9)

Чтобы облегчить задачу, будем считать, что максимум минимизируем ожидаемую вероятность ошибки.

$$q_1 \cdot P(z|1) + q_2 \cdot P(z|2) \quad (10)$$

Для нормальных распред. из (8) и (9) \Rightarrow

1) Принимаем $W^{(1)}$, если $\sum_i z_i x_i \geq \frac{m_1 + m_2}{2} + \ln\left(\frac{q_2}{q_1}\right)$ (11)

2) Принимаем $W^{(2)}$, если $\sum_i z_i x_i < \frac{m_1 + m_2}{2} + \ln\left(\frac{q_2}{q_1}\right)$ (12)

Если $q_1 = q_2 \Rightarrow$ получаем эмпирическое правило.

Всегда же получим решения:

04.12.

$C(2|1)$ - стоимость номера, если принят $W^{(2)}$, а на самом деле имеем $W^{(1)}$

$C(1|2)$ - ... аналогично

Используя, какое минимизируем ожидаемую стоимость номера, т.е.

$$(2|1) \cdot q_1 \cdot P(2|1) + C(1|2) \cdot q_2 \cdot P(1|2)$$

Приходим к след. выражению:

1) Принимаем $W^{(1)}$, если

$$\sum_{i=1}^p z_i x_i \geq \frac{m_1 + m_2}{2} + \ln\left(\frac{q_2 \cdot C(1|2)}{q_1 \cdot C(2|1)}\right)$$

2) Применение $W^{(2)}$, если

$$\sum_i \lambda_i \cdot x_i < \frac{m_1 + m_2}{2} + \ln \left(\frac{q_2 \cdot C(1|2)}{q_1 \cdot C(2|1)} \right)$$

В случае нормального распределения вероятностных оценок критериям имеет вид:

$$P(2|1) = \varphi \left(\frac{K - \frac{1}{2}\Delta^2}{\Delta} \right)$$

$$P(1|2) = \varphi \left(\frac{-K - \frac{1}{2}\Delta^2}{\Delta} \right)$$

$$K = \ln \frac{q_2 \cdot C(1|2)}{q_1 \cdot C(2|1)},$$

Δ^2 -расстояние Махаланобиса

При $C(1|2) = C(2|1)$ возвращается к базисной нормальной задаче.

Если (при этом) $q_1 = q_2$, то получаем эмпирическую проверку \Rightarrow

$$P(2|1) = P(1|2) = \varphi \left(-\frac{\Delta}{2} \right)$$

Устойчивый пример

↑ "нормальную и неизвестную относительно m "

$X = (X_1, X_2)$ - двухмерн. норм. распред. с базисом средних $\mu_1 = (1, 1)$, $\mu_2 = (3, 3)$, $\Sigma = \begin{pmatrix} 1 & 1 \\ 1 & 1,25 \end{pmatrix}$

В случае эмпирической проверки:

$$\begin{cases} d_1 + d_2 = -2 \\ d_1 + 1,25 \cdot d_2 = -2 \end{cases} \Rightarrow \begin{cases} d_2 = 0 \\ d_1 = -2 \end{cases} \Rightarrow Z = -2 \cdot X_1 - \text{дискриминантная функция.}$$

$$\Rightarrow m_1 = -2, m_2 = -6, \sigma_Z^2 = 4$$

$$c = \frac{m_1 + m_2}{2} = -4$$

Таким образом критерий имеет вид:

1) критерий $W^{(1)}$, если

$$-2x_1, x_1 \geq -4 \Leftrightarrow x_1 \leq 2$$

2) критерий $W^{(2)}$, если

$$x_1 > 2.$$

Ограничение обеих критерий задают характеристики $x = (2,5, 1)$.

$$\Rightarrow 2,5 > 2 \Rightarrow \text{надежность критерия } W^{(2)}$$

относится к

§3. Квадратичные в случае неизвестных параметров

Пусть μ_1, μ_2, Σ - неизвестные!

Задача $q_1, q_2, c(2|1), c(1|2)$.

Пусть задача где однородных наблюдений

X_{11}, \dots, X_{1n_1} и X_{21}, \dots, X_{2n_2} ,

где $X_{ij} = (X_{ij}^{(1)}, \dots, X_{ij}^{(p)})^T$ - неавт. и именем многомер. норм.

распред. с параметрами μ_i, Σ

Нахождим оценки:

$$\hat{\mu}_1 = \bar{X}_1 = \frac{1}{n_1} \sum_j X_{1j}$$

$$\hat{\mu}_2 = \bar{X}_2 = \frac{1}{n_2} \sum_j X_{2j}$$

$$\hat{\Sigma} = S = \frac{1}{n_1 + n_2 - 2} [(n_1 - 1)S_1 + (n_2 - 1)S_2]$$

$$\text{где } S_i = \frac{1}{n_i - 1} \sum_j (X_{ij} - \bar{X}_i)(X_{ij} - \bar{X}_i)^T, \quad i=1,2$$

(В случае многомер. норм. распред это несущ. и состоян. оценки.)

Далее применят обобщённого ближайшего соседа, где наименее параметр заменим на их оценки. Алгоритмы решения задачи квадратичны.

1) Решаем систему уравнений

$$z_1 \cdot S_{11} + \dots + z_p \cdot S_{1p} = \hat{\mu}_{1i} - \hat{\mu}_{2i}, \quad i = \overline{1, p}$$

\Rightarrow находим $\hat{z}_1, \dots, \hat{z}_p$.

2) Определяем дискриминантную функцию:

$$\hat{Z} = \hat{z}_1 \cdot X_1 + \dots + \hat{z}_p \cdot X_p$$

3) по однородным наблюдениям вычислим оценки для дискриминантных функций

$$\hat{Z}_{1j}, \hat{Z}_{2j}.$$

4) Заменим n_1, n_2 на их оценки

$$\hat{n}_1 = \frac{1}{n_1} \cdot \sum_{j=1}^{n_1} \hat{Z}_{1j},$$

$$\hat{n}_2 = \frac{1}{n_2} \sum_{j=1}^{n_2} \hat{Z}_{2j}$$

5) Определяем σ_z^2 :

$$\hat{\sigma}_z^2 = S_z^2 = \sum_{i,k=1}^p \hat{z}_i \cdot S_{ik} \cdot \hat{z}_k$$

6) Имеем некий набор измерений $x = (x_1, \dots, x_p)$

a) опишите это в $w^{(1)}$, если

$$\hat{Z} = \sum_{i=1}^p \hat{z}_i x_i > \frac{\hat{m}_1 + \hat{m}_2}{2} + \ln \left(\frac{q_2 \cdot C(1|2)}{q_1 \cdot C(2|1)} \right).$$

Если n_1 и n_2 скажем, то "the economy" = экономическая
представляемая в одинаковых реалиях

Таблица 11. Классический анализ (не знает что такое кванты)

§1. Описание данных

У некоторого множества (некоторой однозначности) отобрано n объектов $I = (I_1, \dots, I_n)$.

У каждого объекта изучено n некотоих характеристик $C = (c_1, \dots, c_p)^T$ — это координаты характеристики.

$\Rightarrow x_{ij}$ — измерение i -ой характеристики у j -го объекта.

$X_j = (x_{1j}, \dots, x_{pj})^T$ — координаты характеристики у объекта I_j .

$\Rightarrow X = (X_1, \dots, X_n)$ — матрица измерений.

§2. Задача классического анализа.

Такие $m < n$. (Дано многие $m < n$)

Предполагается на основе измерений X разделять множество объектов I на m квантов (классов) T_1, \dots, T_m максимум, минимум:

- 1) каждый объект I_j принадлежит только одному кванту,
- 2) каждый объект одного кванта в некотором смысле близок (сходен),
- 3) объекты из разных квантов далеки (несходны).

§3. Меры сходства.

Предположим, что в \mathbb{R}^p задано некоторое расстояние (метрика?). Тогда метрика в соответствии с нашей организацией квантов, но не беседа!

1) Евклидовое расстояние

$$d_2(X_k, X_j) := \left[\sum_{i=1}^p (X_{ik} - X_{ij})^2 \right]^{1/2}$$

2) l_1 -норма

$$d_1(X_k, X_j) = \sum_{i=1}^p |X_{ik} - X_{ij}|$$

3) максимальная норма (норма при $p=\infty$ из l_∞ -нормы)

$$d_\infty(X_k, X_j) = \sup_i |X_{ik} - X_{ij}|$$

4) Векторное расстояние

$$D^2(X_k, X_j) = (X_k - X_j)^T \cdot W^{-1} \cdot (X_k - X_j),$$

где W — матрица весов (обычно задана в задаче)

Задание. Тычи $y = b \cdot x$. Тогда

$D^2(y_k, y_j) = D^2(x_k, x_j)$, т.е. расст. между линейн. обр. различных начальных преобразований.

11.12

$D = (d_{kj})$ - матрица парных расстояний

Оп. Неотрицательная функция $S(x_k, x_j) := s_{kj}$ называется неравенством сходства, если

1) $0 \leq S(x_k, x_j) < 1$, если $x_k \neq x_j$

2) $S(x_k, x_k) = 1$

3) $S(x_k, x_j) = S(x_j, x_k)$

Задача

1) s_{kj} - неравенство сходства неравенств x_k и x_j

2) Если x contains из 0 и 1, то s_{kj} наз-ся "коэф. ассоциации" или "неравенство коэф. соподчиненности".

3) В статистике often расчет в качестве s_{kj} использует модуль коэф-та корреляции

Оп. Тычи X - матрица неравенств. Весимна

$S_d := \frac{1}{2} \sum_{k,j=1}^n d(x_k, x_j)$ - однор. расчение ил-ла обектов I.

Оп. $\bar{S}_d := \frac{S_d}{Nd}$ - среднее расчение,

$Nd = \frac{n(n-1)}{2}$ - это же, но, конечно, не делится

Оп. $S_x := \sum_{j=1}^n (x_j - \bar{x})(x_j - \bar{x})^T$ - матрица расчения, где (1)

$\bar{x} = \frac{1}{n} \sum_{j=1}^n x_j$.

Оп. Весимна

$S_t := \text{tr } S_x = \sum_{j=1}^n \sum_{i=1}^p (x_{ij} - \bar{x}_i)^2$ - статистическое расчение

Оп. Весимна

$\det(S_x)$ - статистическое расчение, соавтоматически определенное

§4. Рассмотрение между классами

Таким же образом находим $I = (I_1, \dots, I_m)$ и

$J = (J_1, \dots, J_{n_2}) \Rightarrow X = (X_1, \dots, X_{n_1})$ и $Y = (Y_1, \dots, Y_{n_2})$.

Оп. Бенуана

$D_1(I, J) := \min_{k, j} d(X_k, Y_j)$ — минимальное локальное расстояние.

Оп. Бенуана

$D_2(I, J) := \max_{k, j} d(X_k, Y_j)$ — максимальное локальное расстояние

Оп. $D_3(I, J) := \sum_k \sum_j d(X_k, Y_j) / n_1 \cdot n_2$ — среднее расстояние между классами

Оп. $D_4(I, J) := (\bar{X} - \bar{Y})^T (\bar{X} - \bar{Y}) / \frac{n_1 \cdot n_2}{n_1 + n_2}$ — стандартированное расстояние между классами

Когда расстояние между классами, то оно

Баиномиил ноль классов $K = I \cup J$.

Обозначим

$$S_K = \sum_{k=1}^{n_1} (X_k - M)(X_k - M)^T + \sum_{j=1}^{n_2} (Y_j - M)(Y_j - M)^T,$$

где $M = \frac{\sum_k X_k + \sum_j Y_j}{n_1 + n_2}$.

Мы видим, что

$$S_K = S_I + S_J + \frac{n_1 \cdot n_2}{n_1 + n_2} (\bar{X} - \bar{Y})(\bar{X} - \bar{Y})^T. \quad (*)$$

Оп. Бенуана $\frac{n_1 \cdot n_2}{n_1 + n_2} (\bar{X} - \bar{Y})(\bar{X} - \bar{Y})^T$ — максимальное межклассовое расстояние.

$\Rightarrow \operatorname{tr} \frac{n_1 \cdot n_2}{n_1 + n_2} (\bar{X} - \bar{Y})(\bar{X} - \bar{Y})^T = \frac{n_1 \cdot n_2}{n_1 + n_2} (\bar{X} - \bar{Y})^T \cdot (\bar{X} - \bar{Y})$ — стандартированное расстояние между $I \cup J$ = межгрупповая сумма квадратов

§5. Понедельное поступление мастеров.

Описание объекта алгоритма:

- 1) Стартует бе можи распаштиваць кога однозначне каспер.
- 2) Відкрити дба чорна: $0 < \frac{t}{s} < \frac{s}{t} < \infty$.
- 3) Еш та рахоманія менше s , то він обєднує в один каспер і спаковується.
- 4) Еш еш та каспера, рахоманія менше s , то обєднує их в один.
- 5) Пересямлює рахоманія відчуті касперов і меншу їхні.
- 6) Приєднуда чорнотається до них чорн, нова рахоманія відчуті каспер не більше t , а меншу каспари не більше s .

Додаток 25.12. 6 10.30 (у пазухи, аудиторію чиєї сам)

Це складає поняття.

Но я не пам'яте про збільшення, а то чиан збільшувати зображені зображені!

Основная учебно-методическая литература

1. Андерсон Т. Введение в многомерный статистический анализ. – М.: Наука, 1963.
2. Johnson R.A. and Winchern D.W. Applied Multivariate Statistical Analysis. – Pearson Prentice Hall, 2007.
3. Хохлов Ю.С. Эконометрика. Вводный курс: Учебное пособие. М.: Издательский отдел факультета ВМиК МГУ, 2006. – 100 с.

Дополнительная учебно-методическая литература

- 1) Факторный, дискриминантный и кластерный анализы. – М.: Финансы и статистика, 1989.
- 2) Шеффе Г. Дисперсионный анализ. – М.: Наука, 1980.
- 3) Болч К., Хуань К.Дж. Многомерные статистические методы для экономики. – М.: Статистика, 1979.
- 4) Hardle W., Simar L. Applied Multivariate Statistical Analysis. - <http://www.xplorestat.de/ebooks/ebooks.html>
- 5) Кулаичев А.П. Методы и средства комплексного анализа данных. – М.: Форум-Инфра-М, 2006.