

Università degli studi di Ferrara
Dipartimento di Matematica
A.A. 2018/2019 – I semestre

STATISTICA MULTIVARIATA

SSD MAT/06

analisi dell'interdipendenza
LEZION 9 – Analisi fattoriale esplorativa

Docente: Valentina MINI

valentina.mini@unife.it

RICEVIMENTO: su appuntamento previa mail

Indice della lezione

1. Introduzione
2. Intuizione
3. Trattazione teorica e formalizzazione
4. Procedimento analitico da seguire nella pratica
5. Esempificazione e esercitazione in R

1 -INTRODUZIONE

Introduzione

- Abbiamo visto la ACP
- Oggi vediamo la AF metodo simile, **basata su un diverso modello** di base chiamato “*modello a fattori comuni*”
- Spesso usata per obiettivi simili a ACP
- Sottolineiamo le differenze
 - Esplicite assunzioni su come ogni variabile del dataset è misurata
 - Modello: la varianza osservata è attribuibile ad un piccolo numero di fattori comuni
- Obiettivo: identificare i fattori comuni e spiegare la loro relazione con i dati osservati

Introduzione

- **obiettivo**: ridurre il numero di variabili esplicative attraverso la creazione di nuove variabili chiamate fattori
- **Metodo**: trasformazione della struttura dei dati osservati in una nuova struttura tale che la variabilità dei dati è spiegata dai fattori

Introduzione

Paradigma comune a molte tecniche di analisi multivariata: modellare l'informazione rilevante (rappresentata da una variabile multivariata \mathbf{X}) con un numero limitato di *fattori* latenti

Esempio

In un'indagine sui consumi delle famiglie, viene registrato il livello dei consumi mensili \mathbf{X} di p beni durevoli.

La variabilità e la covarianza delle p componenti di \mathbf{X} possono essere spiegate da due o tre fattori di comportamento sociale della famiglia: il desiderio di comfort, il tentativo di raggiungere un certo livello sociale, o altri concetti sociali latenti possono spiegare la maggior parte dei comportamenti di consumo.

I sociologi sono più interessati a questi *fattori non osservabili* che alle p variabili osservate \mathbf{X} , perchè forniscono una migliore comprensione del comportamento delle famiglie.

L'analisi dei fattori interessa molti settori: psicologia, marketing, economia, etc.

2 - INTUIZIONE

Proprietà dei fattori:

- Non correlati tra loro
- Variabili latenti non osservate (sconosciute a priori) che riproducono le correlazioni esistenti tra variabili originarie
- Modello sottostante

assunzioni:

- FA può essere applicata ad un set di variabili numeriche **standardizzabili**
- Numero di **unità statistiche dovrebbe essere almeno 5 volte il numero delle variabili originarie**
(per ogni $x \rightarrow$ almeno $5 n$)

Factor Model

- X_1, \dots, X_k **variabile risposta**, tale che
 - $E(X_j) = \mu_j$, $Var(X_j) = \sigma_{jj} = \sigma_j^2$, $Cov(X_j X_r) = \sigma_{jr}$; $j, r = 1, \dots, k$
- $X_j = \lambda_{j1} F_1 + \lambda_{j2} F_2 + \dots + \lambda_{js} F_s + \dots + \lambda_{jq} F_q + U_j + \mu_j$,
 $= \sum_s \lambda_{js} F_s + U_j + \mu_j$, $j=1, \dots, k$
 - $\lambda_{j1}, \lambda_{j2}, \dots, \lambda_{jq}$ ($j=1, \dots, k$): **parametero** (constanti) definito **peso fattoriale**
 - F_1, F_2, \dots, F_q **fattori comuni** (random variables)
 - U_j , **fattore unico o specifico** ($j=1, \dots, k$)

Modello fattoriale

- **assunzioni:**

- ✓ $E(F_s)=0,$ $s=1, \dots, q$
- ✓ $Var(F_s)=1,$ $s=1, \dots, q$
- ✓ $Cov(F_s, F_t)=0,$ $s, t=1, \dots, q; s \neq t$

- ✓ $E(U_j)=0,$ $j=1, \dots, k$
- ✓ $Var(U_j)=\sigma_{jj} = \sigma_j^2$ $j=1, \dots, k$
- ✓ $Cov(U_j, U_r)=0$ $j, r=1, \dots, k; j \neq r$

- ✓ $Cov(F_s, U_j)=0$ $s=1, \dots, q; j=1, \dots, k$

Modello fattoriale

Rappresentazione matriciale:

- $\mathbf{X} = [X_1, \dots, X_k]'$ Vettore di variabili risposta
- $\mathbf{F} = [F_1, \dots, F_q]'$ Vettore di fattori comuni
- $\mathbf{U} = [U_1, \dots, U_k]'$ Vettore di fattori unici
- $\Lambda = [\lambda_{js}]$ $k \times q$ matrice di costanti (parameters)
- $\boldsymbol{\mu} = [\mu_1, \dots, \mu_k]'$ Vettore delle medie

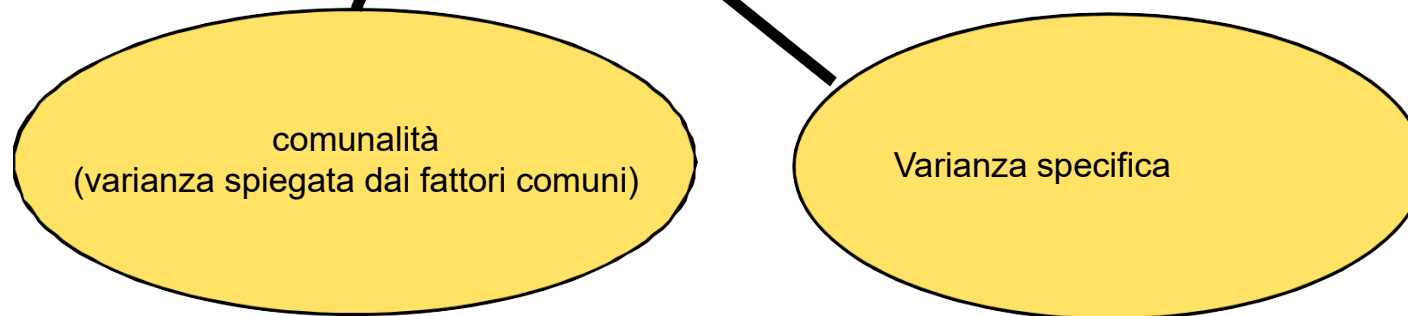
- $\mathbf{X} = \Lambda \mathbf{F} + \mathbf{U} + \boldsymbol{\mu}$

- $E(\mathbf{X}) = \boldsymbol{\mu}$, $\text{Var}(\mathbf{X}) = \Sigma = [\sigma_{jr}]$
- $E(\mathbf{F}) = \mathbf{0}$, $\text{Var}(\mathbf{F}) = I = \text{diag}(1, 1, \dots, 1)$
- $E(\mathbf{U}) = \mathbf{0}$, $\text{Var}(\mathbf{U}) = \text{diag}({}_u\sigma_{11}, \dots, {}_u\sigma_{kk}) = {}_u\Sigma$
- $\text{Cov}(\mathbf{F}, \mathbf{U}) = \mathbf{0}$

Modello fattoriale

Scomposizione della varianza (VD):

$$\begin{aligned} \text{Var}(X_j) &= \sum_s \lambda_{js}^2 + \sigma_{jj} \\ \sigma_{jj} &= h_j^2 + \sigma_{jj}, \quad j=1, \dots, k \end{aligned}$$



$\lambda_{js} = E(X_j, F_s) = \text{Cov}(X_j, F_s) \rightarrow$ misura della dipendenza lineare tra X_j e F_s

Con notazione matriciale: $\Sigma = \Lambda\Lambda' + \sigma_u^2 \Sigma$

Modello fattoriale

- FA applicabile a variabili numeriche standardizzabili
- Il numero di unità deve essere almeno 5 volte il numero delle variabili originarie: $n \geq 5 \times k$
- I fattori comuni devono spiegare almeno il 70% della variabilità totale delle variabili originarie
- Il problema di identificazione di \mathbf{F} e Λ non ha una unica soluzione

Stima dei parametri e rotazione fattoriale FA

Se le assunzioni del modello sono rispettate per F , allora una rotazione di F definisce nuovi fattori F^* per i quali le assunzioni sono ancora vere e per i quali si hanno diversi pesi fattoriali (factor loadings) Λ^* .

Formalmente:

Data la matrice ortogonale $q \times q = G$ (tale per cui $GG' = I$)

$$\begin{aligned} X &= \Lambda F + U + \mu = \\ &= \Lambda GG' F + U + \mu = \\ &= (\Lambda G)(G' F) + U + \mu = \\ &= \Lambda^* F^* + U + \mu \end{aligned}$$

Stima dei parametri e rotazione fattoriale FA

- Metodi di rotazione fattoriale:
 - **Varimax**: rotazione ortogonale
 - Equimax
 - Quartimax
 - ...

3 -TEORIA E FORMALIZZAIZONE

Obiettivo analisi fattoriale

Come nell'ACP, l'intento dell'AF è quello di *ridurre la dimensione dei dati* osservati

La prospettiva però è diversa: si assume che esista un *modello*, detto **Modello Fattoriale**

Il modello assume che le covarianze tra le p variabili di \mathbf{X} possano essere spiegate tramite un numero limitato di fattori latenti.

Obiettivo dell'AF

Costruire un modello statistico che *spieghi* la correlazione tra le variabili osservate in termini di uno o più *fattori latenti*.

Obiettivo analisi fattoriale

Esamineremo per primo il **Modello fattoriale ortogonale**, mostrando che non esiste un'unica soluzione.

Mostreremo come trarre vantaggio da questa non unicità della soluzione per derivare delle tecniche che consentano di ottenere dei risultati più interpretabili.

Vedremo che queste tecniche utilizzano delle **rotazioni** (geometriche) dei fattori.

Mostreremo varie tecniche di stima del modello e definiremo una procedura di rotazione ottimale.

Seguiremo un approccio empirico.

Specificazione del modello fattoriale

Obiettivo dell'AF è descrivere la **matrice di covarianza** delle p variabili in \mathbf{X} in termini di uno o più **fattori inosservabili**.

I fattori sono interpretati come **caratteristiche** latenti (non osservate) **comuni** delle $\mathbf{x} \in \mathbb{R}^p$.

Sia \mathbf{X} un vettore casuale con p componenti, di media $\boldsymbol{\mu}$ e matrice di covarianza $\boldsymbol{\Sigma}$

Il modello fattoriale assume che le \mathbf{X} siano linearmente dipendenti

da (pochi) fattori latenti F_1, \dots, F_m detti **fattori comuni**, p **fattori specifici**, o errori, $\varepsilon_1, \dots, \varepsilon_p$

P.e. \mathbf{X} può essere il vettore relativo a p punteggi di un test di intelligenza, che hanno come fattore latente comune il livello generale di intelligenza.

Nel marketing, le \mathbf{X} possono essere p item di un questionario sul livello di soddisfazione dei consumatori, spiegabili da fattori latenti comuni quali il livello di attrattività del prodotto, l'immagine della marca, ecc.

Modello fattoriale

$$X_1 = \mu_1 + l_{11}F_1 + \dots + l_{1m}F_m + \varepsilon_1$$

$$X_2 = \mu_2 + l_{21}F_1 + \dots + l_{2m}F_m + \varepsilon_2$$

⋮

$$X_p = \mu_p + l_{p1}F_1 + \dots + l_{pm}F_m + \varepsilon_p$$

- l_{kj} sono i *factor loadings*
- F_j sono i *fattori comuni*
- ε_k sono i *fattori specifici*

Modello AF in forma matriciale

$$\begin{aligned}\mathbf{X} &= \boldsymbol{\mu} + \mathbf{L}\mathbf{F} + \boldsymbol{\varepsilon} \\ \mathbf{X} - \boldsymbol{\mu} &= \mathbf{L}\mathbf{F} + \boldsymbol{\varepsilon}\end{aligned}$$

- \mathbf{X} è il vettore matrice $p \times 1$ delle v.c., con vettore delle medie $\boldsymbol{\mu}$
- \mathbf{L} è la matrice $p \times m$ dei coefficienti dei fattori, detti *factor loadings*
- \mathbf{F} è il vettore $m \times 1$ dei *fattori comuni*
- $\boldsymbol{\varepsilon}$ è il vettore $p \times 1$ dei *fattori specifici*, ciascuno associato ad una sola variabile

- Ci sono $m + p$ fattori non osservabili \Rightarrow il modello AF non è direttamente verificabile dai dati.
- \Rightarrow fare ipotesi su \mathbf{F} e $\boldsymbol{\varepsilon}$ che implicano una particolare struttura di $\boldsymbol{\Sigma}$ da verificare tramite i dati.

Assunzioni sui fattori latenti

Nel modello AF ortogonale con m fattori comuni, si fanno le seguenti assunzioni

- I fattori comuni
 - hanno media nulla $E(\mathbf{F}) = \mathbf{0}$ e varianza unitaria
 - sono incorrelati $E(\mathbf{F}\mathbf{F}') = \mathbf{I}$
 - sono incorrelati con i fattori specifici $Cov(\varepsilon, \mathbf{F}) = E(\varepsilon\mathbf{F}') = \mathbf{0}$
- i fattori specifici
 - hanno media nulla $E(\varepsilon) = \mathbf{0}$
 - sono incorrelati con matrice di covarianza

$$\Psi = \begin{pmatrix} \Psi_1 & 0 & \dots & 0 \\ 0 & \Psi_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \Psi_p \end{pmatrix}$$

Modello fattoriale ortogonale

$$\underset{(p \times 1)}{\mathbf{X}} = \underset{(p \times 1)}{\boldsymbol{\mu}} + \underset{(p \times m)}{\mathbf{L}} \underset{(m \times 1)}{\mathbf{F}} + \underset{(p \times 1)}{\boldsymbol{\varepsilon}}$$

μ_k = media della k -ma variabile

ε = k -mo fattore specifico

F_j = j -mo fattore comune

ℓ_{kj} = loading della k -ma variabile sul j -mo fattore

con $j = 1, \dots, m, k = 1, \dots, p$.

- I vettori dei fattori latenti F e ε soddisfano le ipotesi:
 - F e ε sono indipendenti
 - $E(\mathbf{F}) = \mathbf{0}, Cov(\mathbf{F}) = \mathbf{I}$
 - $E(\boldsymbol{\varepsilon}) = \mathbf{0}, Cov(\boldsymbol{\varepsilon}) = \boldsymbol{\Psi} = diag\{\Psi_j\}$
- Il modello e le ipotesi implicano una particolare struttura della matrice di covarianza $\boldsymbol{\Sigma}$ di \mathbf{X}

Struttura della matrice di covarianza

Il modello AF ortogonale implica una particolare struttura della matrice di covarianza Σ

- Ricordando il modello $\mathbf{X} - \boldsymbol{\mu} = \mathbf{LF} + \boldsymbol{\varepsilon}$ possiamo scrivere:

$$\begin{aligned}(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})' &= (\mathbf{LF} + \boldsymbol{\varepsilon})(\mathbf{LF} + \boldsymbol{\varepsilon})' \\ &= \mathbf{LF}(\mathbf{LF})' + \boldsymbol{\varepsilon}(\mathbf{LF})' + \mathbf{LF}\boldsymbol{\varepsilon}' + \boldsymbol{\varepsilon}\boldsymbol{\varepsilon}'\end{aligned}$$

- Considerando i valori attesi:

$$\begin{aligned}\Sigma &= E[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})'] \\ &= \mathbf{L}E[\mathbf{FF}']\mathbf{L}' + E[\boldsymbol{\varepsilon}(\mathbf{F})'](\mathbf{L}') + \mathbf{L}E[\mathbf{F}\boldsymbol{\varepsilon}'] + E[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}'] \\ &= \mathbf{LL}' + \boldsymbol{\Psi}\end{aligned}$$

(ricordare che per il modello AF ortogonale $Cov(\mathbf{F}\boldsymbol{\varepsilon}') = \mathbf{0}$)

Matrice di covarianza dal modello di analisi fattoriale descrittiva

$$\Sigma = \mathbf{LL}' + \Psi = \begin{pmatrix} \sum_{j=1}^m \ell_{1j}^2 + \psi_1 & \sum_{j=1}^m \ell_{1j}\ell_{2j} & \cdots & \sum_{j=1}^m \ell_{1j}\ell_{pj} \\ \sum_{j=1}^m \ell_{2j}\ell_{1j} & \sum_{j=1}^m \ell_{2j}^2 + \psi_2 & \cdots & \sum_{j=1}^m \ell_{2j}\ell_{pj} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{j=1}^m \ell_{pj}\ell_{1j} & \sum_{j=1}^m \ell_{pj}\ell_{2j} & \cdots & \sum_{j=1}^m \ell_{pj}^2 + \psi_p \end{pmatrix}$$

- Il modello fattoriale spiega la maggior parte della **varianza** di \mathbf{X} attraverso un piccolo numero di fattori latenti \mathbf{F} comuni alle sue p componenti: $Var(X_k) = \sum_{j=1}^m \ell_{kj}^2 + \psi_k$,
- e spiega completamente la **covarianza** tra le X_k :
 $Cov(X_k, X_s) = \sum_{j=1}^m \ell_{kj}\ell_{sj}$.
- I fattori specifici consentono di aggiustare il modello per catturare la **variabilità residua** ψ_j , non spiegata dai fattori comuni.

Osservazioni

- Il modello fattoriale si basa sulle assunzioni specificate
- Se le assunzioni non sono valide, l'analisi fornisce dei risultati spuri.
- Sebbene ACP e FA sembrano simili, la loro natura è molto diversa. Le CP:
 - sono *trasformazioni lineari* delle X
 - costruite in modo tale da avere varianza massima
 - e con l'obiettivo di ridurre la dimensione dei dati.
- Nell'**analisi fattoriale**
 - si cerca di spiegare le variazioni di X utilizzando una trasformazione lineare di un numero fisso, limitato, di *fattori latenti*.
 - con l'obiettivo di trovare i *loadings* ℓ_{kj} e le varianze specifiche ψ_j
 - Le stime di L e Ψ sono dedotte dalla struttura di covarianza
$$\Sigma = LL' + \Psi$$

Interpretazione dei fattori

Se un modello fattoriale con m fattori risulta essere ragionevole, cioè se i fattori spiegano la maggior parte della (co)varianza delle p misure considerate, risulta naturale chiedersi cosa rappresentino questi fattori.

Per interpretare i fattori si considerano le correlazioni tra le variabili osservate X_k e i fattori F_j .

osserviamo che in base al modello si ha

$$(\mathbf{X} - \boldsymbol{\mu})\mathbf{F}^j = (\mathbf{L}\mathbf{F} + \boldsymbol{\varepsilon})\mathbf{F}^j = \mathbf{L}\mathbf{F}\mathbf{F}^j + \boldsymbol{\varepsilon}\mathbf{F}^j$$

quindi la covarianza tra fattori latenti e variabili osservate è pari a:

$$\text{Cov}(\mathbf{X}, \mathbf{F}) = \mathbf{L}E(\mathbf{F}\mathbf{F}^j) + E(\boldsymbol{\varepsilon}\mathbf{F}^j) = \mathbf{L}$$

da cui si deriva la correlazione $\text{Corr}(\mathbf{X}, \mathbf{F}) = \mathbf{D}^{-1/2}\mathbf{L}$, con $\mathbf{D} = \text{diag}\{\sigma_{kk}\}$

Ipotesi di linearità

L'ipotesi di *linearità* è fondamentale nel modello fattoriale classico.

- Il modello $\mathbf{X} - \boldsymbol{\mu} = \mathbf{LF} + \boldsymbol{\varepsilon}$ è *lineare* rispetto ai fattori comuni.
- Se le X sono legate ai fattori latenti ma la relazione è NON lineare, p.e. $X_1 - \mu_1 = l_{11}F_1F_3 + \varepsilon_1$, allora la struttura di covarianza $\mathbf{LL} + \boldsymbol{\Psi}$ non è adeguata

Il concetto e la funzione di comunalità

Il modello scompone la varianza di X_k in due quote

$$\underbrace{\sigma_{kk}}_{\text{Var}(X_k)} = \underbrace{l_{k1}^2 + l_{k2}^2 + \dots + l_{km}^2}_{\text{comunalità}} + \underbrace{\psi_k}_{\text{specificità}}$$

- la proporzione di varianza di X_k spiegata dagli m fattori comuni è detta k -ma **comunalità** e si indica con h_k^2
- la proporzione di varianza dovuta al fattore specifico ψ_k è detta **specificità**, o unicità
- $\sigma_{kk} = h_k^2 + \psi_k$
- vediamo un esempio numerico

Scelta del numero dei parametri

Σ ha p varianze e $\frac{p(p-1)}{2}$ covarianze
↓

in Σ ci sono in tutto $p + \frac{p(p-1)}{2} = \frac{2p+p^2-p}{2} = \frac{p(p+1)}{2}$ parametri

il modello fattoriale ha

$p \times m$ parametri in \mathbf{L}

p specificità in Ψ

totale parametri modello fattoriale: $m \times p + p = p(m + 1)$

Numero di variabili = numero di fattori

Caso particolare: se $m = p$ la matrice Σ può essere riprodotta esattamente da $\mathbf{L}\mathbf{L}'$ con $\Psi = \mathbf{0}$

dimostrazione

- Poichè Σ è definita positiva, vale la scomposizione spettrale

$$\Sigma = \sum_{k=1}^p \lambda_k \mathbf{e}_k \mathbf{e}_k' = \mathbf{P}\mathbf{\Lambda}\mathbf{P}'$$

dove $\mathbf{\Lambda} = \text{diag}\{\lambda_k\}$, $\mathbf{P}' = [\mathbf{e}_1, \dots, \mathbf{e}_p]$

- ponendo $\mathbf{L} = \mathbf{P}\mathbf{\Lambda}^{1/2}$, con $\mathbf{\Lambda}^{1/2} = \text{diag}\{\sqrt{\lambda_k}\}$
- si ottiene

$$\Sigma = \mathbf{L}\mathbf{L}' = \mathbf{P}\mathbf{\Lambda}^{1/2}\mathbf{\Lambda}^{1/2}\mathbf{P}' = \mathbf{P}\mathbf{\Lambda}\mathbf{P}'$$

Trade-off tra completezza e modello parsimonioso

Se $m < p$ si ha una descrizione *parsimoniosa* di Σ .

Esempio:

$p = 6$ variabili, $m = 2$ fattori comuni

Σ ha 6 varianze e $(6 \times 5)/2 = 15$ covarianze $\Rightarrow 15 + 6 =$
21 parametri

$LL^j + \Psi$ ha $m \times p = 2 \times 6 = 12$ *loadings* e $p = 6$ specificità
 $\Rightarrow 12 + 6 =$ **18 parametri**

Non sempre è possibile trovare $m < p$ fattori comuni
che rappresentino la struttura di covarianza $\Sigma = LL^j$
 $+ \Psi$

quindi non sempre esiste una soluzione

Quando $m > 1$

Quando $m > 1$ la soluzione del modello fattoriale è sempre *ambigua*.

- Sia \mathbf{T} una matrice ortogonale, cioè: $\mathbf{T}\mathbf{T}' = \mathbf{T}'\mathbf{T} = \mathbf{I}$
- Il modello fattoriale $\mathbf{X} - \boldsymbol{\mu} = \mathbf{L}\mathbf{F} + \boldsymbol{\varepsilon}$ può essere trasformato

$$\mathbf{X} - \boldsymbol{\mu} = \mathbf{L}\mathbf{F} + \boldsymbol{\varepsilon} = \mathbf{L}\mathbf{T}\mathbf{T}'\mathbf{F} + \boldsymbol{\varepsilon} = \mathbf{L}^*\mathbf{F}^* + \boldsymbol{\varepsilon}$$

- dove $\mathbf{L}^* = \mathbf{L}\mathbf{T}$ e $\mathbf{F}^* = \mathbf{T}'\mathbf{F}$
- Questa trasformazione lascia immutate le proprietà del modello, infatti si ha
 - $E(\mathbf{F}^*) = \mathbf{T}'E(\mathbf{F}) = \mathbf{0}$
 - $Cov(\mathbf{F}^*) = \mathbf{T}'Cov(\mathbf{F})\mathbf{T} = \mathbf{T}'\mathbf{T} = \mathbf{I}$



Soluzione indeterminata

È impossibile scegliere tra i *factor loadings* \mathbf{L} e \mathbf{L}^* solo sulla base delle osservazioni $\mathbf{X} \Leftrightarrow$ i fattori comuni \mathbf{F} hanno le stesse proprietà statistiche dei fattori trasformati $\mathbf{F}^* = \mathbf{T}'\mathbf{F}$

non soluzione unica per $m > 1$

I fattori comuni \mathbf{F} hanno le stesse proprietà statistiche dei fattori trasformati $\mathbf{F}^* = \mathbf{T}'\mathbf{F}$

- i *loadings* \mathbf{L} sono in generale diversi dai *loadings* trasformati $\mathbf{L}^* = \mathbf{L}\mathbf{T}$
- ma generano lo *stesso modello* per la matrice di covarianza $\boldsymbol{\Sigma}$!

$$\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}' + \boldsymbol{\Psi} = \mathbf{L}\underbrace{\mathbf{T}\mathbf{T}'}_{\mathbf{I}}\mathbf{L}' + \boldsymbol{\Psi} = (\mathbf{L}^*)(\mathbf{L}^*)' + \boldsymbol{\Psi}$$

- Questa molteplicità di soluzioni possibili giustifica la *rotazione* dei fattori, perchè la trasformazione attraverso una matrice ortogonale corrisponde a una rotazione delle coordinate per \mathbf{X}

Valori fattoriali e comunalità

I *factor loadings* o *valori fattoriali* \mathbf{L} sono determinati attraverso la matrice ortogonale \mathbf{T}

$$\Rightarrow \mathbf{L} \text{ e } \mathbf{L}^* = \mathbf{L}\mathbf{T}$$

\mathbf{L} e \mathbf{L}^* forniscono la stessa rappresentazione di Σ

anche le comunalità generate dalle possibili soluzioni trasformate sono le stesse, infatti le comunalità sono date dalla diagonale di $\mathbf{L}^*(\mathbf{L}^*)'$, e per qualunque matrice ortogonale \mathbf{T} vale:

$$\mathbf{L}^*(\mathbf{L}^*)' = \mathbf{L}\mathbf{T}\mathbf{T}'\mathbf{L}' = \mathbf{L}\mathbf{L}'$$

Come si
procede →

- 1 si impongono dei *vincoli* per ottenere una soluzione unica per \mathbf{L} e Ψ
- 2 Si *ruota* la matrice dei *factor loadings* \mathbf{L} , moltiplicando $\hat{\mathbf{L}}$ per una matrice ortogonale \mathbf{T} : $\mathbf{L}^* = \mathbf{L}\mathbf{T}$
- 3 si sceglie \mathbf{T} in modo da ottenere fattori più facili da interpretare
- 4 Ottenuti i *loadings* e le *specificità*, i fattori risultano identificati e si possono stimare i punteggi fattoriali (*scores*).

Stima del modello di analisi fattoriale

Date n osservazioni p -dimensionali $\mathbf{x}_1, \dots, \mathbf{x}_n$ l'analisi fattoriale cerca di rispondere alla domanda:

Il modello fattoriale ortogonale con un numero $q < p$ può rappresentare adeguatamente i dati?

Da un punto di vista statistico, si tratterà di verificare se vale la relazione

$$\text{Cov}(\mathbf{X}) = \mathbf{\Sigma} = \mathbf{L}\mathbf{L}' + \mathbf{\Psi}$$

- La matrice campionaria \mathbf{S} è uno stimatore di $\mathbf{\Sigma}$
- Il modello fattoriale non è appropriato se le covarianze in \mathbf{S} sono piccole, o analogamente se le correlazioni in \mathbf{R} sono vicine a zero: in tal caso le specificità assumono il peso maggiore e non è possibile determinare pochi fattori comuni rilevanti.
- Se la matrice $\mathbf{\Sigma}$ ha covarianze significativamente diverse da zero, il modello fattoriale è appropriato. Si tratterà prima di tutto di stimare i factor loadings e le specificità.

4 – PROCEDIMENTO ANALITICO in R

ESEMPIO

Label	Description	Coding
ID	Personal ID of the interviewed	Increasing integer number
AgeClass	Age of the person	Age (years)
AGE_CLASS	Age class of the person	1-6
SEX	Sex of the person	M or F
PROV	Province where the interviewed lives	Province code
LIKE_WINE	How much do you like drinking wine?	Integer number from 1 to 7
FREQ_HOME	How often do you drink wine <u>at home</u> with meals?	Integer number from 1 to 5
FREQ_BAR	How often do you drink wine <u>in bars/pubs</u> ?	Integer number from 1 to 5
FREQ_REST	How often do you drink wine <u>at restaurants</u> with meals?	Integer number from 1 to 5
KNOW_PAS	Do you know the wine Passito?	Integer number from 1 to 7
FREQ_PAS	How often do you drink Passito?	Integer number from 1 to 5
FREQ_P_HOL	How often do you drink Passito on holidays and celebrations?	Integer number from 1 to 5
FREQ_P_ALO	How often do you drink Passito when you are alone?	Integer number from 1 to 5
FREQ_P_MEA	How often do you drink Passito at the end of meals?	Integer number from 1 to 5
FREQ_P_OFF	How often do you drink Passito offered by someone?	Integer number from 1 to 5
HOW_MUCH	How much wine do you drink in one year?	Integer number from 1 to 4
LIKE_PAS	How much do you like drinking Passito?	Integer number from 1 to 7
LIKE_AROMA	How much do you like aroma and smell of Passito?	Integer number from 1 to 7
LIKE_SWEET	How much do you like the sweetness of Passito?	Integer number from 1 to 7
LIKE_ALCOHOL	How much do you like the alcohol content of Passito?	Integer number from 1 to 7
LIKE_TASTE	How much do you like the intensity of taste of Passito?	Integer number from 1 to 7
PRICE	How much could you pay for one bottle of Passito? (0.5 litre)	Integer number from 1 to 5

5 - ESERCIZI IN R

Problem 1 – Passito

- Eseguire una ANALISI FATTORIALE ESPLORATIVA sulle 17 variabili risposta che rappresentano il questionario sulle abitudini, il comportamento e le preferenze dei consumatori di vino (dalla variabile LIKE_WINE alla variabile PRICE) per identificare $q < 17$ nuove variabili che “spiegano” i dati

Problem 2 – centro commerciale

- Eseguire una ANALISI FATTORIALE ESPLORATIVA sulle 5 variabili risposta per individuare $q < 5$ nuove variabili che “spiegano” i dati

Problem 3 – abitudini alimentari

- Eseguire una ANALISI FATTORIALE ESPLORATIVA sulle 12 variabili risposta osservate (da *Alcoholic Beverages a Milk*) per individuare $q < 12$ nuove variabili che “spiegano” i dati