

Γραμμική Αλγεβρα II

PCA
(παράδειγμα)

15 Μαΐου 2017

Μας δίνονται τα παρακάτω δεδομένα.

Βάρος (lb)	120	125	125	135	145
Υψος (inch)	61	60	64	68	72

Να γίνει η ανάλυση κυρίων συνιστωσών (PCA).

Λυση:

Ευκολά μπορούμε να υπολογίσουμε το διάνυσμα των μέσων ως εξής

$$M = \frac{1}{5} \begin{bmatrix} \sum x_i \\ \sum y_i \end{bmatrix} = \begin{bmatrix} 130 \\ 65 \end{bmatrix}$$

Μετά αφαιρώντας τον μέσο από κάθε x_i έχουμε

$$\hat{x}_1 = x_1 - M = \begin{bmatrix} -10 \\ -4 \end{bmatrix}, \quad \hat{x}_2 = x_2 - M = \begin{bmatrix} -5 \\ -5 \end{bmatrix}$$

$$\hat{x}_3 = x_3 - M = \begin{bmatrix} -5 \\ -1 \end{bmatrix}, \quad \hat{x}_4 = x_4 - M = \begin{bmatrix} 5 \\ 3 \end{bmatrix}, \quad \hat{x}_5 = x_5 - M = \begin{bmatrix} 15 \\ 7 \end{bmatrix},$$

Αρα τα δεδομένα μας γίνονται ως εξής

$$B = \begin{bmatrix} 10 & -5 & -5 & 5 & 15 \\ -4 & -5 & -1 & 3 & 7 \end{bmatrix}.$$

Απ' τον πίνακα B προκύπτει ο πίνακας διακύμανσης ως εξής

$$\Sigma = \frac{1}{N-1} BB^T = \begin{bmatrix} 100 & 47.5 \\ 47.5 & 25 \end{bmatrix}$$

οπου παρατηρούμε ότι η συνολική διακύμανση είναι ίση με 125. Αρα η μεταβλητή βάρους έχουμε ότι περιγράφει το $\frac{100}{125}\%$ της πληροφορίας που υπάρχει στα δεδομένα ενώ η μεταβλητή υψους περιγράφει το $\frac{25}{125}\%$.

Ας αρχίσουμε τώρα την ανάλυση σε κύριες συνιστώσες, υπολογίζοντας τις ιδιοτιμές του πίνακα Σ . Έχουμε ότι

$$\det(\Sigma - \lambda I) = 0 \Rightarrow \dots \Rightarrow \lambda_1 = 123, \lambda_2 = 2$$

Έτσι προκύπτει ότι η μεταβλητή βάρους έχει την μεγαλύτερη ιδιοτιμή ($\lambda_1 = 123$) και αρα είναι η μεγαλύτερη συνιστώσα. Το αντιστοιχο ιδιοδιάνυσμα της προκύπτει ότι είναι το $u_1 = [0.9 \ 0.436]^T$. Αξίζει να σημειώσουμε ότι η συνολική διακύμανση δεν έχει μεταβληθεί ($\lambda_1 + \lambda_2 = 125$). Ενδεικτικά, αναφέρουμε ότι ένας δείκτης που θα λάμβανε υπόψιν του το βάρος και το υψός που έχουμε διαθέσιμα μέσο αυτού του δείγματος είναι ο

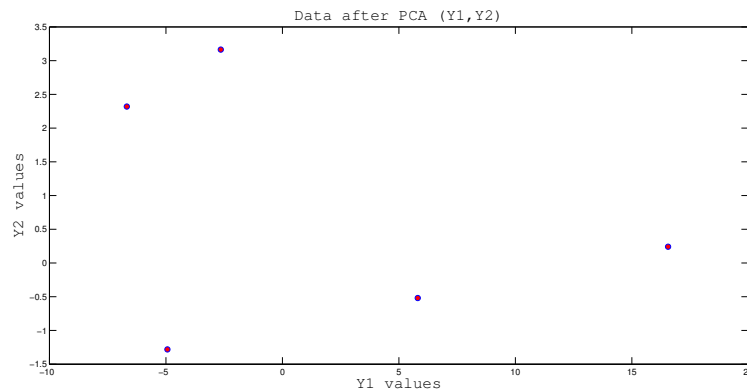
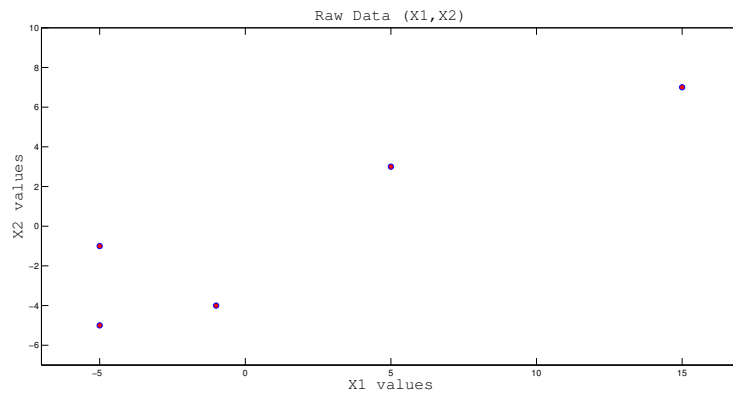
$$y = 0.9x_1 + 0.436x_2$$

Ο νέος πίνακας διακύμανσης θα είναι

$$\hat{\Sigma} = \begin{bmatrix} 123 & 0 \\ 0 & 2 \end{bmatrix}$$

Έτσι τώρα διαπιστώνουμε ότι τελικώς, μετά από την ανάλυση κυρίων συνιστωσών, η πρώτη συνιστώσα έχει το $\frac{123}{125} = 98.4\%$ της πληροφορίας ενώ η δεύτερη συνιστώσα έχει το $\frac{2}{125} = 1.6\%$ της πληροφορίας. Σημειώνουμε ότι μετά την ανάλυση κυρίων συνιστωσών, ο νέος πίνακας διακυμάνσεων, δεν περιέχει συνδιακυμάνσεις, άρα οι μεταβλητές μας έγιναν ασυσχέτιστες. Στη συνέχεια, χρησιμοποιώντας τον ορθογώνιο πίνακα P που περιέχει τα ιδιοδιανύσματα του Σ σε στήλες, θέτουμε $X = PY \rightarrow Y = P^T X$ και παίρνουμε τα νέα δεδομένα Y , τα οποία όπως είπαμε είναι ασυσχέτιστα. Τα παραπάνω φαίνονται στις εικόνες που ακολουθούν:

Σχήμα 1: Τα δεδομένα πριν την Ανάλυση κύριων συνιστωσών



Σχήμα 2: Τα δεδομένα μετά την Ανάλυση κύριων συνιστωσών

Φαίνεται καθαρά ότι τα δεδομένα έχουν χωριστεί σε δυο ομάδες που βρίσκονται πάνω σε δυο ευθείες.