

Ανάλυση Ανεξάρτητων, Θετικών Συνιστώσων

Αρχές και Μέθοδοι Μηχανικής Μάθησης

Σημειώσεις στις σημειώσεις του κ. Διαμαντάρα

Ανεξάρτητοι Παράγοντες

Από το PCA γνωρίζουμε ότι οι παράγοντες y_i είναι **Στατιστικά Ασυσχέτιστοι** μεταξύ τους. Δηλαδή ισχύει:

$$E\{y_i y_j\} = 0, i \neq j \quad (1)$$

Αλλά,

Σε κάποιες εφαρμογές είναι σωστότερο να υποθέσουμε ότι οι **κρυφοί παράγοντες** είναι **στατιστικά ανεξάρτητοι**, δηλαδή έχουμε:

$$p(y_i, y_j) = p(y_i) * p(y_j) \quad (2)$$

Και αυτό γιατί η ανεξαρτησία είναι **πιο ισχυρή** από την έλλειψη συσχέτισης:

$$y_i, y_j \text{ Ανεξάρτητοι} \Rightarrow y_i, y_j \text{ Ασυσχέτιστοι} \quad (3)$$

Δηλαδή οι ασυσχέτιστοι **δεν σημαίνει** απαραίτητα ότι είναι και ανεξάρτητοι.

$$y_i, y_j \text{ Ανεξάρτητοι} \neq y_i, y_j \text{ Ασυσχέτιστοι} \quad (4)$$

Πως μπορεί να χρησιμοποιηθεί

Έστω ότι υπάρχουν 2 παρουσιαστές σε μια σκηνή με 1 μικρόφωνο ο καθένας. Δυστυχώς, επειδή είναι πολύ κοντά, **οι φωνές και των 2 απορροφούνται και από τα 2 μικρόφωνα!** Με το ICA, θα μπορούσαμε αφού έχουμε τα μπλεγμένα δεδομένα, να τα **ξεμπλέξουμε και να ξεχωρίσουμε τις 2 φωνές.**

Ανάλυση Ανεξαρτήτων Συνιστώσων

- Δεδομένα:

- Παρατηρήσεις: $\mathbf{x}(1), \dots, \mathbf{x}(N) \in \mathbb{R}^n$
- Δεν χρησιμοποιούνται στόχοι t (Χωρίς Επίβλεψη)
- Πλήθος παραγόντων n

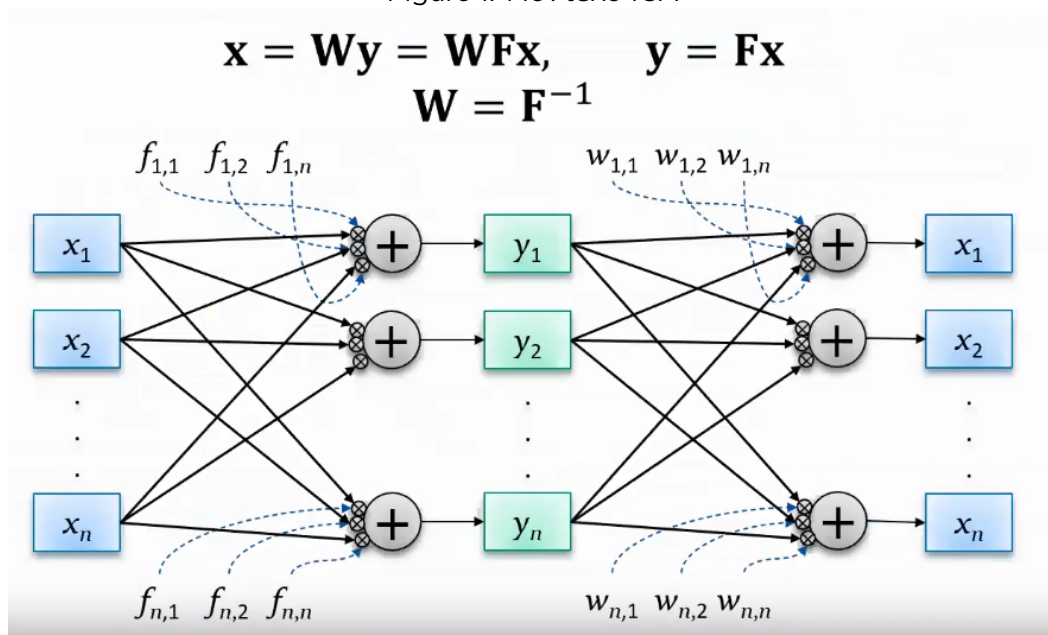
- Πρόβλημα: Βρες τον πίνακα $\mathbf{F} \in \mathbb{R}^n * n$ έτσι ώστε οι παράγοντες $y_i = \mathbf{f}_i^T \mathbf{x}$ να είναι στατιστικά ανεξάρτητοι μεταξύ τους.

- \mathbf{f}_i^T οι γραμμές του \mathbf{F}

- Παρατήρηση: Αν θέσουμε

- $\mathbf{W} = \mathbf{F}^{-1}$
- $\mathbf{y} = \mathbf{F}\mathbf{x}$
- τότε $\mathbf{x} = \mathbf{W}\mathbf{y} = w_1 y_1 + \dots + w_n y_n$

Figure 1: Μοντέλο ICA



Συσσωρεύτριες (cumulants)

Οι Συσσωρεύτριες είναι συναρτήσεις αντίθεσης.

Ορισμός 0.1. Μια συσσωρευτρία k τάξης μιας τυχαίας μεταβλητής x ορίζεται ως

$$c_k(x) = (-j)^k \frac{d^k \phi_x(\omega)}{d\omega^k} \quad (5)$$

Όπου $\phi_x(\omega) = \ln E\{e^{j\omega x}\}$

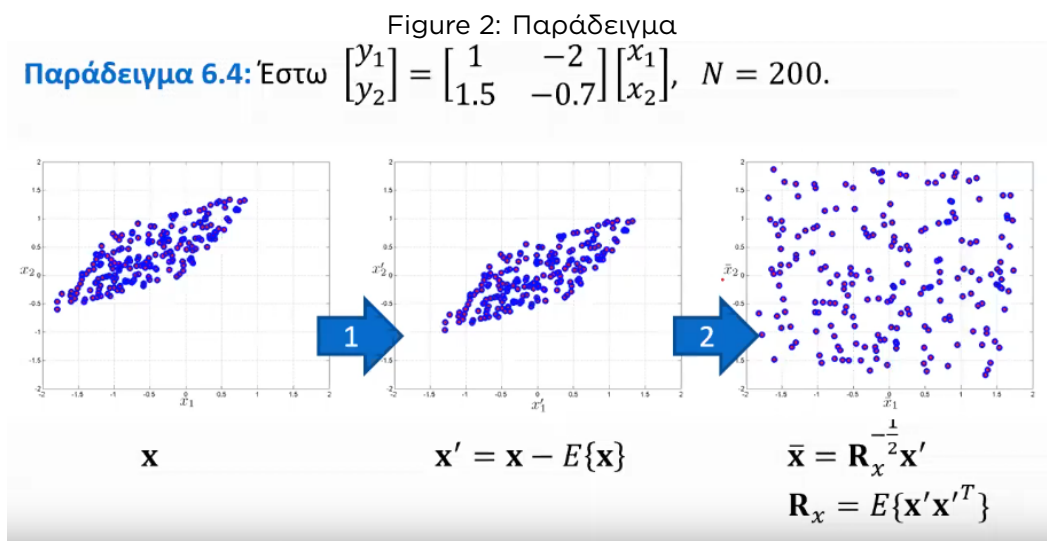
Και πρακτικά έχουμε:

- $c_1 = E\{x\}$
- $c_2 = E\{x^2\}$ (διακύμανση)
- $c_3 = E\{x^3\}$
- $c_4 = E\{x^4\} - 3E\{x^2\}^2$ (κύρτωση, αρκετά χρήσιμη)
- ...κλπ

Προεπεξεργασία Δεδομένων

1. Αφαίρεση μέσου όρου: $\mathbf{x} \Rightarrow \mathbf{x}' : E\{\mathbf{x}'\} = 0$
2. Λεύκανση: $\mathbf{x}' \Rightarrow \bar{\mathbf{x}} : E\{\bar{\mathbf{x}}\bar{\mathbf{x}}^T\} = \mathbf{I}$

Ορισμός 0.2. Λεύκανση ονομάζεται η τετραγωνοποίηση διαγράμματος μιας ομάδας δεδομένων.



Εξαγωγή ενός παράγοντα

Έστω για κάποιο $\mathbf{f} : z(k) = \mathbf{f}^T \mathbf{x}(k)$

θα έχουμε $z(k) = \mathbf{f}^T \mathbf{W} \mathbf{y}(k)$

όπου $\mathbf{f}^T = \mathbf{v}^T \Rightarrow z(k) = \mathbf{v}^T \mathbf{y}(k)$

Θεώρημα 0.1. Αν οι παράγοντες δεν είναι Γκαουσσανές τυχαίες μεταβλητές, τότε η μεγιστοποίηση της κύρτωσης

$J(\mathbf{f}) = |c_4(z)|$ υπό τον περιορισμό $c_2(z) = 1$ επιτυγχάνεται για

$\mathbf{v} = [0 \ 0 \ \dots \ 0 \ \alpha \ 0 \ \dots \ 0]$

Περιορισμοί

- Χάνεται η σειρά των παραγόντων, δηλαδή οι εκτιμώμενοι παράγοντες y_i μπορεί να είναι ανακατωμένοι.
- Οι παράγοντες μπορεί να εξαχθούν με αυθαίρετη κλιμάκωση. Πιθανή **αλλαγή προσήμου**.
- Κάθε ανάλυση ICA σε ίδια δεδομένα μπορεί να εξάγει παράγοντες με διαφορετική σειρά και κλιμάκωση.

PCA εναντίον ICA

Ζουμί	PCA	ICA
Πλήθος παραγόντων(συνιστώσων)	m (≤ διάσταση διανύσματος παρατήρησης)	n (= διάσταση διανύσματος παρατήρησης)
Σχέση παραγόντων	Ασυσχέτιστοι	Ανεξάρτητοι
Συνάρτηση Κόστους	Μέσο τετραγωνικό σφάλμα	Συναρτήσεις αντίθεσεις