

## 23 2.2 Метод на главните компоненти. Факторен анализ

Факторният анализ е възникнал като задача при обработката на данни от психологични анкети.

Данни се представят на данните във формата:

$$X = FL + E$$

Тук  $X \in \mathbb{R}^{n \times m}$  е матрицата на центрираните и нормирани наблюдения, а  $E \in \mathbb{R}^{n \times m}$  грешките. Във факторния анализ са приети следните извания.

- Факторни тежести се наричат коефициентите на разпадане на оригиналните променливи по фактори. те  $L \in \mathbb{R}^{k \times m}$ .
- Факторни стойности се наричат оценките стойности на факторите за всяко наблюдение  $F \in \mathbb{R}^{n \times k}$ .
- Общност-товар е относителната дисперсия на всяка променлива в новото ѝ описание. Т.е. общностите показват доколко стойностите на конкретната променлива могат да бъдат предсказвани или възстановени от факторните стойности.

Основните задачи на факторния анализ са:

- Да се определи бр. на факторите, достатъчни за описание на звестеното (втората размерност на матрицата  $F$ ).
- Да се намерят факторите това означава да се изчислят коефициентите на линейните функции  $L$ , представящи променливите (посредством факторите).

## Корелационна матрица

Корелационна (или ~~кор~~ ковариационна) матрица на случаен нормален вектор в  $m$ -мерното пространство е симетрична и неотрицателно определена. Тя може да бъде представена:

$$S = \sum_{i=1}^m a_i e_i e_i' \quad \dots (1)$$

Тук числата  $\{a_1, a_2, \dots, a_m\}$  са наричани собствени числа, а векторите  $\{e_i, i=1, \dots, m\}$  - собствени вектори, защото удовлетворяват условието:

$$S e_i = a_i e_i \quad \dots (2)$$

Те са ортогонални с единична норма и образуват базис в  $m$ -мерното евклидово пространство

$$S^2 = \sum \|x_i\|^2 = \sum \|P_k x_i\|^2 + \sum \|x_i - P_k x_i\|^2$$

$S^2 = \sum \|x_i\|^2$  - представява центрираност на данните  
то естествено е да търсим това подпространство  
на  $n$  с размерност  $k$  (за фиксирано  $k, 1 \leq k \leq m$ ),  
за което експерименталните точки биха се пренесли  
минимално при своето проектиране  $P$  в оуни.  
Тогава най-малко би се изместила и  $S^2$ , което се  
вижда и от равенството.

Ерой на факторите:

$$\sum_{i=1}^m \|P_k x_i\|^2 \geq 0.95 \sum_{i=1}^m \|x_i\|^2$$

Друг срещан критерий е избора на собствените вектори  
отс. собствени стойности, които са по-големи 1.  
Ако използваме корелационната матрица  $\sum a_i = \text{tr } R = m$  -  
Критерий на Кайзер

Средното кетоди, обекта ваеги интерпретационно на факторите, може да най-популярен е методът варимакс трансформацията. Идеята на варимакс-трансформацията е така да бъдат променени факторите, че те да запазят добрите си свойства (пълното на описание на извадката) и да получат по-добри интерпретативни качества. Препоръчват се факторните тела като големите нарастват, а малките намаляват, което води до това че всеки фактор се обяснява от по-малко на брой променливи. Факторите вече не са толкова неясни. Процесът се обяснява на променливи. Той се състои в максимизиране на функционала:

$$\max_K \text{var}(KL) = \max_K \sum_{i=1}^K \sigma^2(i)$$

по безвзаимните избор на факторите във всеки фиксиран координатен подпространство е.

Тук  $\sigma^2(i)$  е "дисперсията" на квадратичните факторни тела, а  $K \in \mathbb{R}^{K \times K}$  е ортогонална матрица. Телата на променливите в даден фактор са записани по редове в  $L$ . Максимизирането на дисперсията води до увеличаване разликата между големите и малките тела. Така се вижда кои променливи са ясно представени в даден фактор - те получават тела близо до единица. Като максимизиране дисперсията по променливи (по стълбове) но тези тела, съответните методи наричаме кватримакс.

Възникват варианти на едновременно максимизиране претегляна сума на двата функционала.

- бикватримакс - два функционала са с равно тегло
- еквимакс - варимакс функционалът е с тегло  $\frac{K}{2}$

и значимостите на така определените фактори въ-  
з основа на обектите  $F$ .

3. Да се интерпретират получените резултати в  
термини на предметната област на данните.

Методика на факторния анализ.

Ако си представим наблюдавания като облак от точки,  
в пространството, то естествена характеристика  
на този облак биха били неговата форма и раз-  
те се отразяват до известна степен в ковариационната  
матрица. Ковариационната матрица е обект на  
класическия факторен анализ.

Съществува обаче редица варианти на факторния  
анализ, които използват други аналогични матрици,  
които интерпретацията е различна.

Главните компоненти. – Метод на главните компоненти  
позволява да определим до колко определен брой фактори  
отбавят влиянието. ~~на~~ Главните компоненти  
съответстват на осите на елипсото на разсейва-  
нето на точките, представящи обекти (данни) в  
пространството на наблюдаваните променливи.  
Техният брой е равен на броя на променливите.  
Главните компоненти могат да бъдат разглеждани  
като фактори, т.е. като нови променливи, които  
изцяло отбавят извадката и са независими.

Тази независимост позволява да представим общата  
дисперсия на извадката като сума от дисперсиите  
на новите променливи. Отстранявайки едновременно  
факторите тези които имат незначителна дисперсия  
(разсейване) може да дадем описание на свойст-  
вото на извадката с ~~не~~ малък брой <sup>нови</sup> променливи.

Когато броят на факторите е избран, решението поумекно по метода на главните компоненти може да бъде записано:

$$X = FL + E$$

Тук  $X \in \mathbb{R}^{n \times m}$  е матрица на изпитирани и коригирани наблюдения,  $F \in \mathbb{R}^{n \times k}$  - са факторните стойности  $L \in \mathbb{R}^{k \times m}$  - са факторните тежа, а  $E \in \mathbb{R}^{n \times m}$  - грешките.

Решението по метода на Хотелинг дава дава начини на функционала  $\text{tr } E'E$ . Ако разгледаме следното представяне на матрицата  $X = UDV$ . Тук

$U \in \mathbb{R}^{n \times n}$  и  $V \in \mathbb{R}^{m \times m}$  са ортогонални, а  $D \in \mathbb{R}^{n \times m}$  -

"диагонална" с неотрицателни елементи.

$$(d_i, i \geq d_{i+1}, i+1)$$

Това има следното представяне на ковариационната матрица

$$X'X = V'D'U'UDV = V'D'DV = \sum_{i=1}^m d_i^2 v_i v_i'$$

Аналог. е представяне на матрицата  $XX'$ , ко собствени вектори са в друго пространство. Имаше  $d_i^2 i = \lambda$

фиксиране  $k$  и "орязване" първите  $k$  стълба на една матрица. Получаваме  $F = n^{-1/2} (UD)_k$ ,  $L' = (V')_k$

Числата  $d_i^2 i / n$  са дисперсиите на факторите, а

$\sum_{i=1}^k d_i^2 i$  са т.н. общият на променливите. Те

обясняват каква част от дисперсията е обяснена с така подобрение фактори.

Матрицата  $L \in \mathbb{R}^{k \times m}$