

L'Analisi in Componenti Principali (a cura del Dott. Mario Candela):

L'analisi fattoriale è un insieme di tecniche aventi tutte la medesima logica, cambia l'approccio poiché cambia il tipo di dati da analizzare, nel nostro caso immaginiamo una matrice di dati quantitativi.

L'obiettivo di una tecnica dell'analisi fattoriale è di ridurre la dimensionalità del problema, quindi schiacciare lo spazio di analisi in uno spazio a dimensioni ridotte, più facilmente interpretabile.

Il motivo principale per cui nasce l'analisi fattoriale è di costruire indicatori complessi, l'idea è che abbiamo una serie di variabili statistiche, queste ultime esprimono parte della stessa informazione, ovvero sono tra loro correlate.

La domanda a cui risponde l'analisi fattoriale è: "Possiamo eliminare queste variabili iniziali, le quali risultano inefficienti, esprimendo parte della stessa informazione, con una sola variabile o con 2 o 3 variabili? Ovvero, possiamo ridurre il numero di variabili da p variabili molto numerose a un numero ridotto di variabili, esprimendo comunque la medesima informazione (conservando la struttura dei dati) e conservando la variabilità del problema?"

La variabilità è informazione e dobbiamo spiegarla.

Questo fa l'analisi fattoriale.

Quando le nostre variabili sono numeriche, parlando di metodi fattoriali, facciamo riferimento all'Analisi in Componenti Principali.

Il problema, dunque, consiste nel trovare i pesi migliori, e questi ultimi sono quelli che ci danno la migliore ψ .

Alla fine in un metodo statistico si deve trovare la soluzione sulla base di un criterio di ottimo, il criterio di ottimo nel nostro caso è che vogliamo conservare quanta più informazione contenuta nei dati, ovverossia la variabilità; quindi, vogliamo che la nostra variabile di sintesi conservi quanto più possibile la variabilità contenuta nella nube originaria dei punti.

Significa che se 2 individui hanno valori di diversi di x_1 e x_2 , avranno anche punteggi diversi della variabile di sintesi ψ ; invece, se hanno valori simili di x devono avere punteggi simili della ψ .

La 1^a variabile latente che tiriamo fuori coi metodi di analisi fattoriale è sempre quella nella direzione di massimo allungamento della nube, ossia quella che cattura la massima variabilità.

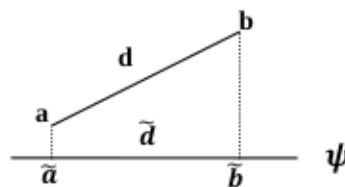
Quindi, andiamo a trovare quei pesi (u) che fanno sì che la varianza della nostra ψ sia massima

$$Var(\psi) = max$$

, questo è il criterio-guida nella scelta dei pesi; così fa l'ACP su variabili quantitative, mentre per l'ACM si tratta di variabili qualitative.

Nell'analisi fattoriale, nel momento in cui schiacciamo lo spazio, necessariamente un'azione di riduzione della dimensionalità (schiacciamento dello spazio del problema) comporta una distorsione nelle distanze tra gli individui, ossia tra le posizioni precedenti rispetto a quelle che otteniamo nello spazio proiettato.

Ad esempio, consideriamo 2 punti, a e b , e d è la distanza originaria tra questi 2 punti; proiettando tali punti per ottenere la variabile di sintesi, ψ , la distanza sarà \bar{d} .



Nell'ACP le distanze si misurano con la distanza euclidea (trattandosi di uno spazio euclideo, considerando il segmento più breve che unisce i 2 punti); dunque, la distanza originaria rappresenta sempre l'ipotenusa, mentre la distanza proiettata è uno dei 2 cateti e, per tale motivo, è più piccola dell'ipotenusa e, quindi, della distanza originaria.

La distanza proiettata è sempre più piccola della distanza originaria.

Massimizzare la varianza di ψ significa la medesima cosa del volere che la distanza proiettata sia la più grande possibile, ossia quanto più vicina a quella originaria, conservando più o meno le posizioni degli individui originari.

Quindi, il problema possiamo vederlo da un punto di vista geometrico, ossia la misura di ψ che vogliamo è quella che schiaccia meno la distanza tra gli individui; il che corrisponde al volere la misura di ψ che conserva la variabilità, poiché schiacciando gli individui diminuisce la variabilità.

Ad esempio, se \mathbf{a} fosse il punto medio, dunque la distanza tra \mathbf{a} e \mathbf{b} è la varianza, perché la distanza si calcola come $(\text{valori di } \mathbf{b} - \text{valori di } \mathbf{a})^2$ se ragioniamo in termini di media, mentre la distanza tra \mathbf{a} e \mathbf{b} è la devianza se ragioniamo in termini di somma; dunque, la distanza tra $\tilde{\mathbf{a}}$ e $\tilde{\mathbf{b}}$ è la varianza o la devianza dei punti proiettati.

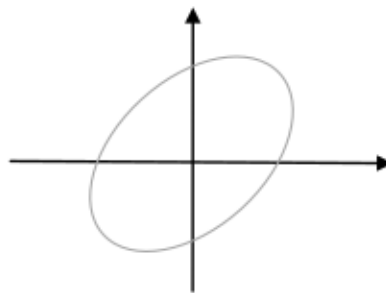
Nel momento in cui vogliamo massimizzare la variabilità di $\boldsymbol{\psi}$, vogliamo massimizzare il segmento $\tilde{\mathbf{d}}$, ossia le differenze di ogni punto proiettato rispetto alla media; se, invece, di confrontare un punto confrontiamo tutti i punti rispetto alla media è chiaro che significhino la stessa cosa.

Il teorema di Huygens, difatti, prevede che massimizzare la distanza proiettata, facendo sì che la variabile conservi le differenze tra gli individui, è la medesima cosa del dire che la sua variabilità è massima.

Nell'ACP, siccome vogliamo fare la somma pesata di variabili, è necessario che le variabili abbiano la medesima unità di misura, ossia tali variabili devono avere la stessa media e stessa varianza (se abbiamo una variabile con varianza molto grande e una variabile con varianza molto piccola, non possiamo confrontare le variazioni delle variabili).

Per poter confrontare variabili, per poter fare una somma delle variabili pesate le variabili devono essere standardizzate, eliminando cioè l'unità di misura di partenza trasformandola in un'unità di misura standard, in σ -*esimi*; fatto ciò le variabili sono confrontabili perché hanno tutte media 0 e varianza 1.

Quindi, quando ragioniamo nell'ACP la 1ª cosa che facciamo è che le \mathbf{x} di partenza sono tutte standardizzate, e avendo tali variabili tutte media 0, la nube è centrata nell'origine degli assi, poiché la media è l'origine



Se facciamo la distanza dall'origine stiamo calcolando la varianza.

Nell'ACP vogliamo massimizzare la varianza di $\boldsymbol{\psi}$ e, dunque, vogliamo trovare quegli \mathbf{u} che massimizzano la varianza di $\boldsymbol{\psi}$; inoltre, quest'ultima è la somma pesata delle variabili originarie per i pesi, in forma matriciale

$$\Psi = \mathbf{x} \cdot \mathbf{u}$$

e la varianza di $\boldsymbol{\psi}$ è uguale a

$$\begin{aligned} \text{Var}(\boldsymbol{\psi}) &= E(\boldsymbol{\psi} - \text{media}(\boldsymbol{\psi}))^2 = \\ &= [(\mathbf{x} \cdot \mathbf{u}) - 0]^2 = \\ &= (\mathbf{x} \cdot \mathbf{u})^2 \end{aligned}$$

Alla fine il criterio che vogliamo massimizzare, in forma matriciale, è

$$\max \text{Var}(\boldsymbol{\psi}) = \mathbf{u}' \cdot \mathbf{x}' \mathbf{x} \cdot \mathbf{u}$$

ossia vogliamo trovare quelle \mathbf{u} che moltiplicate prima e dopo per $\mathbf{x}'\mathbf{x}$, dove \mathbf{x} è la matrice dei dati (ossia tutte le nostre variabili o i nostri item), siano massime.

Abbiamo \mathbf{x} , la matrice dei dati originari, \mathbf{n} individui e \mathbf{x}_p variabili originarie, queste ultime sono standardizzate, avendo

	\mathbf{x}_1	\mathbf{x}_2	\mathbf{x}_p
s_1	$\frac{(x_1 - \mu_1)}{\sigma_1}$...	$\frac{(x_p - \mu_1)}{\sigma_1}$
s_2
s_n	$\frac{(x_1 - \mu_n)}{\sigma_n}$...	$\frac{(x_p - \mu_n)}{\sigma_n}$

Quando facciamo questa matrice per la sua trasposta, moltiplichiamo riga per colonna, ad es. 1^a riga per 1^a colonna; quindi, la matrice che ne deriva è una matrice in cui la dimensione sarà $\mathbf{p} \times \mathbf{p}$ dove le \mathbf{p} sono le variabili.

Il prodotto matriciale si fa 1^a riga per 1^a colonna, ogni elemento l'uno per l'altro facendo, poi, la somma; quindi, nella cella avremo la sommatoria

	\mathbf{x}_1	\mathbf{x}_2	\mathbf{x}_p
\mathbf{x}_1	$\sum \left[\frac{(x_1 - \mu_1)}{\sigma_1} \cdot \frac{(x_1 - \mu_1)}{\sigma_1} \right] = 1$	$\sum \left[\frac{(x_1 - \mu_1)}{\sigma_1} \cdot \frac{(x_2 - \mu_2)}{\sigma_2} \right]$...
\mathbf{x}_2	...	1	...
\mathbf{x}_p	1

Questa è la matrice di correlazione.

$\left[\frac{(x_1 - \mu_1)}{\sigma_1} \cdot \frac{(x_1 - \mu_1)}{\sigma_1} \right]$ è il coefficiente di correlazione di Pearson se avessimo 2 variabili diverse, \mathbf{x}_1 e \mathbf{x}_2 ; difatti, la formula del coefficiente di correlazione di Pearson, al numeratore abbiamo il prodotto degli scarti e al denominatore il prodotto delle deviazioni standard

$$\rho = \frac{(x - \mu_x) \cdot (y - \mu_y)}{\sigma_x \cdot \sigma_y} \Rightarrow \text{coeff. correlazione Pearson}$$

$$(x - \mu_x) \cdot (y - \mu_y) = \sigma_{XY} \text{ \{covarianza tra X e Y\}}$$

Ovviamente, se la variabile è la medesima la correlazione è pari a 1, se le variabili sono diverse abbiamo, ad esempio, $\rho(\mathbf{x}_1, \mathbf{x}_2)$.

Tutto il problema dell'analisi fattoriale è di trovare quei pesi che, moltiplicati per la matrice di correlazione, ci danno la varianza di ψ ; difatti, i pesi ricostruiscono la matrice di correlazione.

In una matrice di correlazione, se il fenomeno è multi-dimensionale, ci potrebbero essere più variabili latenti e non solo una, ossia avremmo più ψ ; poiché, una sola ψ può non bastare a spiegare tutta l'informazione.

Una volta trovata la 1^a ψ , potremmo voler trovare la 2^a ψ , ossia quella che massimizza la variabilità al netto della variabilità della 1^a ψ .

Aggiungiamo, dunque, un altro vincolo nel nostro massimo, che la 2^a ψ resti incorrelata con la 1^a ψ , ovvero che il prodotto del 1^o vettore di pesi col 2^o vettore di pesi sia pari a 0; perché, rispetto alle variabili originarie tali variabili devono essere efficienti (efficienti in questo caso vuol dire incorrelate), ovvero ogni ψ deve spiegare parti diverse dell'informazione.

Da questo sistema viene fuori una serie di vettori di pesi fino a \mathbf{p} , perché potremmo trovare tante variabili ψ quante sono le variabili originarie; l'obiettivo è trovarne meno però, se continuassimo a massimizzare, potremmo trovarne tante quante le variabili originarie, tuttavia le ψ sono efficienti rispetto alle variabili originarie perché sono incorrelate tra loro.

Ogni ψ è un vettore di pesi, il miglior vettore di pesi è quello che massimizza la variabilità; quindi, a questo vettore di pesi si aggiunge un'informazione, il coefficiente λ , che è la varianza di ψ

$$\lambda = \text{Var}(\psi)$$

Immaginiamo ad esempio,

Vettori pesi	
\mathbf{u}_1	$\lambda_1 > \lambda_2$
\mathbf{u}_2	$\lambda_2 > \lambda_3$
\mathbf{u}_3	$\lambda_3 > \lambda_4$
...	...
\mathbf{u}_p	λ_p

Avendo tutti i vettori di pesi, il 1° vettore di pesi, quello migliore, che ci dà la 1ª ψ (la più importante), è quello che ha λ più grande, ossia la varianza più grande. Questa è la 1ª soluzione.

Se la 1ª ψ spiega tutta l'informazione interessante nella nube ci fermiamo, trovando una variabile latente; invece, se 1ª ψ non spiega tutta l'informazione, prendiamo il 2° vettore e otteniamo un'altra λ , ottenendo la 2ª componente principale (ψ_1), poi possiamo ottenere la 3ª e così via.

Per scegliere quante componenti principali ci servono per spiegare una parte sufficiente della informazione contenuta nella nube, consideriamo che la somma dei λ è nota e pari a p ; poiché, se ogni variabile originaria ha varianza 1 e abbiamo 10 variabili, la varianza totale della nube originaria è 10.

Se costruiamo tante variabili latenti ψ quante sono le variabili originarie, ricostruiamo tutta l'informazione della nube, nel caso dell'esempio con \mathbf{x}_1 e \mathbf{x}_2 possiamo trovare 2 ψ , la 1ª ψ è quella della massima direzione di allungamento della nube e la 2ª ψ spiega il resto dell'informazione e, essendo incorrelata con la 1ª, forma un angolo retto con la 1ª componente principale; quindi, la 2ª possibile ψ è per forza perpendicolare, perché abbiamo sostituito le 2 variabili originarie, tra loro correlate e, dunque, inefficienti, con 2 variabili efficienti. Probabilmente, per spiegare l'informazione interessante contenuta nel fenomeno basta una sola ψ .

La somma delle λ fa p perché se le λ sono standardizzate e hanno varianza pari a 1, la loro somma, nel nostro esempio, è pari a 10.

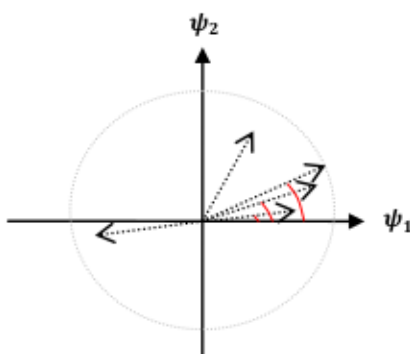
Immaginiamo che $p = 10$, ossia abbiamo 10 variabili originarie, e che $\lambda_1 = 7$, ciò vuol dire che il coefficiente λ_1 sintetizza tanta informazione quanto prima era contenuta in 7 variabili originarie, dunque da sola ψ_1 spiega il 70% dell'informazione originaria ($\frac{7}{10}$); se non ci basta il 70%, prenderemo anche la variabile latente successiva ψ_2 , ad es. $\lambda_2 = 2$, spiega 2 volte quanto spiegava una singola variabile, ottenendo così il 90% dell'informazione originaria ($\frac{2}{10} + \frac{7}{10}$), però originariamente servivano 10 variabili (ossia 10 indicatori statistici).

Il 90% raggiunto con sole 2 variabili di sintesi ci soddisfa, quindi vuol dire che alle variabili \mathbf{x} originarie sostituiamo queste 2.

Questo è un criterio.

Un altro criterio è di prendere le variabili che hanno $\lambda > 1$, siccome λ è la varianza di ψ e la varianza di una variabile originaria è 1, tutti i $\lambda > 1$ vuol dire che le ψ sono variabili che sintetizzano più di una variabile originaria, mentre le ψ che hanno $\lambda < 1$ sintetizzano meno di una variabile originaria, dunque non sono efficienti.

Per l'ACP abbiamo un grafico di questo tipo,



osserviamo che le variabili originarie \mathbf{x}_1 , \mathbf{x}_2 , e \mathbf{x}_3 sono tra loro fortemente correlate, avendo tutte la medesima direzione, e sono correlate con ψ_1 ; quindi, valori alti di ψ_1 stanno sintetizzando \mathbf{x}_1 , \mathbf{x}_2 , e \mathbf{x}_3 .

Anche un'altra variabile, nell'altro quadrante, è fortemente correlata con ψ_1 , però in maniera negativa; quindi, gli individui che hanno un valore elevato di ψ_1 hanno bassi valori di ψ_2 e viceversa.

Le variabili originarie che formano un angolo di quasi 90° con ψ_1 sono incorrelate con quest'ultima e sono sintetizzate da ψ_2 .

Gli angoli che le variabili formano sono le correlazioni, sono i pesi (u) che diamo alle variabili, ossia quanto queste variabili sono correlate rispetto alla variabile di sintesi ψ ; difatti, quella con la maggiore correlazione avrà maggior peso.

Dal grafico dell'ACP vogliamo 2 informazioni, l'angolo ci dice il grado di correlazione con ψ e la lunghezza del segmento di una variabile ci dice quanto di quella variabile è sintetizzata da ψ .