

Ανάλυση κύριων συστατικών (Principal Components Analysis: PCA)
 Απο: Bishop (2006): Pattern recognition and machine Learning, Chapter 12:

- Έστω δαθείω N παρατηρήσεις x_1, \dots, x_N , που η καθεμία έχει μήκος (μ) διαστάσεων: $x_i \in \mathbb{R}^n$

Αν για παράδειγμα το x_i παριστάει την i εικόνα και καθε συνιστώσα x_{ij} του x_i ένα pixel αρα ως εικόνας, με $x_{ij} = 0$ για ασπρο και $x_{ij} = 1$ για μαύρο, τότε αν η εικόνα έχει 100×100 pixel, θα έχουμε $n = 10.000$ συνιστώσες ή κατά x_i .

- Σε κάποιες εφαρμογές τα N διανύσματα του \mathbb{R}^n , x_1, \dots, x_N , βρίσκονται "κοντά" σε έναν υπόχωρο χαμηλότερης διαστάσεως (k) του \mathbb{R}^n : τον χώρο υπόχωρο, U_k .

Από παράδειγμα $n=2, k=1$: Figure 12.2

- Τότε η "πιο κοντινή προσέγγιση" των διανυσμάτων που μπορεί να περιγραφεί "επαρκώς" από την προβολή των x_i στον χώρο υπόχωρο U_k , \tilde{x}_i . Οι ανωκλίσεις $x_i - \tilde{x}_i$ θα μπορούσαν (για κατασκευή ενδογής των k) να είναι "δύο φορές" μικρής σημασίας.

- Κρίσιμος υπόχωρος, διαστάσεως k , θα είναι εκείνος που ~~ελαχιστοποιεί~~ αναλαρισά την περισσότερη διακυμαντικότητα των x_i , δηλαδή εκείνος που ~~ελαχιστοποιεί~~ τις διακυμανσεις των \tilde{x}_i .
- Ισοδύναμα (θα δούμε): εκείνος που ελαχιστοποιεί τις ανώκλίσεις $\sum \|x_i - \tilde{x}_i\|^2$.

PCA χρησιμοποιείται για

- μείωση της διαστάσεως που είναι ανεπάρκεια για να περιγραφούν τα δεδομένα (dimensionality reduction)
- συμπίεση δεδομένων (data compression)
- εξόρυξη χαρακτηριστικών (feature extraction)
- οπτικοποίηση δεδομένων (data visualization)

Για την παρακμπλοή και την εύρεση των καλύτερων αντιστοιχιών υπάρχουν δύο συνάρτες ισοδύναμης διαστάσεως.

- Κατ' αρχάς ως υποθέσουμε ότι τα δεδομένα μας έχουν μηδενικό διαφαντικό μέσο $\bar{X} = \frac{1}{N} \sum X_i = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$. Αν όχι θα αρκεί να αφαιρεθεί το \bar{X} από τα δεδομένα μας, οπότε τα "ρέα" δεδομένα θα είναι μηδενικό διαφαντικό μέσο.

* Πρώτη διαίσθηση: Ξέρουμε υπάρχει U_k , για κάθε ένα δεδομένο διαίσθηση k , που μεγιστοποιεί τη διακύμανση των προβολών $\tilde{X}_i = P_{U_k} X_i$ των δεδομένων.

- Ας υποθέσουμε $k=1$. και $U_1 = \text{Span}(q_1)$, $\|q_1\|=1$ για να βρούμε q_1 : ~~πρέπει~~ να βρούμε ποιο σφαιρικό να είναι το q_1 .

Εστω $\tilde{X}_i = P_{q_1} X_i = (q_1^T X_i) \cdot q_1$

Τότε η "διακυματική διακύμανση" των \tilde{X}_i ορίζεται ως

ως το $\frac{1}{N} \sum_{i=1}^N \|\tilde{X}_i\|^2$, το οποίο ισούται με

$$\frac{1}{N} \sum_{i=1}^N \|\tilde{X}_i\|^2 = \frac{1}{N} \sum_{i=1}^N (X_i^T q_1)^2 = \frac{1}{N} \sum_{i=1}^N q_1^T X_i X_i^T q_1$$

$$= q_1^T S q_1 \quad \text{με} \quad S := \frac{1}{N} \sum X_i X_i^T$$

Αν φάραμε q_1 , με $\|q_1\|=1$ και να μην σκοπεύει να PCB
 $q_1^T S q_1$. Η παραφραση για το μέτρο του Rayleigh (βλ. Strang, ή Σημειώσεις) δίνει την ανάρμοση: Το μέτρο επιτυχίας είναι ότι q_1 είναι το ιδιοδιάνυσμα που αντιστοιχεί στο μέγιστο ιδιοτιμή, λ_1 , του S . Τότε $q_1^T S q_1 = \lambda_1$.

• Αν $k=2$ φάραμε U_2 και έστω τον περιγράψουμε ως $U_2 = \text{Span}(q_1, q_2)$ με q_1, q_2 ορθοκανονικά. Τότε με $\tilde{x}_i = P_{U_2} x_i$

$$\max_{U_2} \frac{1}{N} \sum \| \tilde{x}_i \|^2 = \max_{q_1, q_2 \text{ ορθοκ.}} \frac{1}{N} \sum \| (q_1^T x_i) q_1 + (q_2^T x_i) q_2 \|^2$$

$$= \max_{q_1, q_2 \text{ ορθοκ.}} \left\{ \frac{1}{N} \sum \| (q_1^T x_i) q_1 \|^2 + \frac{1}{N} \sum \| (q_2^T x_i) q_2 \|^2 \right\}$$

$$= \max_{q_1} q_1^T S q_1 + \max_{q_2 \perp q_1} q_2^T S q_2.$$

Το μέτρο αυτό επιτυγχάνεται για q_1, q_2 να είναι τα ιδιοδιάνυσμα που αντιστοιχούν στις δύο μεγαλύτερες ιδιοτιμές του S και ισούται με $\lambda_1 + \lambda_2$.

• Γενικά, $k \geq 2$: Η U_k , ο υπόχωρος που μην σκοπεύει να διατηρήσει διακρίματα των προβολών των x_i στον U_k να είναι ο $U_k = \text{Span}(q_1, \dots, q_k)$ στον q_1, \dots, q_k να είναι τα ιδιοδιάνυσμα των k μεγαλύτερων

ιδιοτιμών του S $\lambda_1 \geq \dots \geq \lambda_k$.

$$\text{Τότε } \tilde{x}_i = (q_1^T x_i) \cdot q_1 + \dots + (q_k^T x_i) q_k$$

* Δοθέντα Στοιχεία: Έχουμε υπόχρημα U_k , για αριθμούς και δεδομένα PCA.
 Διάστασης k , που ελαχιστοποιεί την τετραγωνική απόσταση
 ανάμεσα στα δεδομένα x_i και στην προβολή τους $\tilde{x}_i = P_{U_k} x_i$.
 Ορίζεται οπότε $J = \frac{1}{N} \sum_{i=1}^N \|x_i - \tilde{x}_i\|^2$.

- Έστω p_1, \dots, p_m μια ορθοκανονική βάση του \mathbb{R}^n

$$\text{Τότε } x_i = \sum_{j=1}^m (x_i^T p_j) \cdot p_j$$

- Αν $V_k = \text{Span}(p_1, \dots, p_k)$ τότε το σφαιρίδιο του V_k που
 ανήκει ελάχιστο από το x_i είναι το $\tilde{x}_i = P_{V_k} x_i$

και οπότε $\tilde{x}_i = \sum_{j=1}^k (x_i^T p_j) p_j$ και η απόσταση

$$x_i - \tilde{x}_i = \sum_{j=k+1}^m (x_i^T p_j) \cdot p_j$$

- Τότε θα έχουμε $J = \frac{1}{N} \sum_{i=1}^N \|x_i - \tilde{x}_i\|^2$

$$= \frac{1}{N} \sum_{i=1}^N \sum_{j=k+1}^m (p_j^T x_i)^2$$

$$= \frac{1}{N} \sum_{i=1}^N \sum_{j=k+1}^m p_j^T x_i x_i^T p_j$$

$$= \sum_{j=k+1}^m p_j^T S p_j, \text{ όπου } S = \frac{1}{N} \sum x_i x_i^T$$

Αρα το ζήτημα μας είναι να βρούμε μια βάση P_1, \dots, P_n του \mathbb{R}^n PCA5
 έτσι ώστε να ελαχιστοποιήσουμε k να ελαχιστοποιείται η

$$J = \sum_{j=k+1}^n p_j^T S p_j$$

• Παράδειγμα: ~~n=2~~ $n=2, k=1$. : Ελαχιστοποιούμε την $p^T S p$
 Ανασυνδυάζουμε (αυτίο του Rayleigh) για $p = q_2$ το ιδιοδιάνυσμα
 με μικρότερο ιδιοτιμή λ_2 . Τότε δε έχουμε και $J = q_2^T S q_2 = \lambda_2$
 Αρα $P_1 = q_1$ και $U_1 = \text{Span}(q_1)$, όπως έχουμε
 πριν και με τη μέθοδο ορθογώνιων διασυνδέσεων.

• Γενικός ($k > 1$): Διαλέξτε P_{k+1}, \dots, P_n να είναι τα
 ιδιοδιανύσματα του S που αντιστοιχούν στις μικρότερες
 ιδιοτιμές $\lambda_{k+1} \dots \lambda_n$. Τότε $J = \sum_{j=k+1}^n \lambda_j$
 και $U_k = \text{Span}(q_1, \dots, q_k)$ τα ιδιοδιανύσματα
 που αντιστοιχούν στις μεγαλύτερες ιδιοτιμές $\lambda_1, \dots, \lambda_k$.

Παρατηρήσεις. Αν $k=n$ δεν έχω τίποτα διασυνδέσεις
 αλλα $\tilde{x}_i = x_i$

• Σε πολλές εφαρμογές γίνονται "ομαδικά ήθη"
 ιδιοτιμές $k < n$ για να εξαχθούν οι πρώτοι
 και τα δεδομένα

→ Παρατήρηση.