

Análise Multivariada de Dados - Aula 07

Análise Fatorial

Kaique Matias de Andrade Roberto

Ciências Atuariais

HECSA - Escola de Negócios

FIAM-FAAM-FMU

Conteúdo

- 1. Conceitos que aprendemos em Aulas anteriores
- 2. A Ideia Geral
- 3. Procedimento para uma análise fatorial
- 4. Análise de fatores de componentes principais
- 5. Exemplos
- 6. Comentários Finais
- 7. Referências

Conceitos que aprendemos em

Aulas anteriores

Conceitos que aprendemos em Aulas anteriores

Nas últimas Aulas lidamos em aspectos teóricos e práticos da Análise de Componentes Principais (PCA).

A análise de fatores tem objetivos que são similares àqueles da análise de componentes principais.

A ideia básica é que pode ser possível descrever um conjunto de p variáveis $X_1, X_2, ..., X_p$ em termos de um número menor de índices ou fatores, e no processo obter uma melhor compreensão do relacionamento destas variáveis.

Há, no entanto, uma diferença importante. A análise de componentes principais não é baseada em um modelo estatístico particular, enquanto que a análise de fatores é baseada em um modelo.

O desenvolvimento inicial de análise de fatores é o resultado do trabalho de Charles Spearman.

Enquanto estudava correlações entre escores de testes de estudantes de vários tipos, ele notou que muitas correlações observadas poderiam estar contidas em um modelo simples (Spearman, 1904).

Por exemplo, em um caso ele obteve a matriz de correlações mostrada na Tabela 7.1 (numeração do livro), para alunos em uma escola preparatória para seus escores em testes em clássicos, francês, inglês, matemática, discriminação de tom e música.

Tabela 7.1 Correlações entre escores de testes para meninos em uma escola preparatória

				Discriminação		
	Clássicos	Francês	Inglês	Matemática	de tom	Música
Clássicos	1,00	0,83	0,78	0,70	0,66	0,63
Francês	0,83	1,00	0,67	0,67	0,65	0,57
Inglês	0,78	0,67	1,00	0,64	0,54	0,51
Matemática	0,70	0,67	0,64	1,00	0,45	0,51
Discriminação de tom	0,66	0,65	0,54	0,45	1,00	0,40
Música	-0,63	0,57	0,51	0,51	0,4,0	1,00

Fonte: De Spearman, C. (1904), Am. J. Psychol., 15, 201-293.

Ele notou que esta matriz tinha a interessante propriedade de que quaisquer duas linhas eram quase proporcionais se as diagonais fossem ignoradas. Então para as linhas clássicos e inglês na Tabela 7.1 (numeração do livro), há razões:

$$\frac{0.83}{0.67} \approx \frac{0.70}{0.64} \approx \frac{0.66}{0.54} \approx \frac{0.63}{0.51} \approx 1.2$$

Baseado nesta observação, Spearman sugeriu que os seis escores de testes fossem descritos pela equação

$$X_i = a_i F + e_i$$

em que X_i é o i-ésimo escore depois dele ter sido padronizado para ter uma média zero e um desvio-padrão um para todos os alunos.

Aqui a_i é uma constante; F é um valor "fator", o qual tem média zero e desvio-padrão um para todos os alunos; e e_j é a parte de X_i que é específica para o i-ésimo teste somente.

Spearman mostrou que uma razão constante entre as linhas de uma matriz de correlações segue como uma consequência destas suposições, e que, portanto, este é um modelo plausível para os dados.

Além das razões de correlações constantes, segue também que a variância de X_i é dada por

$$Var(X_i) = Var(a_iF + e_i) = Var(a_iF) + Var(e_i)$$
$$= a_i^2 Var(F) + Var(e_i) = a_i^2 + Var(e_i)$$

porque a_i é uma constante, F e e_j são assumidas independentes, e a variância de F é assumida ser unitária.

Também, porque
$${\sf Var}(X_i)=1,$$

$$1=a_i^2+{\sf Var}(e_i).$$

Portanto a constante a_i a qual é chamada de carga do fator, é tal que seu quadrado é a proporção da variância de X_j que está contida no fator.

Com base no seu trabalho, Spearman formulou sua teoria de dois fatores de testes mentais. De acordo com esta teoria, cada resultado do teste é composto de duas partes, uma que é comum a todos os testes (inteligência geral), e outra que é específica para o teste.

Isto dá o modelo de análise de fatores geral, o qual estabelece que

$$X_i = a_{i1}F_1 + a_{i2}F_2 + ... + a_{im}F_m + e_i$$

em que X_i é o i-ésimo escore do teste com média zero e variância unitária; a_{ij} são as cargas dos fatores para o i-ésimo teste; $F_1, ..., F_m$ são m fatores comuns não correlacionados, cada um com média zero e variância unitária; e e_i é um fator específico somente para o i-ésimo teste que é não correlacionado com qualquer dos fatores comuns e tem média zero.

Com este modelo,

$$Var(X_i) = 1 = a_{i1}^2 Var(F_1) + a_{i2}^2 Var(F_2) + ... + a_{im}^2 Var(F_m)$$

= $a_{i1}^2 + a_{i2}^2 + ... + a_{im}^2 + Var(e_i)$

em que $a_{i1}^2 + a_{i2}^2 + ... + a_{im}^2$ é chamado a **comunalidade** de X (a parte de sua variância que é relacionada aos fatores comuns), e $Var(e_i)$ é chamada **especificidade** de X (a parte de sua variância que não é relacionada aos fatores comuns).

Pode também ser mostrado que a correlação entre X_i e X_j é

$$Correl(X_i, X_j) = a_{i1}a_{j1} + a_{i2}a_{j2} + ... + a_{im}a_{jm}.$$

Portanto dois escores de teste podem somente ser altamente correlacionados se eles têm altas cargas nos mesmos fatores.

Além disso, como a comunalidade não pode exceder um, é preciso que $-1 < a_{ij} < 1$.

Procedimento para uma análise

fatorial

Os dados para uma análise de fatores têm a mesma forma como para uma análise de componentes principais.

Caso	X_1	X ₂	•••	X_{p}
1	X	X ₁₂	•••	X _{1p}
2	χ_{21}	X ₂₂	•••	X_{2p}^{P}
-	•	•	•••	•
-	•	-	•••	•
-			•••	
n	X_{n1}	X _{n2}	•••	X _{np}

Há três estágios para uma análise de fatores. Para começar, cargas de fatores provisórios a_{ij} são determinadas.

Uma abordagem começa com uma análise de componentes principais e negligencia os componentes principais após os primeiros m, os quais são então tomados como sendo os m fatores.

Os fatores encontrados desta maneira são não correlacionados entre si, e são também não correlacionados com os fatores específicos. Isto pode não ser um problema desde que as comunalidades sejam altas.

Qualquer que seja a maneira como as cargas de fatores provisórios são determinadas, é possível mostrar que eles não são únicos.

Se $F_1, F_2, ..., F_m$ são os fatores provisórios, então as combinações lineares deles da forma

$$F_{1}^{*} = d_{11}F_{1} + d_{12}F_{2} + \dots + d_{1m}F_{m}$$

$$F_{2}^{*} = d_{21}F_{1} + d_{22}F_{2} + \dots + d_{2m}F_{m}$$

$$\vdots$$

$$\vdots$$

$$F_{m}^{*} = d_{m1}F_{1} + d_{m2}F_{2} + \dots + d_{mm}F_{m}$$

podem ser construídos de modo a serem não correlacionados e explicar os dados tão bem quanto os fatores provisórios.

De fato, há uma infinidade de soluções alternativas para o modelo de análise de fatores. Isto leva ao segundo estágio na análise, o qual é chamado de rotação de fator.

Neste estágio, os fatores provisórios são transformados a fim de encontrar novos fatores que sejam mais fáceis de interpretar.

Girar ou transformar neste contexto significa essencialmente escolher os valores d_{ij} nas equações já vistas.

O último estágio de uma análise envolve calcular os escores dos fatores. Estes são os valores dos fatores rotacionados $F_1^*, F_2^*, ..., F_m^*$ para cada um dos n indivíduos para os quais os dados estão disponíveis.

Geralmente, o número de fatores (m) depende do analista, apesar de algumas vezes poder ser sugerido pela natureza dos dados.

Quando uma análise de componentes principais é usada para encontrar uma solução provisória, uma regra rústica envolve escolher m como sendo o número de autovalores maiores do que a unidade na matriz de correlações dos escores do teste.

A lógica aqui é a mesma que foi explicada no capítulo anterior sobre análise de componentes principais. Um fator associado com um autovalor menor que a unidade responde por menos variação nos dados do que os escores de teste originais.

Em geral, aumentando m aumenta as comunalidades das variáveis. Entretanto, comunalidades não são alteradas por rotação de fator.

Rotação de fatores pode ser ortogonal ou oblíqua. Com rotação ortogonal, os novos fatores são não correlacionados, como os fatores provisórios. Com rotação oblíqua, os novos fatores são correlacionados.

Qualquer que seja o tipo de rotação usada, é desejável que as cargas de fator para os novos fatores sejam ou próximas de zero ou muito diferentes de zero.

Um a_{ij} próximo de zero significa que X_i não é fortemente relacionado com o fator F_j . Um grande valor positivo ou negativo de a_{ij} significa que X_i , é determinado em grande parte por F_j .

Se cada escore de teste é fortemente relacionado com alguns fatores, mas não relacionado com outros, então isso toma os fatores mais fáceis de serem identificados do que o seria em outro caso.

Um método de rotação de fatores ortogonal que é muitas vezes usado é chamado de rotação varimax.

Este é baseado na suposição de que a interpretabilidade do fator j pode ser medida pela variância dos quadrados de suas cargas de fator, i.e., a variância de $a_{1j}^2, a_{j2}^2, ..., a_{mj}^2$.

Se esta variância é grande, então os valores a_{ij} tendem a ser ou próximos de zero ou próximos da unidade. A rotação varimax, portanto, maximiza a soma destas variâncias para todos os fatores.

Inúmeros outros métodos de rotação ortogonal têm sido propostos. Entretanto, rotação varimax parece ser uma boa abordagem padrão.

Algumas vezes analistas de fatores são preparados para desistir da ideia dos fatores serem não correlacionados a fim de tomar as cargas de fator tão simples quanto possível.

Uma rotação oblíqua pode então dar uma melhor solução do que uma ortogonal. Novamente, há numerosos métodos disponíveis para fazer a rotação oblíqua.

Um método para calcular os escores de fator para indivíduos, baseado nos componentes principais, é descrito na próxima seção. Existem outros métodos disponíveis, de modo que aquele escolhido para uso dependerá do pacote computacional que está sendo usado na análise.

Análise de fatores de

componentes principais

Foi observado anteriormente que uma maneira de fazer uma análise de fatores é começar com uma análise de componentes principais e usar os primeiros componentes principais como fatores não rotacionados.

Isto tem a virtude da simplicidade, apesar que, devido aos fatores específicos $e_1, e_2, ..., e_p$, serem correlacionados, o modelo de análise de fatores não é muito correto.

Algumas vezes analistas de fatores fazem primeiro uma análise de fatores de componentes principais e então, após isto, tentam uma outra abordagem.

O método para encontrar os fatores não rotacionados é como segue. Com p variáveis, haverá o mesmo número de componentes principais.

Estes são combinações lineares das variáveis originais

$$Z_{1} = b_{11}X_{1} + b_{12}X_{2} + ... + b_{1p}X_{p}$$

$$Z_{2} = b_{21}X_{1} + b_{22}X_{2} + ... + b_{2p}X_{p}$$

$$\vdots$$

$$\vdots$$

$$Z_{p} = b_{p1}X_{1} + b_{p2}X_{2} + ... + b_{pp}X_{p}$$

em que os valores b_{ij} são dados pelos auto vetores da matriz de correlações.

Esta transformação dos valores X para valores Z é ortogonal $(A^{-1}=A^t)$, de modo que o relacionamento inverso é simplesmente

Para uma análise de fatores, somente m das componentes principais são retidas, assim as últimas equações se tornam

em que e_i é uma combinação linear dos componentes principais Z_{m+1} a Z_p .

Tudo que é preciso ser feito agora é escalonar os componentes principais $Z_1, Z_2, ..., Z_m$ para terem variâncias unitárias, como requerido pelos fatores.

Para fazer isto, Z_i precisa ser dividido pelo seu desvio-padrão, o qual é $\sqrt{\lambda_i}$, a raiz quadrada do correspondente autovalor na matriz de correlações.

As equações então se tornam

$$\begin{split} X_1 &= \sqrt{\lambda_1} b_{11} F_1 + \sqrt{\lambda_2} b_{12} F_2 + \ldots + \sqrt{\lambda_m} b_{m1} F_m + e_1 \\ X_2 &= \sqrt{\lambda_1} b_{12} F_1 + \sqrt{\lambda_2} b_{22} F_2 + \ldots + \sqrt{\lambda_m} b_{m2} F_m + e_2 \\ & \cdot \\ & \cdot \\ X_p &= \sqrt{\lambda_1} b_{1p} F_1 + \sqrt{\lambda_2} b_{2p} F_2 + \ldots + \sqrt{\lambda_m} b_{mp} F_m + e_p \end{split}$$

em que
$$F_i = Z_i/\sqrt{\lambda_i}$$
.

O modelo de fatores não rotacionado é então

onde $a_{ij} = \sqrt{\lambda_j} b_{ji}$.

Após uma rotação varimax ou outro tipo de rotação, uma nova solução tem a forma

em que F_i^* representa o novo i-ésimo fator.

Os valores do *i*-ésimo fator não rotacionado são justamente os valores do *i*-ésimo componente principal após eles terem sido escalonados para terem uma variância um.

Tabela 7.2 Autovalores e autovetores para dados de emprego europeu da Tabela 1.5 Autovetores X_1 X_2 X_3 X_5 X_6 X_7 X_8 X₉ X_4 Autovalores 'AGR MIN FAB FEA CON SER FIN SSP TC -0,2053.111 0.512 0,375 -0.246-0,315 -0.222-0.382-0,131-0.428-0.0240.000 0,432 0.109 -0.242-0.408-0.5530.055 0,516 1,809 1.495 -0,2780,516 -0,503-0,2920,071 0,064 -0.0960.360 0.413 -0.0421.063 0.016 0.113 0,058 0,023 0,783 0,169 -0.489-0.3170,025 -0,3450.231 -0.854-0.0640.269 -0.1330.046 0,023 0,705 -0.0280,208 -0.399-0.167-0,1360.311 -0.0450,203 -0,5030.674 0.293 0.166 -0.212-0,2380,065 0,014 -0.165-0,4630,619 -0,492-0.4310.157 0.030 0.203 -0.026-0.0450,504 0,203 0.539 -0.447-0.447-0.030-0.129-0.245-0.191-0.410-0.0610.000 -0.582-0,419

Aqui, os valores entre parênteses são as comunalidades. Por exemplo, a comunalidade para a variável X_1 é $(0,90)^2 + (-0,03)^2 + (-0,34)^2 + (0,02)^2 = 0,93$.

As comunalidades são bastante altas para todas as variáveis exceto X_4 (FEA, fornecimento de energia e água). Grande parte da variância para as outras oito variáveis originais está, portanto, contida nos quatro fatores comuns.

Cargas de fator que são 0,50 ou mais (ignorando o sinal) estão sublinhadas nas equações acima. Estas cargas grandes ou moderadas indicam como as variáveis estão relacionadas com os fatores.

Uma indesejável propriedade desta escolha de fatores é que cinco das nove varáveis X são fortemente relacionadas a dois dos fatores. Isto sugere que uma rotação de fatores pode fornecer um modelo mais simples para os dados.

As comunalidades não mudaram e os fatores são ainda não correlacionados. No entanto, esta é uma solução um pouco melhor do que a anterior, pois somente X_9 é apreciavelmente dependente de mais do que um fator.

Neste estágio, é usual tentar colocar rótulas aos fatores. É honesto dizer que isto muitas vezes requer um grau de criatividade e imaginação!

Tabela 7.3 Escores de fatores rotacionados para 30 países europeus

País	Fator 1	Fator 2	Fator 3	Fator 4
Bélgica	-0,97	-0,56	-0,10	-0,48
Dinamarca	-0,89	-0.47	-0,03	-0,67
França	-0,56	-0.78	-0.15	-0,25
Alemanha	0,05	-0,57	-0,47	0,58
Grécia	0,48	0,19	-0,23	0,02
Irlanda	0,28	-0,60	-0,36	0,03
Itália	0,25	-0,13	0,17	1,00
Luxemburgo	-0,46	-0,36	0,02	0,92
Países Baixos	-1,36	-1,56	-0,03	-2,09
Portugal	0,66	-0,45	-0,37	0,64
Espanha	0,23	-0,11	-0,09	0,93
Reino Unido	-0,50	-1,14	-0,35	-0.04
Áustria	0,18	0,05	-0.71	0,56
Finlândia	-0.78	-0,20	-0,21	-0.52
Islândia	-0.18	-0.04	-0,06	0,46
Noruega	-1,36	-0,17	0,20	-0.42
Suécia	-1,20	-0,52	0,04	-0.74
Suíça	0,12	-0,67	0,01	0,65
Albânia	3,16	-1,82	1,76	-1,78
Bulgária	0,47	1,56	-0,57	-0,65
Repúblicas Tcheca/Eslováquia	-0,26	1,45	3,12	0,44
Hungria	-1,05	1,70	2,82	-0,15
Polônia	0,97	0,71	-0,37	-0,42
Romênia	1,11	1,73	-1,69	-0,81
USSR (antiga)	0,08	2,09	-0,11	0,14
Iugoslávia (antiga)	0,13	1,48	-1,70	0,17
Cingapura	0,46	-0,32	0,03	1,08
Gibraltar	-0,05	-1,05	0,08	3,26
Malta	-1,17	0,49	-0,79	-1,31
Turquia	2,15	0,07	0,15	-0,56

Comentários Finais

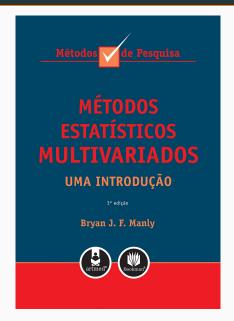
Comentários Finais

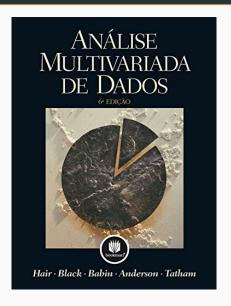
Em resumo, na aula de hoje nós:

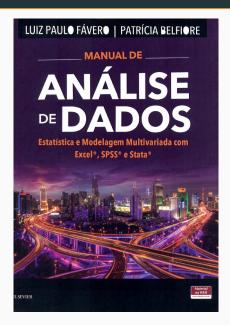
- aprendemos o que é a Análise Fatorial;
- executamos no Excel um exemplo proveniente do livro-texto.

Comentários Finais

Na próxima aula nós lidaremos com a Análise de Variância (ANOVA).







Bons Estudos!

