

## BRIEF ARTICLE

THE AUTHOR

### 1. ITALIAN STUFF USED AS DRAFT

Si tratta di un metodo di analisi che è stato introdotto nel 1933 da Harold Hotelling [?] per analizzare dati di tipo psicometrico. Nella letteratura scientifica in lingua inglese tale metodo prende il nome di *Principal Component Analysis*, il cui acronimo PCA è molto popolare. Il modo più compatto di introdurre la PCA è attraverso una serie di considerazioni di tipo algebrico svolte sulla matrice di varianza-covarianza. Per prima cosa richiamiamo la definizione della matrice di varianza-covarianza e cerchiamo di riscriverla in modo simmetrico:

Tale identità algebrica mette in luce la prima proprietà della matrice di varianza-covarianza, ovvero il fatto che si tratta di una matrice simmetrica – un fatto di cui ci siamo già accorti tramite la definizione dei suoi elementi di matrice, si veda l'Eq. (??). Infatti, se introduciamo la matrice  $\mathbf{A} = \chi / (N_o - 1)$ , avremo per la matrice  $\mathbf{\Sigma}$ :

$$(1) \quad \mathbf{\Sigma} = \mathbf{A}^t \mathbf{A}$$

Valutando la trasposta di  $\mathbf{\Sigma}$  otteniamo:

$$(2) \quad \mathbf{\Sigma}^t = (\mathbf{A}^t \mathbf{A})^t = \mathbf{A}^t (\mathbf{A}^t)^t = \mathbf{A}^t \mathbf{A} = \mathbf{\Sigma},$$

che dimostra ancora una volta la simmetria della matrice di varianza-covarianza. Un'altra proprietà significativa della matrice di varianza-covarianza è che si tratta di una matrice definita positiva<sup>1</sup>, che quindi possiede autovalori  $\geq 0$  [?].

Vediamo ora in che modo si possa introdurre una trasformazione lineare che permetta di passare dagli indici di variabile  $v$  a nuovi indici  $s$ , che chiameremo indici delle componenti principali. L'idea è quella di sostituire al set di variabili con cui abbiamo operato per raccogliere le osservazioni un nuovo set di variabili "migliori", per due motivi: (1) la loro mutua indipendenza (ortogonalità) e (2) il loro numero, potenzialmente più ridotto del numero di variabili originarie. Tale set di variabili "migliori" è il set delle cosiddette componenti principali, da associare agli indici  $s$ . Detto in altri termini, *tramite le componenti principali diventa possibile descrivere le osservazioni sperimentali utilizzando un set ridotto di grandezze tra di loro indipendenti*, con vantaggi evidenti.

---

<sup>1</sup>Questo risultato è semplice da dimostrare utilizzando la forma  $\mathbf{\Sigma} = \mathbf{A}^t \mathbf{A}$ . Una matrice semi definita positiva  $\mathbf{M}$  è tale per cui dato un vettore generico  $\mathbf{x}$  risulta  $\mathbf{x}^t \mathbf{M} \mathbf{x} \geq 0$ . Nel nostro caso, sfruttando la definizione  $\mathbf{\Sigma} = \mathbf{A}^t \mathbf{A}$  abbiamo  $\mathbf{x}^t \mathbf{\Sigma} \mathbf{x} = \mathbf{x}^t \mathbf{A}^t \mathbf{A} \mathbf{x} = (\mathbf{A} \mathbf{x})^t (\mathbf{A} \mathbf{x}) = |\mathbf{A} \mathbf{x}|^2$ . L'ultimo passaggio dimostra che otteniamo il modulo al quadrato di un generico vettore  $\mathbf{A} \mathbf{x}$ ; tale modulo al quadrato è certamente una grandezza  $\geq 0$ .

Dal punto di vista tecnico, per implementare quest'idea basta utilizzare la decomposizione spettrale della matrice di covarianza, che è agevole perchè la matrice  $\Sigma$  è simmetrica. A quel punto si cercherà di riportare la decomposizione spettrale in una forma equivalente a quella dell'Eq. 1 in cui tuttavia il primo membro sarà in forma diagonale. Il riconoscimento di una struttura del tipo  $A^t A$  nel secondo membro fornirà la definizione del dataset nelle nuove variabili semplificatrici, ovvero nelle componenti principali  $s$ . È chiaro che, il requisito di *indipendenza* delle componenti principali risulterà soddisfatto se assoceremo a ciascuna di questa uno degli autovettori della matrice di covarianza. Infatti autovettori diversi saranno ortogonali tra di loro (*vide infra*) e forniranno i coefficienti della trasformazione lineare che permette di passare dalle variabili di indice  $v$  ad una certa componente principale di indice  $s$ .

Esaminiamo quindi i dettagli algebrici in modo dettagliato. Cominciamo con l'equazione agli autovalori per la matrice di covarianza  $\Sigma$ , in cui associamo ogni coppia autovalore/autovettore ad un indice di componente principale  $s$ :

$$(3) \quad \Sigma_{vv} L_{vs} = L_{vs} \sigma_{ss}$$

Nell'equazione (3)  $\sigma_{ss}$  è la matrice diagonale degli autovalori di  $\Sigma_{vv}$ .  $L_{vs}$  è la matrice degli autovettori di  $\Sigma_{vv}$ . Dato che  $\Sigma_{vv}$  è simmetrica,  $L_{vs}$  è una matrice ortogonale, per la quale valgono le relazioni:

$$(4) \quad \begin{aligned} L_{vs} L_{sv} &= \mathbf{1}_{vv} \\ L_{sv} L_{vs} &= \mathbf{1}_{ss}. \end{aligned}$$

Moltiplicando da sinistra l'Eq. (3) per la matrice  $L_{sv}$ , e sfruttando le sue proprietà di ortogonalità, si ricava la *decomposizione spettrale* della matrice di varianza-covarianza:

$$(5) \quad L_{sv} \Sigma_{vv} L_{vs} = \sigma_{ss}$$

Introduciamo ora nel secondo membro dell'Eq. (5) la definizione di  $\Sigma_{vv} = \chi_{vo} \chi_{ov} / (N_o - 1)$  (cfr. Eq. ??):

$$(6) \quad \sigma_{ss} = \frac{1}{N_o - 1} L_{sv} \chi_{vo} \chi_{ov} L_{vs}$$

Con un piccolo sforzo possiamo riconoscere al secondo membro una struttura data dal prodotto di una nuova matrice  $S$  e della sua trasposta:

$$(7) \quad \sigma_{ss} = \left[ \frac{1}{\sqrt{N_o - 1}} L_{sv} \chi_{vo} \right] \left[ \chi_{ov} L_{vs} \frac{1}{\sqrt{N_o - 1}} \right] = S_{so} S_{os} = S^t S.$$

Le righe di tale matrice, detta matrice degli *scores* ( $S = S_{os}$ ), definiscono le osservazioni  $o$  in funzione delle componenti principali  $s$ :

$$(8) \quad S_{os} \equiv \left[ \frac{1}{\sqrt{N_o - 1}} \chi_{ov} L_{vs} \right]$$

La diagonalizzazione della matrice di covarianza fornisce l'insieme degli autovalori delle componenti principali, rappresentati in quello che viene solitamente chiamato *screeplot*. Tale grafico fornisce l'andamento decrescente delle varianze principali ( $\sigma_s$ ) in funzione

dell'indice  $s$  della componente principale. Da questo grafico si può rapidamente giudicare l'importanza relativa delle componenti principali nel descrivere la varianza totale dei dati del dataset.

La matrice  $\mathbf{L}_{vs}$  degli autovettori della matrice di covarianza  $\mathbf{\Sigma}_{vv}$  è detta matrice dei *loadings*. Essa rappresenta il legame tra le variabili ( $v$ ) e le componenti principali ( $s$ ). Nel caso di dataset formati da spettri, la rappresentazione dei *loadings* di una componente principale in funzione delle variabili (e.g., lunghezze d'onda), è affine a quella di uno spettro. I picchi (positivi o negativi) nella rappresentazione di uno dei *loadings* mostrano in quali regioni spettrali sono osservate le maggiori variazioni (crescita/decrecita) del segnale all'interno del dataset.

Per quanto riguarda invece la matrice  $\mathbf{S}$  degli *scores*, il grafico cartesiano delle prime due colonne della matrice  $\mathbf{S}_{os}$  (pensate come coordinate  $x, y$ ) fornisce la posizione delle osservazioni del dataset rispetto al sistema di riferimento ortogonale fornito dalle componenti principali  $s_1, s_2$ . Tale grafico è talvolta definito *scatterplot* (grafico di dispersione dei dati). È utile osservare che talvolta si utilizza una matrice di *scores*  $\mathbf{S}'_{os}$  normalizzata rispetto alle deviazioni standard principali:

$$(9) \quad \mathbf{S}'_{os} \equiv \mathbf{S}_{os} \boldsymbol{\sigma}_{ss}^{-1/2}$$

In questo modo nella presentazione dello *scatterplot* ci si può ricondurre ad un sistema di riferimento adimensionale (standardizzato). Infatti, le varianze calcolate rispetto alla matrice  $\mathbf{S}'_{os}$  risultano date da un prodotto uguale alla matrice identità:

$$(10) \quad (\mathbf{S}')^t \mathbf{S}' = \left( \boldsymbol{\sigma}_{ss}^{-1/2} \mathbf{S}_{so} \right) \left( \mathbf{S}_{os} \boldsymbol{\sigma}_{ss}^{-1/2} \right) = \boldsymbol{\sigma}_{ss}^{-1/2} \boldsymbol{\sigma}_{ss} \boldsymbol{\sigma}_{ss}^{-1/2} = \mathbf{1}.$$