



Εγγύτητα και αποστάσεις

Ιωάννης Μαδεμλής

Ομοιότητα και ανομοιότητα

- Κατά την προεπεξεργασία, ενδεχομένως να απαιτείται η αποτίμηση της **εγγύτητας** μεταξύ των στοιχείων/προτύπων/διανυσμάτων χαρακτηριστικών του συνόλου δεδομένων.
- Παράδειγμα: ορισμένοι αλγόριθμοι εξόρυξης εργάζονται μόνο με τις τιμές εγγύτητας μεταξύ των προτύπων.
 - Από τη στιγμή που αυτές υπολογιστούν, δεν είναι πλέον απαραίτητη η αποθήκευση των αρχικών δεδομένων.
- Ως εγγύτητα ορίζεται η **ομοιότητα** ή η **ανομοιότητα**:
 - Αριθμητικά μέτρα του βαθμού στον οποίον δύο πρότυπα είναι παρόμοια ή διαφέρουν, αντιστοίχως.

Εγγύτητα

Ομοιότητα και ανομοιότητα

- Σημαντική υποπερίπτωση των μέτρων ανομοιότητας είναι τα μέτρα **απόστασης**.
- Αρκετά μέτρα εγγύτητας λαμβάνουν πραγματικές τιμές στο κλειστό διάστημα $[0, 1]$.
 - 1 σημαίνει μέγιστη ομοιότητα (ή ανομοιότητα).
 - 0 σημαίνει ελάχιστη.
- Συνήθως, τα μέτρα απόστασης λαμβάνουν τιμές στο $[0, +\infty]$.
- Τι γίνεται αν το εκάστοτε τρέχον μέτρο εγγύτητας λαμβάνει τιμές σε πεπερασμένο διάστημα διάφορο του $[0, 1]$;

Εγγύτητα

Ομοιότητα και ανομοιότητα

- Αυθαίρετης κλίμακας τιμές εγγύτητας για το τρέχον σύνολο δεδομένων εύκολα μετασχηματίζονται ώστε να εμπίπτουν στο $[0, 1]$:

- **Κανονικοποίηση min-max:** έστω ότι η τρέχουσα τιμή εγγύτητας μεταξύ δύο συγκεκριμένων προτύπων είναι s , ενώ στο ολικό σύνολο δεδομένων η μέγιστη/ελάχιστη εμφανιζόμενη τιμή είναι s_{max}/s_{min} , αντιστοίχως.

- $$s' = \frac{s - s_{min}}{s_{max} - s_{min}}$$

- Συνήθως, όχι όμως πάντα, μία τιμή ομοιότητας s εξάγεται μέσω μαθηματικού μετασχηματισμού της τιμής ανομοιότητας d :

- Π.χ., $s = 1 - d$, αν η ανομοιότητα λαμβάνει τιμές στο $[0, 1]$, ή εναλλακτικά $s = -d$.

Εγγύτητα

Ομοιότητα και ανομοιότητα

- Προσοχή: όταν υπολογίζουμε το s με βάση ένα d το οποίο λαμβάνει τιμές στο $[0, +\infty]$, τότε το s ίσως να μην εμπίπτει στο $[0, 1]$.
- Αν είναι κρίσιμο το s να εμπίπτει στο $[0, 1]$, μπορεί πρώτα να εφαρμοστεί κανονικοποίηση min-max στις τιμές ανομοιότητας:
 - $s = 1 - \frac{d - d_{min}}{d_{max} - d_{min}}$.
- Είναι σημαντικό να εξασφαλίσουμε ότι δεν χάνεται ή αλλοιώνεται πολύτιμη πληροφορία.
 - Εξαρτάται από την εφαρμογή και τα δεδομένα.
 - Π.χ., με την κανονικοποίηση min-max ορισμένες εγγύτητες μεταξύ διαφορετικών προτύπων ίσως γίνουν ίσες λόγω σφαλμάτων στρογγύλευσης.

Εγγύτητα

Υπολογισμός ανομοιότητας

- Η ανομοιότητα δύο στοιχείων διάστασης 1 (με ένα μόνο γνώρισμα) μπορεί να οριστεί ως:
 - α) η **απόλυτη διαφορά** των τιμών των στοιχείων σε αυτό το γνώρισμα, αν το τελευταίο είναι *αριθμητικό* ή *διατακτικό*,
 - β) να τεθεί ανομοιότητα 1 για δύο στοιχεία με διαφορετική τιμή στο γνώρισμα και ανομοιότητα 0 για ίδια τιμή, σε περίπτωση *ονομαστικού* γνωρίσματος.
- Στην περίπτωση διατακτικού γνωρίσματος μπορεί προαιρετικά να διαιρεθεί η απόλυτη διαφορά με το πλήθος των δυνατών διαφορετικών τιμών του γνωρίσματος.
 - Έτσι, η τελική ανομοιότητα θα λαμβάνει τιμές στο διάστημα $[0, 1]$.

Εγγύτητα

Υπολογισμός ανομοιότητας

- Συνήθως, τα πρότυπα των ενδιαφερόντων συνόλων δεδομένων έχουν περισσότερες από μία διαστάσεις/γνωρίσματα.
- Τότε υπάρχουν ποικίλοι τρόποι σύνθεσης των διαφορών μεταξύ των τιμών δύο στοιχείων σε όλα τα γνωρίσματα, προκειμένου να προκύψει εν τέλει μόνο μία, βαθμωτή τιμή εγγύτητας.
- Τι γίνεται αν τα m -διάστατα πρότυπα έχουν γνωρίσματα διαφορετικού τύπου;
 - Π.χ., ονομαστικά, διατακτικά και αναλογίας;

Εγγύτητα

Υπολογισμός ανομοιότητας

- Παράδειγμα: κάθε πρότυπο περιέχει τρία γνωρίσματα ενός υπαλλήλου κάποιου οργανισμού:
 - *Αύξων αριθμός* (ονομαστικό),
 - *Μισθός* (αναλογίας),
 - *Επίπεδο μόρφωσης* (διατακτικό).
- Τότε η ανομοιότητα μεταξύ δύο προτύπων μπορεί να βρεθεί ως εξής:
 - Υπολογίζουμε ξεχωριστά την ανομοιότητά τους για κάθε γνώρισμα, και κατόπιν
 - λαμβάνουμε τη μέση τιμή των m επιμέρους ανομοιοτήτων.
 - Κατά τον υπολογισμό του μέσου όρου, μπορούμε να προβούμε σε ποικίλους, επιπρόσθετους χειρισμούς, αναλόγως με το πρόβλημα και τα δεδομένα.

Εγγύτητα

Υπολογισμός ανομοιότητας

- Κατά τον υπολογισμό του μέσου όρου:
 - Μπορούμε να αγνοήσουμε γνωρίσματα για τα οποία δεν έχει νόημα η σύγκριση μεταξύ προτύπων (π.χ., αύξων αριθμός).
 - Να αγνοήσουμε ενδεχόμενα ασύμμετρα δυαδικά γνωρίσματα όπου και τα δύο πρότυπα έχουν μηδενική τιμή.
 - Να αγνοήσουμε ενδεχόμενα γνωρίσματα όπου απουσιάζει η τιμή του ενός από τα δύο πρότυπα.
 - Αν κάποια γνωρίσματα είναι πιο σημαντικά από τα υπόλοιπα, μπορούμε να υπολογίσουμε **σταθμισμένο μέσο όρο**.
 - Τα βάρη στάθμισης των γνωρισμάτων εξαρτώνται από το εκάστοτε πρόβλημα και σύνολο δεδομένων.

Εγγύτητα

Υπολογισμός απόστασης

- Συνήθως όμως όλα τα γνωρίσματα είναι ίδιου τύπου.
- Τότε, ο πιο κοινός τρόπος υπολογισμού ανομοιότητας μεταξύ δύο πολυδιάστατων προτύπων είναι τα **μέτρα απόστασης**.
- Πρόκειται για μέτρα ανομοιότητας, για τα οποία ισχύουν τα εξής:
 - Η τιμή τους είναι πάντοτε μεγαλύτερη ή ίση του μηδενός (ίση για πανομοιότυπα διανύσματα/πρότυπα).
 - Είναι συμμετρικά: $d(\mathbf{x}, \mathbf{z}) = d(\mathbf{z}, \mathbf{x})$.
 - Ισχύει η τριγωνική ανισότητα:
 - $d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$.
 - Μπορούμε να φτάσουμε από το \mathbf{x} στο \mathbf{z} μέσω μίας παράκαμψης διά του \mathbf{y} , αντί για απευθείας, αλλά δεν υπάρχει περίπτωση έτσι η διαδρομή να γίνει συντομότερη.

Εγγύτητα

Υπολογισμός απόστασης

- Το βασικότερο παράδειγμα είναι η **απόσταση Μινκόφσκι**.
 - Κάθε απόσταση στηριγμένη στη νόρμα \mathcal{L}_p της διαφοράς δύο διανυσμάτων/προτύπων.
- Ειδικές περιπτώσεις της απόστασης Μινκόφσκι είναι:
 - Η απόσταση Μανχάταν (\mathcal{L}_1).
 - Για δύο διανύσματα χαρακτηριστικών/πρότυπα/στοιχεία με *μόνο δυαδικά γνωρίσματα*, η απόσταση Μανχάταν ταυτίζεται με την **απόσταση Hamming**: το πλήθος των διαφορετικών bit τους.
 - Η ευκλείδεια απόσταση (\mathcal{L}_2).
 - Η απόσταση μεγίστου (\mathcal{L}_{sup}).
 - Η απόσταση μεγίστου ισούται με τη μέγιστη απόλυτη διαφορά τιμής μεταξύ αντίστοιχων γνωρισμάτων/συνιστωσών των δύο διανυσμάτων/προτύπων/στοιχείων.

Εγγύτητα

Υπολογισμός ομοιοτήτων

- Για τα **μέτρα ομοιότητας** συνήθως ισχύει η *συμμετρία* και η *θετικότητα*, όχι όμως και η *τριγωνική ανισότητα*.
 - Όμως η τριγωνική ανισότητα μπορεί να αυξήσει την αποδοτικότητα ορισμένων αλγορίθμων εξόρυξης οι οποίοι την εκμεταλλεύονται.
 - Π.χ., κάποιες υλοποιήσεις των k -μέσων.
 - Έτσι, δεν είναι σπάνιο κατά την προεπεξεργασία να μετατρέπουμε υπολογισμένες τιμές ομοιότητας σε αποστάσεις.
- Παρακάτω αναλύονται ορισμένα διαδομένα μέτρα ομοιότητας για πολυδιάστατα στοιχεία/πρότυπα/διανύσματα χαρακτηριστικών.

Εγγύτητα

Simple Matching Coefficient

- Για δύο διανύσματα με **δυαδικά** γνωρίσματα, συνηθίζεται η χρήση του συντελεστή ομοιότητας **Simple Matching Coefficient (SMC)**.
- Είναι το πηλίκο του πλήθους των γνωρισμάτων/bit για τα οποία τα δύο διανύσματα έχουν ίδια τιμή, διά το ολικό πλήθος των γνωρισμάτων.
 - Το ποσοστό των ταυτώσεων (με κοινή τιμή) bit.
- Σε περίπτωση *ασύμμετρων δυαδικών* γνωρισμάτων χρησιμοποιείται μία παραλλαγή του SMC, ο **συντελεστής Jaccard**:
 - Στον τύπο υπολογισμού του δεν λαμβάνονται υπόψη τα γνωρίσματα/bit όπου και τα δύο πρότυπα έχουν τιμή 0.
 - Έτσι αποφεύγεται η ανάθεση υψηλής ομοιότητας σε δύο διανύσματα μόνο και μόνο επειδή έχουν και τα δύο πάρα πολλά μηδενικά bit.

Εγγύτητα

Συντελεστής συνημιτόνου

- Στην περίπτωση διανυσμάτων με **μη δυαδικά ασύμμετρα** γνωρίσματα, εφαρμόζεται συνήθως ο **συντελεστής συνημιτόνου**:
$$c = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\|_2 \|\mathbf{y}\|_2}$$
 - Επίσης λαμβάνει υπ' όψη την ασυμμετρία των τιμών, αλλά τα γνωρίσματα δεν χρειάζεται να είναι δυαδικά.
 - Δεν συνυπολογίζει τα μηδενικά γνωρίσματα.
- Εκφράζει το συνημίτονο της γωνίας που σχηματίζουν τα δύο διανύσματα στο διανυσματικό τους χώρο.
 - Λαμβάνει τιμές στο $[-1, 1]$. Δύο διανύσματα έχουν τιμή c :
 - -1 αν είναι αντιπαράλληλα, 0 αν είναι ορθογώνια, 1 αν είναι παράλληλα.
 - Τα μέτρα των διανυσμάτων/προτύπων αγνοούνται.
 - Ο συντελεστής συνημιτόνου δύο \mathcal{L}_2 -κανονικοποιημένων διανυσμάτων (μοναδιαίου ευκλείδειου μέτρου) ισούται με το εσωτερικό τους γινόμενο.

Εγγύτητα

Συσχέτιση

• Έναλλακτικό μέτρο ομοιότητας μεταξύ δύο διανυσμάτων/προτύπων \mathbf{x} και \mathbf{y} είναι η **συσχέτιση**:

- Το πηλίκο της συνδιακύμανσής τους (v_{xy}) διά το παρακάτω γινόμενο: τυπική απόκλιση του \mathbf{x} (v_x) επί την τυπική απόκλιση του \mathbf{y} (v_y).
- Ο υπολογισμός της τυπικής απόκλισης κάθε διανύσματος γίνεται επί όλων των τιμών των γνωρισμάτων/συνιστωσών του.

$$v_{xy} = \frac{1}{m-1} \sum_{k=1}^m (x_k - \bar{\mathbf{x}})(y_k - \bar{\mathbf{y}}).$$

$$v_x = \left[\frac{1}{m-1} \sum_{k=1}^m (x_k - \bar{\mathbf{x}})^2 \right]^{1/2}.$$

$$v_y = \left[\frac{1}{m-1} \sum_{k=1}^m (y_k - \bar{\mathbf{y}})^2 \right]^{1/2}.$$

- Η συνδιακύμανση μετρά την τάση ταυτόχρονης αυξομείωσης των τιμών των \mathbf{x} και \mathbf{y} , καθώς διατρέχουμε τα m γνωρίσματά τους.
- Η συσχέτιση είναι απλώς κανονικοποιημένη συνδιακύμανση.

Εγγύτητα

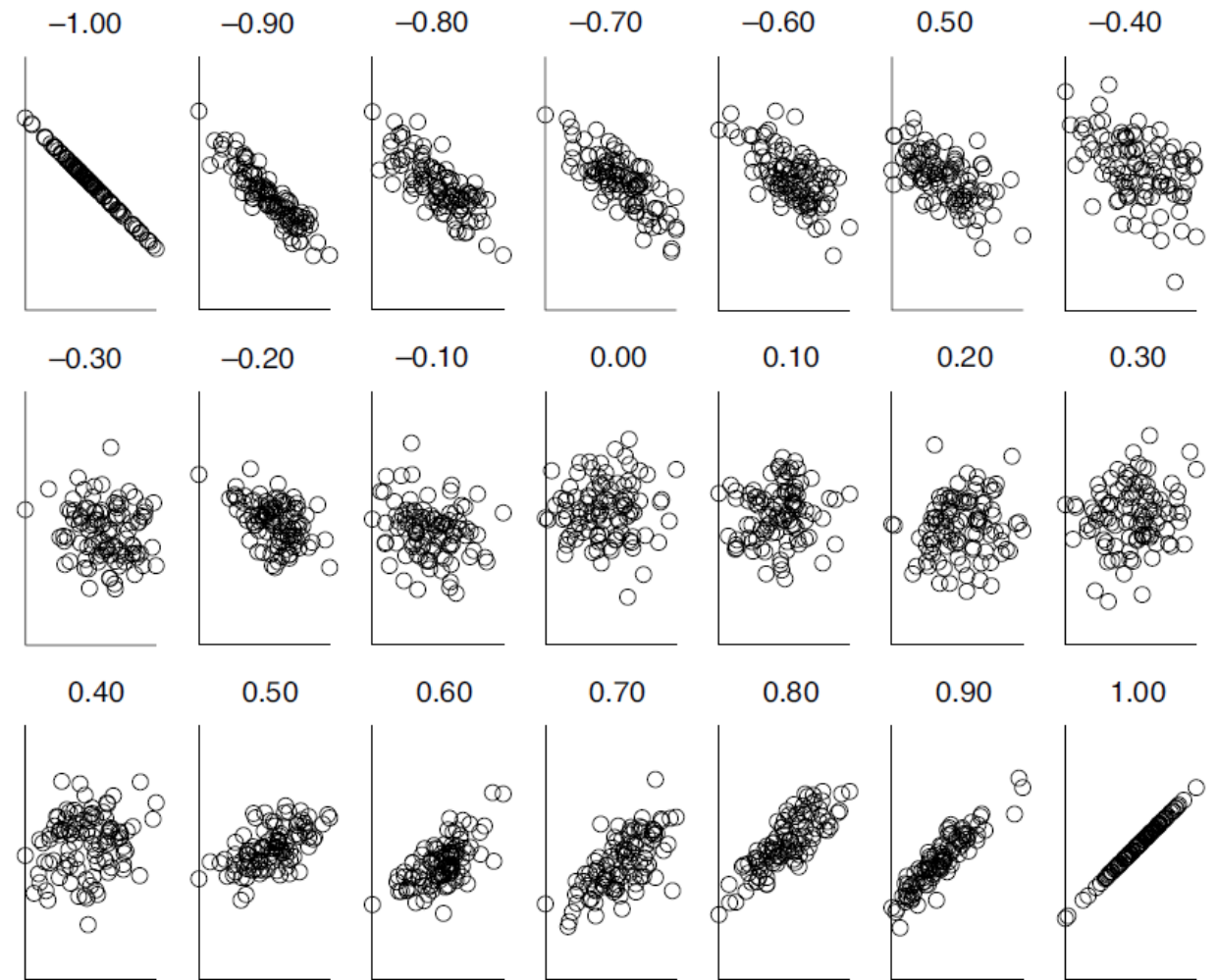
Συσχέτιση

- Η συσχέτιση είναι μία πραγματική τιμή στο $[-1, 1]$ η οποία συλλαμβάνει τη γραμμική σχέση μεταξύ των δύο στοιχείων:
 - Είναι $-1/1$ αν το ένα διάνυσμα είναι βαθμωτό πολλαπλάσιο του δεύτερου με αρνητικό/θετικό συντελεστή κλιμάκωσης, αντιστοίχως.
 - Αν η συσχέτισή τους είναι 0 , δεν υπάρχει καμία ανιχνεύσιμη γραμμική σχέση μεταξύ των τιμών των γνωρισμάτων των δύο διανυσμάτων.
- Μπορούμε να οπτικοποιήσουμε τα m γνωρίσματα των δύο διανυσμάτων ως m σημεία σε μία διδιάστατη γραφική παράσταση.
 - Ο οριζόντιος/κατακόρυφος άξονας εκφράζει το πρώτο/δεύτερο διάνυσμα, αντιστοίχως.
 - Η ισχυρή αρνητική ή θετική συσχέτιση είναι εμφανής ως μία ισχυρή τάση διασποράς των m σημείων κατά μήκος μίας συγκεκριμένης διεύθυνσης.

Εγγύτητα

Συσχέτιση

Διαγράμματα
γνωρισμάτων δύο
διανυσμάτων για
ποικίλες
διαφορετικές
τιμές συσχέτισης
στο $[-1,1]$.



Εγγύτητα

Συσχέτιση

- Συσχέτιση όμως μπορεί να εμφανίζεται μεταξύ όχι δύο συγκεκριμένων προτύπων, αλλά μεταξύ κάποιων γνωρισμάτων, αν αυτά υπολογιστούν επί όλων των προτύπων του συνόλου δεδομένων.
 - Συντελεστής Pearson ρ .
- Πώς υπολογίζονται τα ρ ;
 - Θεωρούμε καθένα από τα N m -διάστατα πρότυπα του συνόλου δεδομένων μας ως ξεχωριστή παρατήρηση ενός τυχαίου διανύσματος των m τυχαίων μεταβλητών.
 - Υπολογίζουμε ξεχωριστά τη συσχέτιση μεταξύ κάθε δυνατού ζεύγους γνωρισμάτων.
 - Τα αθροίσματα στους τύπους τώρα διατρέχουν τα διαφορετικά πρότυπα, όχι τα διαφορετικά γνωρίσματα.

Εγγύτητα

Συσχέτιση

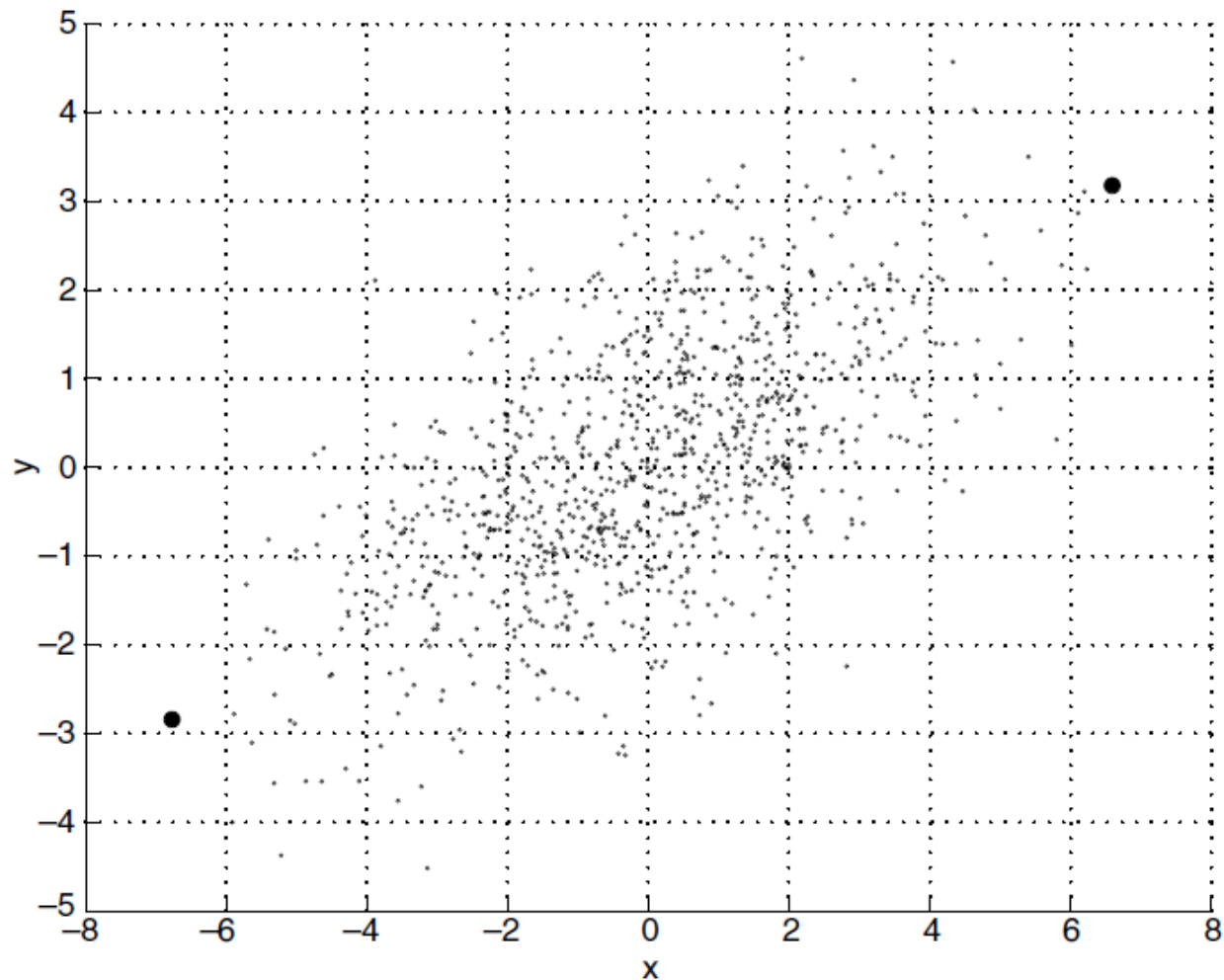
- Αν υπάρχουν υψηλές συσχετίσεις μεταξύ ορισμένων γνωρισμάτων, τα N m -διάστατα πρότυπα παρουσιάζουν μία ισχυρή τάση **διασποράς** στον διανυσματικό τους χώρο κατά μήκος μίας συγκεκριμένης διεύθυνσης.
- Σε περίπτωση 2-διάστατων ή 3-διάστατων προτύπων, αυτό είναι και άμεσα οπτικά εμφανές.
 - Οπτικοποίηση των N διανυσμάτων με ένα διάγραμμα όπου καθένας από τους m άξονες εκφράζει ένα γνώρισμα.
- Σε περίπτωση δεδομένων υψηλότερης διάστασης, η PCA μας δίνει τους άξονες της υψηλής διασποράς.
 - Υπολογίζει το ελλειψοειδές το οποίο περικλείει βέλτιστα τα πρότυπα στον m -διάστατο χώρο.
 - Επιστρέφει τους άξονες του ελλειψοειδούς.

Εγγύτητα

Συσχέτιση

Απεικόνιση των προτύπων ενός συνόλου διδιάστατων δεδομένων. Είναι εμφανής η υψηλή συσχέτιση μεταξύ των δύο γνωρισμάτων.

Εγγύτητα



Αντιμετώπιση συσχέτισης

- Η ύπαρξη αξόνων ασυνήθιστα υψηλής διασποράς, λόγω συσχετίσεων μεταξύ γνωρισμάτων, μπορεί να οδηγήσει σε προβλήματα κατά την ανάλυση των δεδομένων.
 - Υπέρμετρη έμφαση ενός υπολογισμού απόστασης στην απόσταση κατά μήκος της διεύθυνσης μέγιστης διασποράς.
 - Αποτέλεσμα: **απώλεια πληροφορίας** όσον αφορά τη *γεωμετρία* των δεδομένων, δηλαδή τη διασπορά τους στον διανυσματικό τους χώρο.
 - Π.χ.: δύο πρότυπα x και y τα οποία απέχουν εξίσου από το κέντρο του συνόλου δεδομένων, όπου τοποθετείται το μέσο πρότυπο μ , ενδέχεται να έχουν πολύ διαφορετικές ιδιότητες:
 - Το x έχει εγγύς του γειτονικά στοιχεία, διότι είναι τοποθετημένο επάνω στον άξονα μέγιστης διασποράς, ενώ το y όχι, διότι είναι τοποθετημένο επάνω σε έναν ορθογώνιο άξονα.
 - Αποτέλεσμα: ο απλός υπολογισμός απόστασης **δεν θα συλλάβει αυτή την ποιοτική διαφορά** μεταξύ x και y .

Εγγύτητα

Αντιμετώπιση συσχέτισης

- Λύση: αν οι συνδιακυμάνσεις οργανωθούν σε έναν συμμετρικό πίνακα $\Sigma \in \mathbb{R}^{m \times m}$ μπορούμε να υπολογίσουμε τον αντίστροφό του Σ^{-1} .
- **Απόσταση Mahalanobis:**
 - $M(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^T \Sigma^{-1} (\mathbf{x} - \mathbf{y})}$.
 - Για ταυτοτικό Σ , ισούται με την ευκλείδεια απόσταση!
- Ο υπολογισμός της απόστασης Mahalanobis μεταξύ δύο προτύπων *ισοδυναμεί* με τα εξής βήματα:
 - Υπολογισμός των αξόνων μέγιστης διασποράς και των μηκών τους, μέσω PCA.
 - Ιδιοδιανύσματα και ιδιοτιμές του Σ .
 - Επανεκφραση των δεδομένων σε μία νέα διανυσματική βάση, μέσω διαδοχικής **περιστροφής** και **κλιμάκωσης**.
 - Υπολογισμός ευκλείδεια απόστασης μεταξύ των μετασχηματισμένων προτύπων.

Εγγύτητα

Αντιμετώπιση συσχέτισης

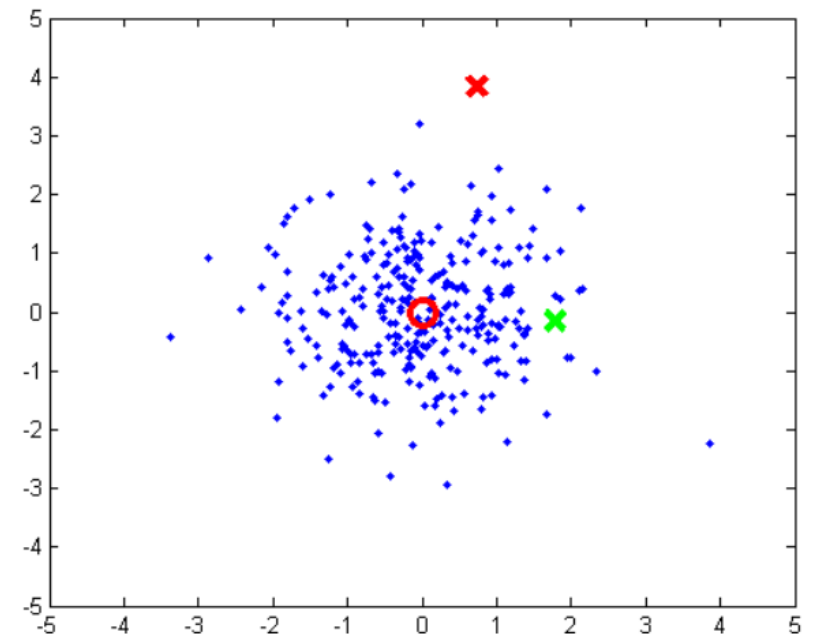
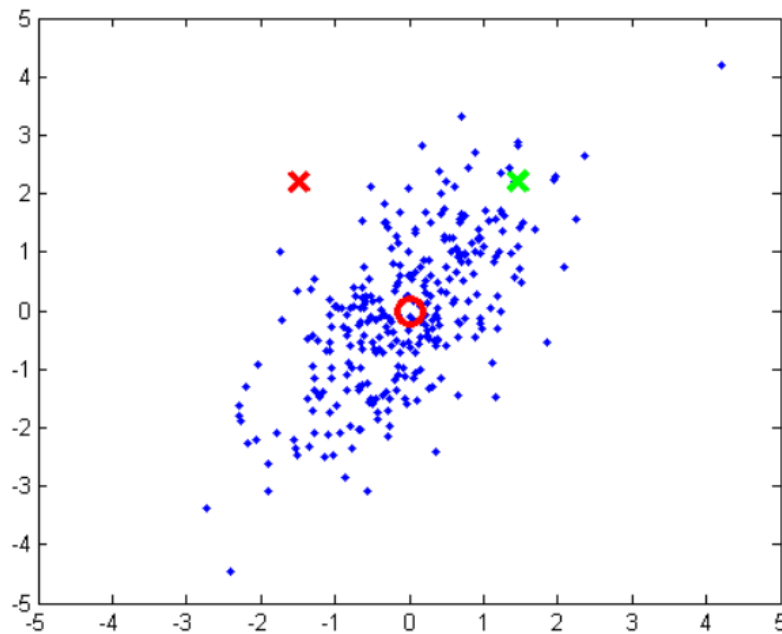
- Η ουσία είναι στην επανέκφραση των δεδομένων μέσω διαδοχικής **περιστροφής** και **κλιμάκωσης**:
 - Οι άξονες του νέου συστήματος συντεταγμένων είναι οι κύριες συνιστώσες των αρχικών δεδομένων (από PCA).
 - Η κλίμακα καθενός από τους νέους άξονες είναι πλέον τέτοια ώστε τα δεδομένα να έχουν **μοναδιαία διακύμανση** κατά μήκος του.
 - Ο πίνακας συνδιακύμανσης των μετασχηματισμένων δεδομένων είναι πλέον ο **ταυτοτικός πίνακας**.
 - Άρα, τα μετασχηματισμένα δεδομένα έχουν πλέον **σφαιρική** διασπορά στον m -διάστατο χώρο!
 - Ο πίνακας συνδιακύμανσης περιγράφει πλήρως τη διασπορά (άρα το ολικό γεωμετρικό σχήμα στον m -διάστατο χώρο) μίας γκαουσιανής προσέγγισης του συνόλου δεδομένων.

Εγγύτητα

Αντιμετώπιση συσχέτισης

- Οι ευκλείδειες αποστάσεις του κόκκινου και του πράσινου προτύπου από το μ (κύκλος στο κέντρο) είναι ίσες.
- Όμως η απόσταση Mahalanobis δίνει μεγαλύτερη απόσταση από το μ για το κόκκινο πρότυπο, απ' ότι για το πράσινο!
- Η απόσταση Mahalanobis ισοδυναμεί με υπολογισμό της ευκλείδειας στα μετασχηματισμένα (σφαιρικά) δεδομένα.

Εγγύτητα



Αντιμετώπιση συσχέτισης

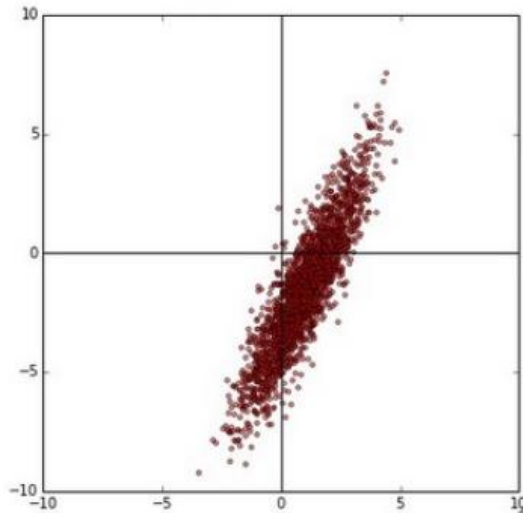
- Ο μετασχηματισμός των δεδομένων έτσι ώστε να έχουν σφαιρική διασπορά καλείται **λεύκανση** (whitening):
 - Έχει λάβει το όνομά της από τον *λευκό θόρυβο*.
 - Υπάρχουν και άλλοι τρόποι λεύκανσης, πέρα από τη χρήση PCA.
- Γιατί τα μετασχηματισμένα από τον Σ^{-1} δεδομένα έχουν πλέον σφαιρική διασπορά;
 - Λύση PCA: $\Sigma = \mathbf{U}\mathbf{S}\mathbf{U}^T$, όπου:
 - Ο \mathbf{U} είναι **ορθογώνιος** πίνακας και περιέχει τις κύριες συνιστώσες.
 - Ο \mathbf{S} είναι διαγώνιος πίνακας και περιέχει τα αντίστοιχα μήκη τους.
 - Ο αντίστροφος \mathbf{A} ενός διαγώνιου πίνακα \mathbf{B} είναι διαγώνιος πίνακας όπου $A_{ii} = \frac{1}{B_{ii}}$.

Εγγύτητα

Αντιμετώπιση συσχέτισης

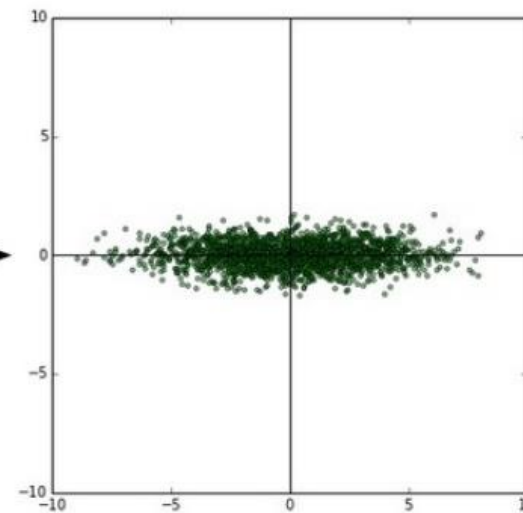
- Κατά τη λεύκανση των δεδομένων, σύνηθες στάδιο προεπεξεργασίας είναι το **κεντράρισμα** των δεδομένων.
 - Από κάθε πρότυπο αφαιρούμε το μ .
- Η λεύκανση είναι πιο ισχυρή από την απλή αποσυσχέτιση.

Αρχικά δεδομένα



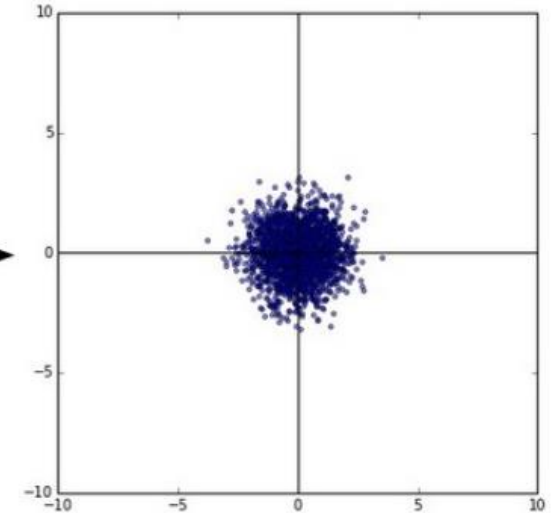
Πυκνός Σ

Αποσυσχετισμένα δεδομένα



Διαγώνιος Σ

Σφαιρικά δεδομένα



Ταυτοτικός Σ

Εγγύτητα

Thank you for your attention!

Q & A

Contact: imademlis@aueb.gr