

# A BRIEF GUIDE TO PRINCIPAL COMPONENT ANALYSIS AS IMPLEMENTED IN OPENPCA

MATTEO TOMMASINI

The Principal Component Analysis (PCA) was introduced in 1933 by Harold Hotelling [?] in the context of psicometric data analysis. Since its birth, PCA has been widely applied to many fields where multivariate datasets have to be dealt with. The easier approach to introduce PCA, by also taking into consideration its numerical implementation in Matlab, is through a fully algebraic approach that focuses on the variance-covariance matrix of the dataset and its spectral decomposition. Let us introduce first the multivariate dataset matrix  $\mathbf{X}_{ov}$ , which along each row stores the results of one multivariate observation along a given number of variables ( $N_v$ ). The adopted notation for the dataset matrix highlights the different role of row *vs.* column indexes. The different observations are identified in the  $\mathbf{X}_{ov}$  matrix by the row index ( $o$ ), whereas the different variables of each multivariate measurement (observation) are identified by the column index ( $v$ ). In the context of spectroscopy, each row represents one spectrum, and the different variables are the wavenumbers at which the instrument has recorded a given spectral intensity (*e.g.*, Raman intensity, or absorbance). Hence, because of the adopted notation, we have the following identities:

$$(1) \quad \begin{aligned} \mathbf{X} &= \mathbf{X}_{ov} \\ \mathbf{X}_{vo} &= (\mathbf{X}_{ov})^t = \mathbf{X}^t, \end{aligned}$$

where  $^t$  indicates matrix transposition. As described later, the variance-covariance matrix among the variables of the dataset can be straightforwardly introduced through the matrix of the centered dataset,  $\chi_{ov}$ :

$$(2) \quad \chi_{ov} = \mathbf{X}_{ov} - \mathbf{X}_{\langle o \rangle v},$$

where  $\mathbf{X}_{\langle o \rangle v}$  represents the row vector of the average values of the variables over the number of  $N_o$  observations, and its elements are given by:

$$(3) \quad X_{\langle o \rangle v} = \frac{1}{N_o} \sum_o X_{ov}$$

We adopt in Eq. (2) the same abuse of notation used in Matlab: by subtracting a row vector to a matrix actually one subtracts the given row vector to each row of the matrix. Hence Eq. (2) is implemented in Matlab as simply as `chi = X - mean(X)`, because the Matlab function `mean(X)` gives the row vector corresponding to the average of all the rows of the X matrix – which effectively corresponds to averaging out with respect to the available observations (see above). The variance-covariance matrix among the variables

of the dataset  $(\mathbf{\Sigma}_{vv})$  is defined as follows:

$$(4) \quad \Sigma_{v_1 v_2} = \frac{1}{N_o - 1} \sum_o (X_{ov_1} - X_{\langle o \rangle v_1}) (X_{ov_2} - X_{\langle o \rangle v_2}) = \\ = \frac{1}{N_o - 1} \sum_o \chi_{ov_1} \chi_{ov_2}$$

The last equality shows that, with matrix notation, the variance covariance matrix is simply given by:

$$(5) \quad \mathbf{\Sigma}_{vv} = \frac{1}{N_o - 1} (\mathbf{\chi}_{ov})^t \mathbf{\chi}_{ov} = \frac{1}{N_o - 1} \mathbf{\chi}_{vo} \mathbf{\chi}_{ov}.$$

Clearly, by definition,  $\mathbf{\Sigma}$  is a symmetric matrix, and it is positive definite<sup>1</sup>. Therefore it admits spectral decomposition by the orthogonal matrix of its eigenvectors, and the eigenvalues are positive quantities [?]. The matrix eigenvalue problem of the variance-covariance matrix is written as:

$$(6) \quad \mathbf{\Sigma}_{vv} \mathbf{L}_{vs} = \mathbf{L}_{vs} \mathbf{\sigma}_{ss}$$

In Eq. (6)  $\mathbf{\sigma}_{ss}$  is the diagonal matrix of the eigenvalues of  $\mathbf{\Sigma}_{vv}$  and  $\mathbf{L}_{vs}$  is the orthogonal matrix of the eigenvectors of  $\mathbf{\Sigma}_{vv}$ . The orthogonality of  $\mathbf{L}_{vs}$  implies:

$$(7) \quad \mathbf{L}_{vs} \mathbf{L}_{sv} = \mathbf{1}_{vv} \\ \mathbf{L}_{sv} \mathbf{L}_{vs} = \mathbf{1}_{ss}.$$

Therefore, by left-multiplying Eq. (14) by  $\mathbf{L}_{sv}$ , and by considering its orthonormality, one obtains the spectral decomposition of the variance-covariance matrix:

$$(8) \quad \mathbf{L}_{sv} \mathbf{\Sigma}_{vv} \mathbf{L}_{vs} = \mathbf{\sigma}_{ss}$$

By substituting in the right-hand side of Eq. (16) the definition of  $\mathbf{\Sigma}_{vv} = \mathbf{\chi}_{vo} \mathbf{\chi}_{ov} / (N_o - 1)$  (cfr. Eq. ??), one obtains:

$$(9) \quad \mathbf{\sigma}_{ss} = \frac{1}{N_o - 1} \mathbf{L}_{sv} \mathbf{\chi}_{vo} \mathbf{\chi}_{ov} \mathbf{L}_{vs}$$

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

$$(10) \quad \mathbf{\sigma}_{ss} = \left[ \frac{1}{\sqrt{N_o - 1}} \mathbf{L}_{sv} \mathbf{\chi}_{vo} \right] \left[ \mathbf{\chi}_{ov} \mathbf{L}_{vs} \frac{1}{\sqrt{N_o - 1}} \right] = \mathbf{S}_{so} \mathbf{S}_{os} = \mathbf{S}^t \mathbf{S}.$$

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

$$(11) \quad \mathbf{S}_{os} \equiv \left[ \frac{1}{\sqrt{N_o - 1}} \mathbf{\chi}_{ov} \mathbf{L}_{vs} \right]$$

---

<sup>1</sup>This is straightforward to show by using the expression  $\mathbf{\Sigma} = \mathbf{A}^t \mathbf{A}$ , with  $\mathbf{A} = \mathbf{\chi} / \sqrt{N_o - 1}$ . For  $\mathbf{\Sigma}$  to be positive defined it should be  $\mathbf{x}^t \mathbf{\Sigma} \mathbf{x} \geq 0$  for a generic  $\mathbf{x}$  column vector. By considering that  $\mathbf{\Sigma} = \mathbf{A}^t \mathbf{A}$ , such requirements becomes  $0 \leq \mathbf{x}^t \mathbf{\Sigma} \mathbf{x} = \mathbf{x}^t \mathbf{A}^t \mathbf{A} \mathbf{x} = (\mathbf{A} \mathbf{x})^t (\mathbf{A} \mathbf{x}) = |\mathbf{A} \mathbf{x}|^2$ . The last step clearly proves the statement, because it represents the squared modulus of the generic vector  $\mathbf{A} \mathbf{x}$  (certainly a non-negative quantity).

## 1. STUFF

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  $\mathbf{A} = \boldsymbol{\chi}/(N_o - 1)$ , avremo per la matrice  $\boldsymbol{\Sigma}$ :

$$(12) \quad \boldsymbol{\Sigma} = \mathbf{A}^t \mathbf{A}$$

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

$$(13) \quad \boldsymbol{\Sigma}^t = (\mathbf{A}^t \mathbf{A})^t = \mathbf{A}^t (\mathbf{A}^t)^t = \mathbf{A}^t \mathbf{A} = \boldsymbol{\Sigma},$$

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>2</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.

Dal punto di vista tecnico, per implementare quest'idea basta utilizzare la decomposizione spettrale della matrice di covarianza, che è agevole perchè la matrice  $\boldsymbol{\Sigma}$  è simmetrica. A quel punto si cercherà di riportare la decomposizione spettrale in una forma equivalente a quella dell'Eq. 12 in cui tuttavia il primo membro sarà in forma diagonale. Il riconoscimento di una struttura del tipo  $\mathbf{A}^t \mathbf{A}$  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  $s$ .

---

<sup>2</sup>Questo risultato è semplice da dimostrare utilizzando la forma  $\boldsymbol{\Sigma} = \mathbf{A}^t \mathbf{A}$ . 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  $\boldsymbol{\Sigma} = \mathbf{A}^t \mathbf{A}$  abbiamo  $\mathbf{x}^t \boldsymbol{\Sigma} \mathbf{x} = \mathbf{x}^t \mathbf{A}^t \mathbf{A} \mathbf{x} = (\mathbf{A} \mathbf{x})^t (\mathbf{A} \mathbf{x}) = |\mathbf{A} \mathbf{x}|^2$ . 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  $\geq 0$ .

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  $s$ :

$$(14) \quad \Sigma_{vv} \mathbf{L}_{vs} = \mathbf{L}_{vs} \sigma_{ss}$$

Nell'equazione (14)  $\sigma_{ss}$  è la matrice diagonale degli autovalori di  $\Sigma_{vv}$ .  $\mathbf{L}_{vs}$  è la matrice degli autovettori di  $\Sigma_{vv}$ . Dato che  $\Sigma_{vv}$  è simmetrica,  $\mathbf{L}_{vs}$  è una matrice ortogonale, per la quale valgono le relazioni:

$$(15) \quad \begin{aligned} \mathbf{L}_{vs} \mathbf{L}_{sv} &= \mathbf{1}_{vv} \\ \mathbf{L}_{sv} \mathbf{L}_{vs} &= \mathbf{1}_{ss}. \end{aligned}$$

Moltiplicando da sinistra l'Eq. (14) per la matrice  $\mathbf{L}_{sv}$ , e sfruttando le sue proprietà di ortogonalità, si ricava la *decomposizione spettrale* della matrice di varianza-covarianza:

$$(16) \quad \mathbf{L}_{sv} \Sigma_{vv} \mathbf{L}_{vs} = \sigma_{ss}$$

Introduciamo ora nel secondo membro dell'Eq. (16) la definizione di  $\Sigma_{vv} = \chi_{vo} \chi_{ov} / (N_o - 1)$  (cfr. Eq. ??):

$$(17) \quad \sigma_{ss} = \frac{1}{N_o - 1} \mathbf{L}_{sv} \chi_{vo} \chi_{ov} \mathbf{L}_{vs}$$

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

$$(18) \quad \sigma_{ss} = \left[ \frac{1}{\sqrt{N_o - 1}} \mathbf{L}_{sv} \chi_{vo} \right] \left[ \chi_{ov} \mathbf{L}_{vs} \frac{1}{\sqrt{N_o - 1}} \right] = \mathbf{S}_{so} \mathbf{S}_{os} = \mathbf{S}^t \mathbf{S}.$$

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

$$(19) \quad \mathbf{S}_{os} \equiv \left[ \frac{1}{\sqrt{N_o - 1}} \chi_{ov} \mathbf{L}_{vs} \right]$$

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  $\Sigma_{vv}$  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:

$$(20) \quad \mathbf{S}'_{os} \equiv \mathbf{S}_{os} \boldsymbol{\sigma}_{ss}^{-1/2}$$

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:

$$(21) \quad (\mathbf{S}')^t \mathbf{S}' = \left( \boldsymbol{\sigma}_{ss}^{-1/2} \mathbf{S}_{so} \right) \left( \mathbf{S}_{os} \boldsymbol{\sigma}_{ss}^{-1/2} \right) = \boldsymbol{\sigma}_{ss}^{-1/2} \boldsymbol{\sigma}_{ss} \boldsymbol{\sigma}_{ss}^{-1/2} = \mathbf{1}_{ss}.$$