



**ΔΙΕΘΝΕΣ ΠΑΝΕΠΙΣΤΗΜΙΟ ΕΛΛΑΔΟΣ**

**ΣΧΟΛΗ ΜΗΧΑΝΙΚΩΝ**

**ΤΜΗΜΑ ΜΗΧΑΝΙΚΩΝ ΠΛΗΡΟΦΟΡΙΚΗΣ  
ΥΠΟΛΟΓΙΣΤΩΝ & ΤΗΛΕΠΙΚΟΙΝΩΝΙΩΝ**

**ΥΠΟΛΟΓΙΣΤΙΚΗ ΝΟΗΜΟΣΥΝΗ**

**Σημειώσεις**

Χαράλαμπος Π. Στρουθόπουλος  
Καθηγητής

ΣΕΡΡΕΣ ΙΑΝΟΥΑΡΙΟΣ 2024

# ΠΕΡΙΕΧΟΜΕΝΑ

ΚΕΦΑΛΑΙΟ 1.....	3
ΕΙΣΑΓΩΓΗ.....	3
ΚΕΦΑΛΑΙΟ 2.....	6
ΠΕΡΙΓΡΑΦΗ ΤΩΝ ΠΡΟΤΥΠΩΝ ΚΑΙ ΟΡΙΣΜΟΙ.....	6
2.2 Εσωτερικό γινόμενο.....	9
2.3 Αποστάσεις.....	10
ΚΕΦΑΛΑΙΟ 3.....	12
ΜΕΘΟΔΟΙ ΤΑΞΙΝΟΜΗΣΗΣ – ΕΚΠΑΙΔΕΥΣΗ ΜΕ ΕΠΟΠΤΗ.....	12
3.1. Ταξινόμηση με βάση την απόσταση από τους K-γείτονες .....	12
3.2 Αναγνώριση με βάση τα κέντρα των τάξεων .....	13
3.3 Ταξινόμηση σε δύο κλάσεις με γραμμική διακριτική συνάρτηση, ο νευρώνας Perceptron .....	13
3.4.1. Η περίπτωση πολλών κλάσεων. ....	19
3.4.2 Η περίπτωση της XOR, ταξινομητές πολλών επιπέδων.....	21
3.5 Πολυεπίπεδοι ταξινομητές – Διόρθωση σφάλματος με οπισθοδιάδοση (Back Error Propagation) .....	24
Εκπαίδευση με οπισθοδρόμηση του σφάλματος (back-error propagation).....	27
3.6 Δένδρα απόφασης.....	32
3.7 Συνελκτικά Νευρωνικά Δίκτυα .....	34
Επίπεδο Συνέλιξης (Convolution) .....	36
Επίπεδο δειγματοληψίας (Pooling Layer).....	39
Επίπεδο Ταξινόμησης .....	40
Υπολογισμός του πλήθους βαρών και εκπαίδευση στα ΣΝΔ.....	41
Εκπαίδευση των Συνελκτικών Νευρωνικών Δικτύων.....	42
Εκπαίδευση στο επίπεδο πλήρους διασύνδεσης.....	42
Αλγόριθμος οπισθοδρόμησης σφάλματος στο επίπεδο δειγματοληψίας. ....	44
Αλγόριθμος οπισθοδρόμησης σφάλματος στο επίπεδο συνέλιξης.....	45
ΚΕΦΑΛΑΙΟ 4.....	48
ΕΚΠΑΙΔΕΥΣΗ ΧΩΡΙΣ ΕΠΟΠΤΗ .....	48
4.1 Ο Αλγόριθμος ISODATA ή K-Μέσων (k-means ή c-means) .....	48
4.2 Απεικόνιση αλυσίδας κοντινών γειτόνων.....	51
4.3 Χάρτης απεικόνισης χαρακτηριστικών με αυτό-οργάνωση .....	53
ΚΕΦΑΛΑΙΟ 5.....	58
ΑΝΑΛΥΣΗ ΧΑΡΑΚΤΗΡΙΣΤΙΚΩΝ.....	58

5.1 Ανάλυση χαρακτηριστικών στην εκπαίδευση με επόπτη .....	58
5.2 Ανάλυση χαρακτηριστικών στην εκπαίδευση χωρίς επόπτη .....	60
3.7 Ταξινομητές με βάση τον κανόνα πιθανοτήτων του Bayes .....	65
ΠΑΡΑΡΤΗΜΑ Α: ΒΑΣΙΚΕΣ ΠΡΑΞΕΙΣ ΓΡΑΜΜΙΚΗΣ ΑΛΓΕΒΡΑΣ.....	71
ΠΑΡΑΡΤΗΜΑ Γ: Στατιστικά χαρακτηριστικά του χώρου των προτύπων. ....	71
ΠΑΡΑΡΤΗΜΑ Γ:ΤΑ ΙΔΙΟΔΙΑΝΥΣΜΑΤΑ ΤΟΥ ΠΙΝΑΚΑ ΣΥΝΔΙΑΣΠΟΡΑΣ .....	72

# ΚΕΦΑΛΑΙΟ 1

## ΕΙΣΑΓΩΓΗ

Η Υπολογιστική Νοημοσύνη είναι ένα υποσύνολο της τεχνητής νοημοσύνης (AI) που δομεί συστήματα ικανά να μαθαίνουν και να βελτιώνονται από την εμπειρία δηλαδή από τα δεδομένα εισόδου εξόδου, χωρίς να είναι ρητά δηλωμένος ο μηχανισμός που διέπει την σχέση τους. Για παράδειγμα μία λογική πύλη (π.χ. της διάζευξης, XOR) περιγράφεται ρητά με την άλγεβρα του *Boole*, εν τούτοις είναι δυνατόν μόνο από τα δεδομένα του πίνακα αληθείας της να δημιουργήσουμε ένα σύστημα Μηχανικής Μάθησης που θα προσομοιώνει την λειτουργία της. Αναφέρεται το σύστημα αυτό και ως εμπειρικό μοντέλο. Η αδυναμία ρητής περιγραφής ενός συστήματος και η ανάγκη επιστράτευσης της μηχανικής μάθησης μπορεί να προέρχεται είτε από την πολυπλοκότητά του είτε από την αδυναμία γλωσσικής περιγραφής. Για παράδειγμα δεν μπορούμε να διδάξουμε σε κάποιον χρησιμοποιώντας την φυσική γλώσσα πως να κάνει ποδήλατο, πρέπει να δοκιμάσει και να μάθει εμπειρικά.

Βασικές έννοιες και όροι που εμπλέκονται συχνά στην Υπολογιστική Νοημοσύνη είναι επιγραμματικά:

- Τα δεδομένα
- Οι αλγόριθμοι ταξινόμησης, παλινδρόμησης, ομαδοποίησης, ενισχυτικής μάθησης.
- Η Εκπαίδευση και ή δοκιμή
- Η Μάθηση, εποπτευόμενη, μη εποπτευόμενη και ενισχυτική
- Η Ανάλυση Χαρακτηριστικών

Εφαρμογές Μηχανικής Μάθησης:

Η Υπολογιστική Νοημοσύνη εφαρμόζεται σε διάφορους τομείς που αναφέρονται ενδεικτικά παρακάτω:

- Επεξεργασία Φυσικής Γλώσσας (NLP)
- Αναγνώριση εικόνας και ομιλίας
- Προβλεπτική Ανάλυση
- Συστήματα συστάσεων
- Υγεία και Ιατρική
- Χρηματοδότηση
- Αυτόνομα Οχήματα

Ένα σύνολο όμοιων προτύπων, όχι αναγκαστικά πανομοιότυπων, ορίζει μία σύλληψη (concept). Τα νοήμονα έμβια παρατηρούν και μαθαίνουν τα όντα (πρότυπα) του περιβάλλοντος κόσμου και δημιουργούν γι' αυτά συλλήψεις. Για παράδειγμα το σύνολο των αλόγων που γνωρίσατε σας έμαθε την γενική σύλληψη (ιδέα) άλογο (κάντε συσχετισμό με τον κόσμο των ιδεών του Πλάτωνα). Στον χώρο των γραμμάτων της Ελληνικής γλώσσας υπάρχουν 24 συλλήψεις που είναι τα 24 γράμματα της Ελληνικής αλφαβήτου. Κάθε γράμμα της Ελληνικής γλώσσας που γράφτηκε αποτελεί μία πραγματοποίηση ή συμβάν (event) μιας από τις 24 συλλήψεις. Για παράδειγμα οι μορφές:

β β β β **B** β β β **B**

δ Δ δ Δ δ Δ δ Δ δ

αποτελούν συμβάντα (πρότυπα) των συλλήψεων των γραμμάτων ΒΗΤΑ και ΔΕΛΤΑ. Τα συγκεκριμένα πρότυπα μπορεί να διαφέρουν στο μέγεθος, την γραμματοσειρά, να είναι κεφαλαία ή μικρά, ή να έχουν μικροδιαφορές που δημιουργήθηκαν κατά τη στιγμή της εκτύπωσης (θόρυβος). Αντικείμενα (πρότυπα) που ανήκουν στην ίδια σύλληψη αποτελούν μία κλάση ή τάξη (class). Είναι προφανές πως οι συλλήψεις και τα οικεία πρότυπά τους πρέπει να διαθέτουν χαρακτηριστικά που να τις διακρίνουν από άλλες συλλήψεις και πρότυπα. Η εύρεση και μέτρηση τέτοιων χαρακτηριστικών (features) των προτύπων είναι πρωταρχικής σημασίας για την περιγραφή και την αναγνώρισή τους.

Η ταξινόμηση των αντικειμένων σε κλάσεις λέγεται *αναγνώριση προτύπων (pattern recognition)* και μπορεί να αντικαταστήσει πληθώρα επίπονων ανθρώπινων εργασιών, να βελτιώσει το εργασιακό περιβάλλον και να αυξήσει την παραγωγικότητα. Για παράδειγμα η μετατροπή τυπωμένων εγγράφων σε ηλεκτρονική μορφή με κωδικοποίηση κατά ASCII, μπορεί να γίνει από τον άνθρωπο δακτυλογράφο που διαβάζει το κείμενο, αναγνωρίζει τους χαρακτήρες και τους

πληκτρολογεί, μπορεί όμως ακόμη να γίνει με την σάρωση και ψηφιοποίηση (scanning) του εγγράφου από κατάλληλο πρόγραμμα οπτικής αναγνώρισης χαρακτήρων (OCR: *Optical Character Recognition*) που θα αναγνωρίζει την ψηφιακή μορφή κάθε χαρακτήρα και θα τον κωδικοποιεί με το αντίστοιχο κώδικα ASCII.

Ένα σύστημα (αλγόριθμος) που αναγνωρίζει πρότυπα λέγεται *ταξινομητής (classifier)*. Σε ένα ταξινομητή επιτελούνται δύο κύριες εργασίες που είναι η *εκπαίδευση (training)* και η *ταξινόμηση*. Κατά την εκπαίδευση διαχωρίζονται κατάλληλα οι κλάσεις ή προσδιορίζονται οι συγκεντρώσεις των προτύπων και ρυθμίζονται οι παράμετροι του συστήματος ταξινόμησης με βάση τις οποίες θα είναι αυτό ικανό να ταξινομήσει τα πρότυπα σωστά στις κλάσεις. Με την ταξινόμηση αποδίδεται ένα πρότυπο σε μία κλάση ή σε μία συγκέντρωση. Οι μέθοδοι εκπαίδευσης χωρίζονται σε δύο βασικές κατηγορίες:

A) Εκπαίδευση με επόπτη (*supervised learning*) κατά την οποία είναι γνωστή η κλάση στην οποία ανήκει κάθε πρότυπο του συνόλου εκπαίδευσης και επιδιώκεται η ορθή απόδοση των προτύπων που δεν ανήκουν στο σύνολο εκπαίδευσης, σε μία από τις προκαθορισμένες κλάσεις. Σε κάθε πρότυπο του συνόλου εκπαίδευσης αντιστοιχίζουμε το όνομα της κλάσης του με μία ετικέτα (*label*).

B) Εκπαίδευση χωρίς επόπτη (*unsupervised learning*) κατά την οποία δεν είναι γνωστή η κλάση στην οποία ανήκει κάθε πρότυπο του συνόλου εκπαίδευσης και αναζητούνται βασικά οι συγκεντρώσεις των προτύπων (*clustering*).

Τα πρότυπα που χρησιμοποιούνται για την εκπαίδευση του ταξινομητή αποτελούν το *σύνολο εκπαίδευσης (training set)*. Συχνά χρησιμοποιούμε και ένα σύνολο προτύπων για την επικύρωση των αποτελεσμάτων το οποίο ονομάζουμε *σύνολο επικύρωσης (validation set)*.

Προτιμότερο είναι η εκπαίδευση ενός ταξινομητή να μην εκτελείται μόνο μία αρχική φορά, αλλά να επαναλαμβάνεται ώστε η μηχανή να προσαρμόζεται στις αλλαγές των προτύπων που ταξινομεί. Οι ταξινομητές είναι ουσιαστικά το αντικείμενο με το οποίο θα ασχοληθούμε περαιτέρω και θα αναλυθούν στα επόμενα κεφάλαια.

## ΚΕΦΑΛΑΙΟ 2

### ΠΕΡΙΓΡΑΦΗ ΤΩΝ ΠΡΟΤΥΠΩΝ ΚΑΙ ΟΡΙΣΜΟΙ

Για κάθε πρότυπο μετρούμε ένα πεπερασμένο πλήθος χαρακτηριστικών του (π.χ. βάρος, χρώμα, ύψος, πλάτος, μήκος). Για ένα ποιοτικό χαρακτηριστικό που δεν μετριέται άμεσα με αριθμητικές τιμές όπως για παράδειγμα το χρώμα, χρησιμοποιούμε κάποια κωδικοποίηση που αντιστοιχεί αριθμούς στις διάφορες μορφές του (π.χ. 0 για το μαύρο, 1 για το μπλε, 2 για το κόκκινο κ.λ.π.) Αν  $\Omega$  είναι ένα σύνολο που αποτελείται από  $K$  το πλήθος πρότυπα  $\Pi_k$ ,  $k=1,2,3,\dots,K$ , δηλαδή  $\Omega=\{\Pi_1,\Pi_2,\dots,\Pi_k,\dots,\Pi_K\}$  και μετρούμε  $N$  το πλήθος κοινά χαρακτηριστικά για κάθε πρότυπο, τότε το  $\Pi_k$  πρότυπο μπορεί να περιγραφεί από ένα διάνυσμα (πίνακα στήλης)  $\mathbf{x}_k$  σε ένα χώρο  $N$  διαστάσεων σύμφωνα με την σχέση

$$\mathbf{x}_k = \begin{bmatrix} x_{1k} \\ x_{2k} \\ \vdots \\ x_{vk} \\ \vdots \\ x_{Nk} \end{bmatrix} = |x_{1k} \quad x_{2k} \quad \dots \quad x_{vk} \quad \dots \quad x_{Nk}|^T \quad (2.1.1)$$

Κάθε πρότυπο θεωρείται ξεχωριστό από κάθε άλλο ακόμη και όταν όλες οι τιμές των αντίστοιχων χαρακτηριστικών τους είναι ίσες. Δηλαδή το πρότυπο  $\Pi_\mu$  είναι άλλο στοιχείο του συνόλου  $\Omega$ , από το πρότυπο  $\Pi_\lambda$  αν και  $\mathbf{x}_\mu=\mathbf{x}_\lambda$  ( $\mu, \lambda \in \{1,2,\dots,K\}$ , και  $\Pi_\mu, \Pi_\lambda \in \Omega$ ). Ως εκ τούτου το  $\Omega$  είναι σύνολο με την μαθηματική έννοια. Όλους τους πίνακες θα τους συμβολίζουμε με παχιά γράμματα ή με υπογράμμιση.

Αν  $K$  το πλήθος των προτύπων,  $N$  το πλήθος των χαρακτηριστικών που μετράμε σε κάθε πρότυπο και  $T$  το πλήθος των τάξεων, μπορούμε να χρησιμοποιήσουμε ένα δείκτη  $k=1,2,\dots,K$  για κάθε πρότυπο, ένα δείκτη  $v=1,2,\dots,N$  για κάθε χαρακτηριστικό και ένα δείκτη  $t=1,2,\dots,T$  για κάθε κλάση. Ο διανυσματικός

χώρος των προτύπων έχει διάσταση  $N$ , ( $R^N$ ). Για την τιμή του  $n$ -οστού χαρακτηριστικού του  $k$  προτύπου που ανήκει στην κλάση  $\tau$  γράφουμε

$$x_{nk}^{\tau}.$$

Συχνά κάποιοι δείκτες δεν χρησιμοποιούνται όταν έχουν προφανή τιμή ή όταν η αναφορά τους δεν είναι απαραίτητη. Η μαθηματική περιγραφή των πραγμάτων στην αναγνώριση προτύπων γίνεται με την χρήση όρων της γραμμικής άλγεβρας. Το διάνυσμα  $\vec{x}$  που περιέχει  $N$  χαρακτηριστικά ενός προτύπου θα θεωρείται πίνακας με  $N$  γραμμές και μία στήλη και θα συμβολίζεται  $\mathbf{x}_{N \times 1}$  ή χάριν συντομίας  $\mathbf{x}$  ή  $\underline{x}$ . Συχνά θα χρησιμοποιούμε τον ανάστροφο πίνακα  $\mathbf{A}^T$  ενός πίνακα  $\mathbf{A}$  (Παράρτημα Α). Για εξοικονόμηση χώρου στην σελίδα, τον πίνακα στήλης  $\mathbf{x}$  ενός προτύπου θα τον γράφουμε ανεστραμμένο  $\mathbf{x}^T = [x_1, x_2, \dots, x_n, \dots, x_N]$ . Με βάση τον ορισμό του πολλαπλασιασμού δύο πινάκων (Παράρτημα Α) η Ευκλείδεια απόσταση  $d_E(\mathbf{x}, \mathbf{y})$  των πινάκων στήλης  $\mathbf{x}$ ,  $\mathbf{y}$  δύο προτύπων  $P_x, P_y$  δίνεται από την σχέση

$$d_E(\mathbf{x}, \mathbf{y}) = [(\mathbf{x}-\mathbf{y})^T(\mathbf{x}-\mathbf{y})]^{1/2} \quad (2.1.2)$$

Αν  $N \leq 3$  τότε μπορούμε να παραστήσουμε τα διανύσματα των προτύπων στον φυσικό χώρο. Ακολούθως θα δούμε ένα παράδειγμα προτύπων που είναι άτομα για τα οποία μετράμε δύο χαρακτηριστικά που είναι το βάρος και το ύψος. Οι τιμές των χαρακτηριστικών για όλα τα πρότυπα που μετρήθηκαν δείχνονται στο ακόλουθο πίνακα 2.1.

<b>Π</b> (πρότυπο)	<b>Υ</b> (ύψος σε εκ.)	<b>Β</b> (βάρος σε κιλά)
Π1	185	85
Π2	185	82
Π3	182	115
Π4	180	80
Π5	178	78
Π6	175	112
Π7	171	114
Π8	160	58
Π9	160	55
Π10	155	62
Π11	150	85
Π12	150	82
Π13	145	84
Π14	143	80

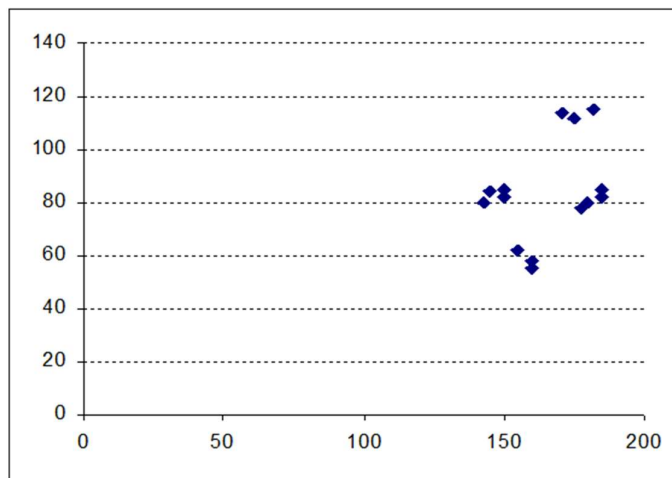
Πίνακας 2.1



Τα αντίστοιχα διανύσματα των προτύπων γράφονται ως πίνακες στήλης όπως παρακάτω:

$$\begin{aligned} \mathbf{x}_1 &= \begin{bmatrix} 185 \\ 85 \end{bmatrix}, \quad \mathbf{x}_2 = \begin{bmatrix} 185 \\ 82 \end{bmatrix}, \quad \mathbf{x}_3 = \begin{bmatrix} 182 \\ 115 \end{bmatrix}, \quad \mathbf{x}_4 = \begin{bmatrix} 180 \\ 80 \end{bmatrix}, \\ \mathbf{x}_5 &= \begin{bmatrix} 178 \\ 78 \end{bmatrix}, \quad \mathbf{x}_6 = \begin{bmatrix} 175 \\ 112 \end{bmatrix}, \quad \mathbf{x}_7 = \begin{bmatrix} 171 \\ 114 \end{bmatrix}, \quad \mathbf{x}_8 = \begin{bmatrix} 160 \\ 58 \end{bmatrix}, \\ \mathbf{x}_9 &= \begin{bmatrix} 160 \\ 55 \end{bmatrix}, \dots, \quad \mathbf{x}_{14} = \begin{bmatrix} 143 \\ 80 \end{bmatrix} \end{aligned} \quad (2.1.3)$$

Στο Σχ.2.1 που ακολουθεί φαίνονται τα σημεία τέλους των διανυσμάτων που αντιστοιχούν σε κάθε πρότυπο.



Σχήμα.2.1.

Έστω ότι επιθυμούμε να κατατάξουμε τα πρότυπα μας σε τάξεις ως εξής: Βαρείς αν  $B > 80$ , ελαφρούς αν  $B \leq 80$ , ψηλούς αν  $Y > 170$  και κοντούς αν  $Y \leq 170$ . Παρατηρούμε ότι τα άκρα των διανυσμάτων των προτύπων συγκεντρώνονται σε ομάδες. Το γεγονός αυτό οφείλεται στο ότι τα όμοια πρότυπα γειτνιάζουν επειδή τα αντίστοιχα χαρακτηριστικά τους έχουν κοντινότερες τιμές (ομόλογες συντεταγμένες των ανυσμάτων τους) σε σχέση με τις τιμές των χαρακτηριστικών των άλλων προτύπων. Για παράδειγμα η ευκλείδεια απόσταση  $D_{12}$  των προτύπων  $\Pi_1, \Pi_2$  είναι

$$D_{12} = \sqrt{(x_{11} - x_{12})^2 + (x_{21} - x_{22})^2} = \sqrt{(185 - 185)^2 + (85 - 82)^2} = 3 \quad (2.1.4)$$

και η ευκλείδεια απόσταση των προτύπων  $\Pi_1, \Pi_5$  είναι

$$D_{15} = \sqrt{(x_{11} - x_{15})^2 + (x_{21} - x_{25})^2} = \sqrt{(185 - 178)^2 + (85 - 78)^2} = 9\sqrt{2} \quad (2.1.5)$$

Άρα η απόσταση των ανυσμάτων των προτύπων είναι ένα μαθηματικό κριτήριο της ομοιότητας; δύο προτύπων. Οι ομάδες που σχηματίζουν τα πρότυπα λέγονται

συγκεντρώσεις (clusters). Οι συγκεντρώσεις δημιουργούνται από πρότυπα με μορφολογική ομοιότητα. Στο παράδειγμα μας οι τέσσερις συγκεντρώσεις δημιουργήθηκαν από άτομα α) με  $B > 100$  κ. β) με  $B < 70$  κ. γ) με  $Y < 1.80\mu$  και  $70 < B < 100$  κ. και δ) με  $Y > 1.80\mu$  και  $70 < B < 100$ κ.

Ο εντοπισμός των συγκεντρώσεων (clusters) προτύπων που δημιουργούνται στον χώρο των προτύπων είναι σημαντικός για την αναγνώριση προτύπων και αναφέρεται ως πρόβλημα εύρεσης των συγκεντρώσεων (clustering). Η σωστή επιλογή των χαρακτηριστικών που μετράμε για τα πρότυπα (feature selection) είναι καθοριστικής σημασίας για την εύρεση των ομάδων των προτύπων. Συχνά ενδιαφερόμαστε για πρότυπα που έχουν κάποια ιδιότητα και βρίσκονται σε περισσότερες από μία ομάδες. Για παράδειγμα τα πρότυπα των ατόμων με ύψος μεγαλύτερο από  $1,70 \mu$  ανήκουν σε δύο συγκεντρώσεις. Ανάλογα συμβαίνει και με τα πρότυπα των πεζών και κεφαλαίων γραμμάτων όπως των 'Α', 'α' ή 'Ω', 'ω' κ.α. Η εύρεση της ομάδας ή των ομάδων των προτύπων που έχουν μία τελική κοινή χαρακτηριστική ιδιότητα (ταυτότητα: label) αποτελεί την ταξινόμηση των προτύπων και οδηγεί στην αναγνώρισή τους. Στην γενική περίπτωση το πρόβλημα της ταξινόμησης είναι δύσκολο κυρίως λόγω α) κακής επιλογής των μετρούμενων χαρακτηριστικών, β) μεγάλης ποικιλομορφίας των προτύπων, γ) θορύβου, δ) της ασάφειας των ορίων των κλάσεων στον χώρο των προτύπων.

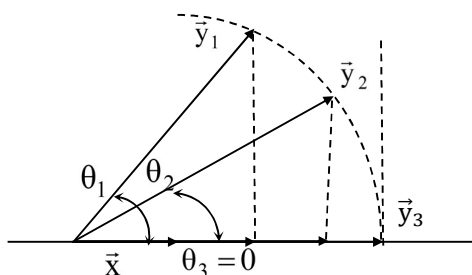
## 2.2 Εσωτερικό γινόμενο

Για δύο πίνακες  $\mathbf{x} = [x_1, x_2, \dots, x_n, \dots, x_N]^T$ , και  $\mathbf{y} = [y_1, y_2, \dots, y_n, \dots, y_N]^T$  το εσωτερικό γινόμενό τους ορίζεται από τη σχέση

$$\mathbf{x}^T \cdot \mathbf{y} = \sum_{v=1}^N x_v \cdot y_v = x_1 \cdot y_1 + x_2 \cdot y_2 + \dots + x_n \cdot y_n + \dots + x_N \cdot y_N = |\mathbf{x}| |\mathbf{y}| \cos(\theta) \quad (2.2.1)$$

Αν  $N=2$  ή  $N=3$ ,  $\theta$  είναι η κυρτή γωνία που σχηματίζουν τα δύο διανύσματα με την γεωμετρική έννοια. Αν  $N > 3$ ,  $\theta = \cos^{-1} \frac{\mathbf{x}^T \cdot \mathbf{y}}{|\mathbf{x}| |\mathbf{y}|}$  με  $|\mathbf{x}| = \sqrt{\mathbf{x}^T \cdot \mathbf{x}}$  το μέτρο του  $\mathbf{x}$ .

Γεωμετρικά αν το διάνυσμα  $\vec{\mathbf{X}}$  έχει μέτρο (μήκος) ίσο με τη μονάδα, το γινόμενο  $\vec{\mathbf{X}} \cdot \vec{\mathbf{y}}$  είναι η αλγεβρική τιμή της προβολής του  $\vec{\mathbf{y}}$  επάνω στην ευθεία που διέρχεται από το διάνυσμα  $\vec{\mathbf{X}}$  (φορέας του  $\vec{\mathbf{X}}$ ). Όταν η γωνία  $\theta$  των  $\vec{\mathbf{X}}$  και  $\vec{\mathbf{y}}$  μικραίνει το εσωτερικό τους γινόμενο  $\vec{\mathbf{X}} \cdot \vec{\mathbf{y}}$  αυξάνει και μεγιστοποιείται όταν τα  $\vec{\mathbf{X}}$  και  $\vec{\mathbf{y}}$  βρίσκονται στην ίδια ευθεία και έχουν την ίδια φορά (Σχ.2.2).



Σχήμα 2.2.

Αυτό σημαίνει ότι τότε ισχύει η σχέση

$$\vec{y} = \lambda \cdot \vec{x}, \quad \mathbf{y} = \lambda \cdot \mathbf{x} \quad \text{όπου } \lambda \in \mathbb{R} \quad \text{και} \quad \frac{y_1}{x_1} = \frac{y_2}{x_2} = \dots = \frac{y_v}{x_v} = \dots = \frac{y_N}{x_N} = \lambda \quad (2.2.2)$$

Η σχέση (2.2.2) είναι σχέση αναλογίας και το εσωτερικό γινόμενο μπορεί να χρησιμοποιηθεί ως κριτήριο ομοιότητας (όχι υποχρεωτικά ισότητας) δύο προτύπων. Αν θεωρήσουμε  $\mathbf{x} = [x_1, x_2, \dots, x_v, \dots, x_N]^T$ ,  $\mathbf{y} = [y_1, y_2, \dots, y_v, \dots, y_N]^T$ , τότε το εσωτερικό γινόμενο τους δίνεται από τη σχέση  $\mathbf{x}^T \mathbf{y} = \mathbf{y}^T \mathbf{x}$ .

### 2.3 Αποστάσεις

Εκτός από την Ευκλείδεια απόσταση υπάρχουν και άλλοι τύποι αποστάσεων σε διανυσματικούς χώρους. Αν  $d(\vec{x}, \vec{y})$  είναι απόσταση θα πρέπει για οποιοδήποτε διάνυσμα  $\mathbf{x}, \mathbf{y}, \mathbf{z}$  να ικανοποιούνται οι σχέσεις:

$$\begin{aligned} d(\vec{x}, \vec{y}) &= d(\vec{y}, \vec{x}) \\ d(\vec{x}, \vec{y}) &\leq d(\vec{x}, \vec{z}) + d(\vec{z}, \vec{y}) \\ d(\vec{x}, \vec{y}) &\geq 0 \\ \text{Αν } d(\vec{x}, \vec{y}) &= 0 \Leftrightarrow \vec{x} = \vec{y} \end{aligned} \quad (2.3.4)$$

Συνήθεις τύποι αποστάσεων είναι οι παρακάτω:

α) Minkowski τάξης s

$$d_\mu(\vec{x}, \vec{y}) = [\sum_{v=1}^N |x_v - y_v|^s]^{1/s} \quad (2.3.5)$$

β) City Block

Είναι ειδική περίπτωση της Minkowski με s=1

$$d_c(\vec{x}, \vec{y}) = \sum_{v=1}^N |x_v - y_v| \quad (2.3.6)$$

γ) Ευκλείδεια

Είναι ειδική περίπτωση της Minkowski για  $s=2$

$$d_\varepsilon(\vec{x}, \vec{y}) = \left[ \sum_{v=1}^N (x_v - y_v)^2 \right]^{1/2} = [(\mathbf{x} - \mathbf{y})^T \cdot (\mathbf{x} - \mathbf{y})]^{1/2} \quad (2.3.7)$$

δ) Chebychev

$$d_T(\vec{x}, \vec{y}) = \max(|x_v - y_v|) \quad (2.3.8)$$

ε) Mahalanobis

$$d_R(\vec{x}, \vec{y}) = d_R(\mathbf{x}, \mathbf{y}) = (\mathbf{x} - \mathbf{y})^T \cdot \text{Cov}^{-1} \cdot (\mathbf{x} - \mathbf{y}) \quad (2.3.9)$$

όπου Cov ο πίνακας συμμεταβλητότητας των  $\mathbf{x}$  και  $\mathbf{y}$ .

στ) Μη γραμμική (Non Linear)

$$d_{NL}(\vec{x}, \vec{y}) = \begin{cases} 0 & \alpha \nu \ d(x,y) \leq T \\ H & \alpha \nu \ d(x,y) > T \end{cases} \quad (2.3.10)$$

όπου  $H, T \in \mathbb{R}$  παράμετροι της απόστασης και  $d(\vec{x}, \vec{y})$  μια άλλη απόσταση του χώρου.

## ΚΕΦΑΛΑΙΟ 3

### ΜΕΘΟΔΟΙ ΤΑΞΙΝΟΜΗΣΗΣ – ΕΚΠΑΙΔΕΥΣΗ ΜΕ ΕΠΟΠΤΗ

#### 3.1. Ταξινόμηση με βάση την απόσταση από τους K-γείτονες

Μία απλή και μέθοδος ταξινόμησης ενός προτύπου  $\mathbf{x}$  βασίζεται στην εύρεση του κοντινότερου γείτονά του που ανήκει στο αρχικό σύνολο προτύπων γνωστής κλάσης. Σύμφωνα με την μέθοδο αν  $\mathbf{x}_c$  είναι ο γείτονας αυτός και ανήκει στην κλάση C, τότε το  $\mathbf{x}$  ταξινομείται στην κλάση C. Η μέθοδος ισχυροποιείται αν λάβουμε υπόψη περισσότερους από έναν γείτονες σε περιπτώ κατά προτίμηση αριθμό K και αποδώσουμε στο  $\mathbf{x}$  την κλάση της πλειοψηφίας αυτών.

```
% K nearest neighbor's classification example in Matlab R2018a
load fisheriris
x = [4.8 3.5 1.5 0.2];
[idx D] = knnsearch(meas,x,'k',3);
species(idx)
D
% Result Presentation
plot(meas( 1: 50,1),meas( 1: 50,2),'ro');hold on;
plot(meas( 51:100,1),meas( 51:100,2),'go')
plot(meas(101:150,1),meas(101:150,2),'bo')
plot(x(1),x(2),'kx')
plot(meas( idx,1),meas(idx,2),'k+')
```

Ο αλγόριθμος των K πλησιέστερων γειτόνων (KNN: K-Nearest Neighbors) είναι ακριβής και απλός στην υλοποίησή του. Σημαντικά μειονεκτήματά του είναι η διατήρηση στην μνήμη των προτύπων του συνόλου εκπαίδευσης και ο υπολογισμός των αποστάσεων του υπό ταξινόμηση προτύπου από αυτά του συνόλου εκπαίδευσης. Υπάρχουν προτάσεις για την βελτίωση των δύο αυτών μειονεκτημάτων. Ειδικά για

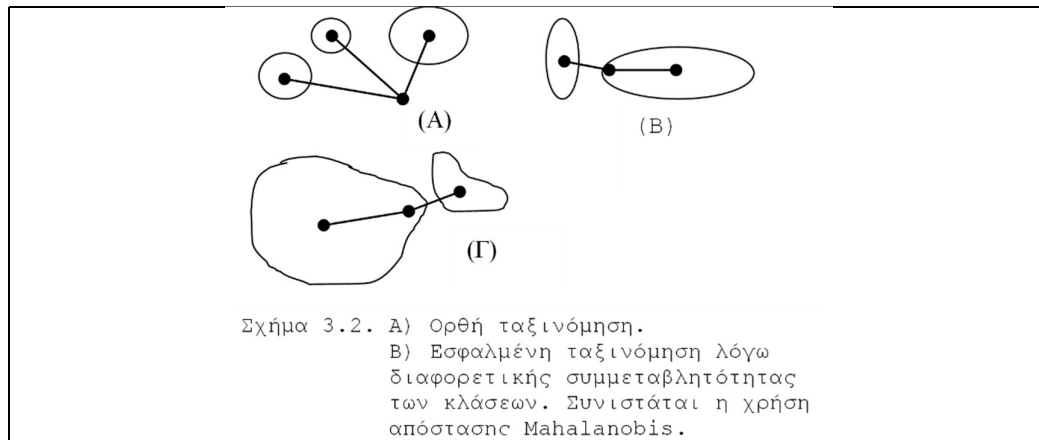
τον υπολογισμό της εύρεσης των κοντινότερων γειτόνων επιστρατεύεται ο αλγόριθμος *k-d tree* .

### 3.2 Αναγνώριση με βάση τα κέντρα των τάξεων

Είναι μια απλοϊκή μέθοδος αναγνώρισης με επόπτη σύμφωνα με την οποία κατά την εκμάθηση υπολογίζεται η μέση τιμή της κάθε κλάσης (κέντρο της κλάσης). Κάθε νέο πρότυπο ταξινομείται στην κλάση που το κέντρο της απέχει λιγότερο από το πρότυπο. Αν  $T$  το πλήθος των τάξεων με κέντρα  $\mu_t, t=1,2,\dots,T$  το πρότυπο με πίνακα  $\mathbf{x}$  αποδίδεται στην κλάση  $\min_t(d(\mathbf{x}, \mu_t))$ . Η ευκλείδεια απόσταση  $d_E(\mathbf{x}, \mu_i)^2 = (\mathbf{x} - \mu_i)^T(\mathbf{x} - \mu_i) = \mathbf{x}^T\mathbf{x} - 2\mathbf{x}^T\mu_i + \mu_i^T\mu_i$ . Για την σύγκριση των αποστάσεων του  $\mathbf{x}$  από τις κλάσεις  $i$  και  $j$  αρκεί ο υπολογισμός της διαφοράς τους

$$d_{ij} = d_E(\mathbf{x}, \mu_i)^2 - d_E(\mathbf{x}, \mu_j)^2 = 2\mathbf{x}^T(\mu_j - \mu_i) + \mu_i^T\mu_i - \mu_j^T\mu_j. \quad (\alpha)$$

Έτσι η σύγκριση απαιτεί κάθε φορά κυρίως τον υπολογισμό του εσωτερικού γινομένου  $\mathbf{x}^T(\mu_j - \mu_i)$ . Αν  $d_{ij}=0$  τότε η σχέση (α) περιγράφει το υπερ-επίπεδο που είναι μεσοκάθετο στο ευθύγραμμο τμήμα με άκρα τα κέντρα των κλάσεων. Η μέθοδος αυτή μπορεί να είναι αποδοτική όταν οι κλάσεις έχουν τη μορφή υπερσφαιρών με ακτίνες μικρότερες της ημιαπόστασης των κέντρων τους (Σχ.3.2).



### 3.3 Ταξινόμηση σε δύο κλάσεις με γραμμική διακριτική συνάρτηση, ο νευρώνας Perceptron.

Ο υπολογισμός των αποστάσεων είναι υπολογιστικά δαπανηρός και στην περίπτωση του KNN μνημοβόρος. Μία άλλη προσέγγιση υπολογιστικά ταχύτερη με μικρότερες απαιτήσεις μνήμης για την διαδικασία ταξινόμησης μπορεί να βασισθεί στον διαχωρισμό των κλάσεων με γραμμικές διακριτικές συναρτήσεις.

Αν  $\mathbf{x}$  το διάνυσμα στήλης που περιγράφει ένα πρότυπο και  $\mathbf{w}$  διάνυσμα παραμέτρων,  $\mathbf{x}, \mathbf{w} \in \mathcal{R}^N, w_0 \in \mathcal{R}$  με  $N$  το πλήθος των χαρακτηριστικών, η σχέση

$$d(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

αποτελεί μία γραμμική διακριτική συνάρτηση. Αν  $d(\mathbf{x})=0$  τότε  $\mathbf{x}$  βρίσκεται σε μία ευθεία αν  $N=2$ , σε ένα επίπεδο εάν  $N=3$  ή ένα υπερεπίπεδο για  $N>3$ . Η  $d(\mathbf{x})=0$  χωρίζει τον  $\mathcal{R}^N$  σε δύο μέρη. Το  $H^+ = \{\mathbf{x} / \mathbf{x}, d(\mathbf{x}) > 0\}$  που περιλαμβάνει όλα τα  $\mathbf{x}$  που καθιστούν την  $d(\mathbf{x}) > 0$  και το  $H^-$  που περιλαμβάνει όλα τα  $\mathbf{x}$  που καθιστούν την  $d(\mathbf{x}) < 0$ . Αν το άνυσμα  $\mathbf{w}$  των παραμέτρων έχει κατάλληλες τιμές έτσι ώστε τα πρότυπα της μίας κλάσης να βρίσκονται στον ένα μέρος (πχ  $H^+$ ) και της άλλης στο άλλο, τότε ένα άγνωστο πρότυπο  $\mathbf{x}'$  ταξινομείται με κριτήριο την τιμή  $d(\mathbf{x}')$ , αν δηλαδή  $d(\mathbf{x}') > 0$  ή  $d(\mathbf{x}') < 0$ .

Ακολούθως θα παρουσιάσουμε τα παραπάνω και ακολούθως τον αλγόριθμο εύρεσης των  $\mathbf{w}$  και  $w_0$  για  $N=2$  ώστε να είναι δυνατή η σχηματική και γεωμετρική αναπαράσταση στο επίπεδο. Αυτό δεν καταργεί την γενικότητα για  $N>2$  και βοηθά στην κατανόηση και θεώρηση του πράγματος γεωμετρικά.

Η ευθεία ( $\varepsilon$ ) στο Σχ.3.3.1 χωρίζει το επίπεδο σε δύο ημιεπίπεδα κάθε ένα των οποίων περιέχει και μια κλάση. Κάθε σημείο της ευθείας είναι το πέρας ενός διανύσματος  $\vec{\mathbf{X}}(x_1, x_2)$  που ικανοποιεί την εξίσωση της ευθείας που είναι

$$w_1 \cdot x_1 + w_2 \cdot x_2 + w_0 = 0 \quad (3.3.1)$$

ή

$$d(\mathbf{x}) = 0 \text{ όπου}$$

$$d(\mathbf{x}) = w_1 \cdot x_1 + w_2 \cdot x_2 + w_0 = \mathbf{w}^T \mathbf{x} + w_0 \quad (3.3.2)$$

με  $\mathbf{x} = [x_1, x_2]^T$  και  $\mathbf{w} = [w_1, w_2]^T$

Τα διανύσματα είναι ελεύθερα και ανήκουν σε κλάσεις ισοδυναμίας. Αν  $\mathbf{x}_A$  και  $\mathbf{x}_B$  ανήκουν στην ευθεία θα ισχύουν οι σχέσεις

$$\mathbf{w}^T \mathbf{x}_A + w_0 = 0 \text{ και} \quad (3.3.3)$$

$$\mathbf{w}^T \mathbf{x}_B + w_0 = 0 \quad (3.3.4)$$

από όπου με αφαίρεση κατά μέλη προκύπτει ότι

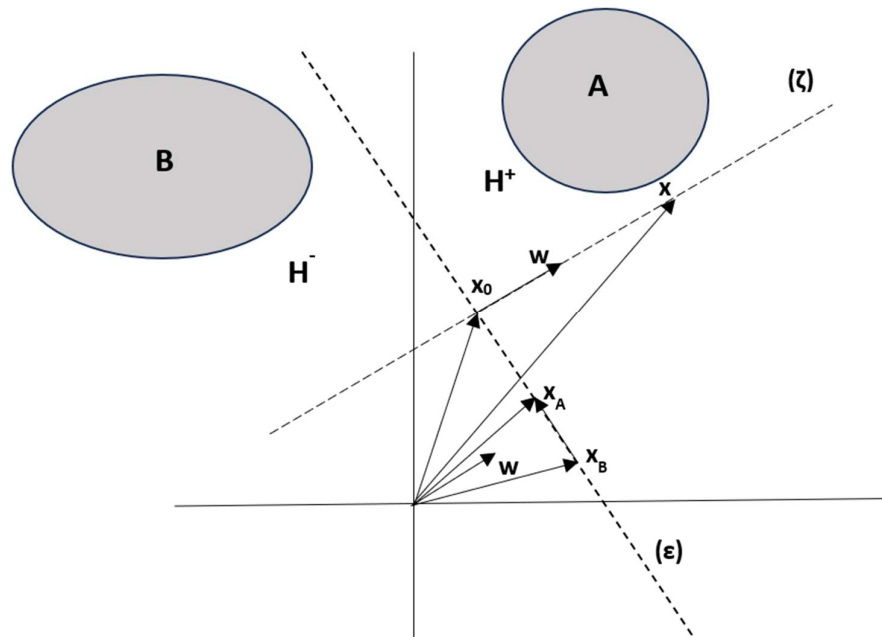
$$\mathbf{w}^T (\mathbf{x}_A - \mathbf{x}_B) = 0 \quad (3.3.5)$$

Η σχέση αυτή δείχνει ότι το άνυσμα  $\mathbf{w}$  είναι κάθετο στην ευθεία (ή το επίπεδο αν το  $\mathbf{x}$  ανήκει στον τρισδιάστατο χώρο). Όλα τα σημεία του ενός ημιεπιπέδου καθιστούν την  $d(\mathbf{x}) > 0$  και του άλλου ημιεπιπέδου την  $d(\mathbf{x}) < 0$ . Πράγματι, αν  $\mathbf{x}$  διάνυσμα με πέρας στο ένα ημιεπίπεδο, η ευθεία ( $\zeta$ ) που διέρχεται από το  $\mathbf{x}$  και είναι κάθετη στην ευθεία ( $\varepsilon$ ) δίνεται από την σχέση

$$(\zeta) \rightarrow \mathbf{x} = \mathbf{x}_0 + \lambda \cdot \mathbf{w}, \lambda \in \mathcal{R}$$

όπου  $\mathbf{x}_0$  είναι το διάνυσμα με πέρας το σημείο τομής της ( $\zeta$ ) με την ( $\varepsilon$ )

$$\mathbf{x} - \mathbf{x}_0 = \lambda \mathbf{w}$$



Σχήμα 3.3.1

Η ποσότητα  $D = \lambda |\mathbf{w}|$  είναι η προσημασμένη απόσταση του  $\mathbf{x}$  από την  $(\varepsilon)$  μετρημένη σε  $|\mathbf{w}|$ . Αν  $\lambda > 0 \Leftrightarrow D > 0$  και το  $\mathbf{x}$  βρίσκεται στο ημιεπίπεδο που δείχνει το άνωσμα  $\mathbf{w}$  και αν  $\lambda < 0 \Leftrightarrow D < 0$  και το  $\mathbf{x}$  βρίσκεται στο άλλο ημιεπίπεδο.

$$d(\mathbf{x}_0) = \mathbf{w}^T \mathbf{x}_0 + w_0 = 0$$

$$d(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = \mathbf{w}^T (\mathbf{x}_0 + \lambda \mathbf{w}) + w_0 = \quad (3.3.6)$$

$$\mathbf{w}^T \mathbf{x}_0 + w_0 + \lambda \mathbf{w}^T \mathbf{w} = \lambda |\mathbf{w}|^2 \Rightarrow$$

$$\frac{d(\mathbf{x})}{|\mathbf{w}|} = \lambda \cdot |\mathbf{w}| = D \quad (3.3.7)$$

Συνεπώς αν  $d(\mathbf{x}) > 0$  το πέρας του  $\mathbf{x}$  βρίσκεται στο ημιεπίπεδο που δείχνει το άνωσμα  $\mathbf{w}$  ( $H^+$ ) ενώ αν  $d(\mathbf{x}) < 0$  το πέρας του  $\mathbf{x}$  βρίσκεται στο άλλο ημιεπίπεδο ( $H^-$ ).

Είναι επίσης προφανές πως αν  $|\mathbf{w}| = 1$  (από την 3.3.7), τότε η ποσότητα  $d(\mathbf{x})$  είναι η προσημασμένη απόσταση του  $\mathbf{x}$  από την ευθεία  $(\varepsilon)$ .

Ζητούμε συνεπώς να βρούμε τα κατάλληλα  $\mathbf{w}$  και  $w_0$  για τα δοθέντα γνωστά  $\mathbf{x}$  που ανήκουν στις δύο κλάσεις. Ένας τρόπος να γίνει αυτό είναι μία επαναληπτική διαδικασία διόρθωσης σφάλματος που θα περιγράψουμε ακολούθως και λέγεται εκπαίδευση του ταξινομητή.

Για την απλούστερη μαθηματική διατύπωση του μηχανισμού διόρθωσης θα θεωρήσουμε τους επαυξημένους πίνακες των προτύπων. Η γραμμική διακριτική συνάρτηση θα γράφεται ως ακολούθως

$$d(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = [\mathbf{w}^T, w_0] \cdot \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$



$$d(\mathbf{x}) = \tilde{\mathbf{w}}^T \cdot \tilde{\mathbf{x}}$$

$$\tilde{\mathbf{w}} = \begin{bmatrix} \mathbf{w} \\ w_0 \end{bmatrix}, \quad \tilde{\mathbf{x}} = \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

Αν δηλαδή ο χώρος των πινάκων των προτύπων έχει  $N$  διαστάσεις, οι σχέσεις θα γραφούν στον χώρο  $R^{N+1}$  που έχει  $N+1$  διαστάσεις.

Εστω ότι ο επόπιτης επιθυμεί  $d(\mathbf{x}^A) > 0$  και  $d(\mathbf{x}^B) < 0$ . Αρχικά στο  $\mathbf{w}$  αποδίδονται τυχαίες τιμές και πιθανότατα η  $d(\mathbf{x})$  δεν διαχωρίζει τις κλάσεις. Η βασική ιδέα της εκμάθησης είναι η διόρθωση του υπερεπιπέδου  $\tilde{\mathbf{w}}^T \cdot \tilde{\mathbf{x}} = 0$  κάθε φορά που κάποιο πρότυπο ταξινομείται εσφαλμένα ώστε τελικά να αφήνει το πρότυπο προς την πλευρά των ορθά ταξινομημένων προτύπων. Η διόρθωση αυτή γίνεται σύμφωνα με την σχέση

$$\tilde{\mathbf{w}}(t+1) = \tilde{\mathbf{w}}(t) \pm \rho \cdot \tilde{\mathbf{x}} \quad (3.3.12)$$

Το πρόσθετο (+) χρησιμοποιείται όταν ταξινομείται λάθος ένα πρότυπο της κλάσης A και (-) της κλάσης B. Η ποσότητα  $\rho$  λέγεται παράμετρος ή ρυθμός εκμάθησης και παίρνει μικρές τιμές στο διάστημα  $(0,1)$ , π.χ.  $\rho=0.2$ . Ο  $t$  είναι ένας μετρητής επανάληψης της διαδικασίας εκμάθησης,  $t=0,1,2,3 \dots$ .

Ο μηχανισμός βασίζεται στο γεγονός ότι η ποσότητα  $\rho \cdot \tilde{\mathbf{x}}^T \tilde{\mathbf{x}}$  είναι πάντα θετική

$$d_{t+1}(\mathbf{x}) = \tilde{\mathbf{w}}(t+1)^T \cdot \tilde{\mathbf{x}} = (\tilde{\mathbf{w}}(t) \pm \rho \cdot \tilde{\mathbf{x}})^T \cdot \tilde{\mathbf{x}} = d_t(\mathbf{x}) \pm \rho \cdot \tilde{\mathbf{x}}^T \tilde{\mathbf{x}} \quad (3.3.12)$$

η διόρθωση μεταβάλλει την εσφαλμένη απόσταση του προτύπου προς την πλευρά του ημιεπιπέδου της κλάσης και οδηγεί επαναληπτικά στο επιθυμητό αποτέλεσμα. Η παράμετρος  $\rho$  ρυθμίζει το βήμα διόρθωσης και την ταχύτητα σύγκλισης (Σχ. 3.3.2).

Τα παραπάνω παρουσιάζονται ακολούθως βηματικά.

Βήμα 1ο: Αρχικοποιήσεις:

Ορίζουμε ένα μετρητή επανάληψης  $t=0,1,2,3, \dots$  και του αποδίδουμε αρχικά την τιμή μηδέν. Αποδίδουμε τυχαίες (συνήθως μικρές θετικές) τιμές στο άνωμα  $\mathbf{w}_0 = [w_1, w_2, \dots, w_N, w_0]$ , δημιουργούμε τα επαυξημένα διανύσματα του συνόλου εκπαίδευσης  $\tilde{\mathbf{x}} = [\mathbf{x}, 1]$  και θεωρούμε την γραμμική διακριτική συνάρτηση  $D_t(\tilde{\mathbf{x}}) = \mathbf{w}_t^T \cdot \tilde{\mathbf{x}}$ . Ζητούμε κατά σύμβαση  $D(\tilde{\mathbf{x}}^A) > 0$  και  $D(\tilde{\mathbf{x}}^B) < 0$

Βήμα 2ο: Αυξάνουμε τον μετρητή επανάληψης και επιλέγουμε ένα τυχαίο πρότυπο με πίνακα  $\tilde{\mathbf{X}}$ , από το σύνολο εκπαίδευσης και υπολογίζουμε την ποσότητα  $D_t(\tilde{\mathbf{X}})$ . Αν το πρότυπο είναι σωστά ταξινομημένο επαναλαμβάνουμε το Βήμα 2 αλλιώς διορθώνουμε τα βάρη σύμφωνα με την σχέση

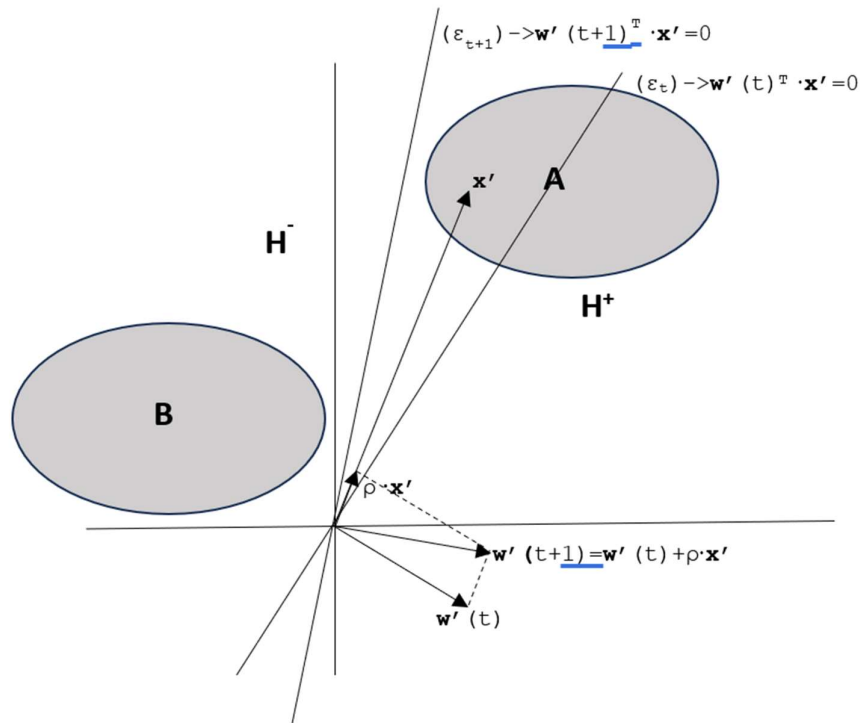
$$\tilde{\mathbf{w}}(t+1) = \tilde{\mathbf{w}}(t) \pm \rho \cdot \tilde{\mathbf{x}} \quad (3.3.13)$$

Η χωρίς την χρήση επαυξημένων διανυσμάτων

$$\mathbf{w}(t+1) = \mathbf{w}(t) \pm \rho \cdot \mathbf{x} \quad \text{και} \quad w_0(t+1) = w_0(t) \pm \rho$$

Με (+) αν το πρότυπο ανήκει στην κλάση A, με (-) στην κλάση B.

Βήμα 3ο: Αυξάνεται η τιμή του  $t$  κατά ένα και η διαδικασία επαναλαμβάνεται από το Βήμα 2 εωσότου όλα τα πρότυπα να ταξινομηθούν σωστά.



Σχήμα 3.3.2

( Το σύμβολο (´) έχει αντικαταστήσει το σύμβολο (~) στα  $\mathbf{x}$  και  $\mathbf{w}$ )

Στον αλγόριθμο το άνωσμα  $\mathbf{x}$  επιλέγεται τυχαία κάθε φορά και γι αυτό θεωρείται στοχαστικός (stochastic). Αποδεικνύεται ότι για δύο γραμμικά διαχωρίσιμες κλάσεις ο αλγόριθμος οδηγεί πάντοτε σε λύση για κατάλληλα μικρή τιμή του  $\rho$ .

Υλοποίηση με MatLab.

```
function [ weights ] =...
perceptron_stochastic( trset,class_labels,...% +1 or -1
                      learning_rate,error_tolerance )
    ftrnum=size(trset,1); % feature number
    N=size(trset,2);     % vector number

    w=rand(ftrnum,1); w0=rand(1); % Weights initialization
    while sum( sign(w'*trset+w0)==class_labels ) < N-error_tolerance
        i=randi(N)
        x = trset(:,i);
        d = w'*x + w0;
        if( d*class_labels(i) < 0);
            w = w + class_labels(i)*learning_rate*x;
            w0= w0+ learning_rate*class_labels(i);
        end
    end
    weights = [w; w0];
end
```

Σε παρόμοιο αποτέλεσμα οδηγούμαστε εάν θεωρήσουμε το πρόβλημα εύρεσης των βαρών ως πρόβλημα βελτιστοποίησης (*optimization*) μιας συνάρτησης κόστους  $K(\tilde{\mathbf{w}})$ . Μία τέτοια συνάρτηση είναι η

$$K(\tilde{\mathbf{w}}): \mathbb{R}^{N+1} \rightarrow \mathbb{R}^+ \text{ με } K(\tilde{\mathbf{w}}) = \sum_{\mathbf{x} \in S} \delta(\mathbf{x}) \tilde{\mathbf{w}}^T \tilde{\mathbf{x}} \quad (3.3.14)$$

$$\delta(\mathbf{x}) = \begin{cases} -1 & \text{αν } \mathbf{x} \in A \\ +1 & \text{αν } \mathbf{x} \in B \end{cases}$$

όπου  $S$  το σύνολο των διανυσμάτων που ταξινομήθηκαν λάθος. Η συνάρτηση  $K(\tilde{\mathbf{w}})$  εκφράζει ουσιαστικά το συνολικό σφάλμα ταξινόμησης και είναι κατά τμήματα γραμμική συνάρτηση. Η βέλτιστη λύση του προβλήματος είναι η εύρεση ανύσματος βαρών  $\tilde{\mathbf{w}}$  ώστε  $K(\tilde{\mathbf{w}}) = 0$ . Αν  $t=0,1,2,\dots$  και  $\tilde{\mathbf{w}}(0)$  είναι μία αρχική τιμή, η ελάχιστη τιμή του  $\tilde{\mathbf{w}}$  μπορεί να προσεγγισθεί επαναληπτικά με την με την μέθοδο της *επικλινούσς κατάβασης (gradient descent)* σύμφωνα με την σχέση

$$\tilde{\mathbf{w}}(t+1) = \tilde{\mathbf{w}}(t) - \rho \frac{\partial K}{\partial \tilde{\mathbf{w}}} \Big|_{\tilde{\mathbf{w}}(t)} = \tilde{\mathbf{w}}(t) - \rho \sum_{\mathbf{x} \in S} \delta(\tilde{\mathbf{x}}) \tilde{\mathbf{x}} \quad (3.3.15)$$

Η διόρθωση γίνεται αφού ληφθεί υπόψη το σύνολο  $S$  και η διαδικασία εκπαίδευσης χαρακτηρίζεται ως εκπαίδευση *δέσμης (batch)*.

Υλοποίηση του Perceptron με εκπαίδευση κατά δέσμες

```
function [ w ] = ...
perceptron_batch(trset,class_labels,...% +1 or -1
                learning_rate,error_tolerance )
X=[trset; ones(1,size(trset,2))]; % vector augmentation
w = rand(size(trset,1)+1,1)-0.5 ; % weight initialization
while true
    misclsf = sign(w'*X)~=class_labels;
    if sum( misclsf )<=error_tolerance; break; end
    s = sum( (X.*(misclsf)).*(-class_labels), 2);
    w = w- learning_rate *s;
end
end
```

Είναι χρήσιμο να εφαρμόσουμε στην τιμή του αθροιστή  $\sigma=D(\cdot)$  την συνάρτηση *προσήμου* ή την *βηματική*

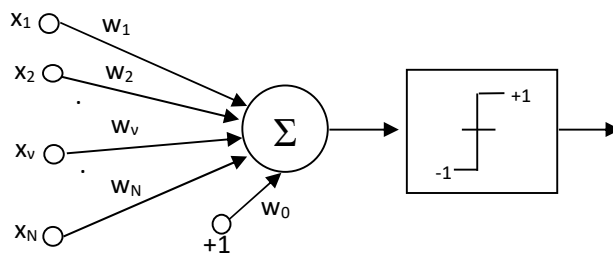
Βηματική 0/1 (step function 0/1):

$$f(\sigma) = \begin{cases} 0 & \text{αν } \sigma < 0 \\ 1 & \text{αν } \sigma \geq 0 \end{cases}$$

Προσήμευ -1/1 (sign function -1/1):

$$f(\sigma) = \begin{cases} -1 & \text{αν } \sigma < 0 \\ +1 & \text{αν } \sigma \geq 0 \end{cases}$$

ώστε να κωδικοποιείται κάθε κλάση π.χ. με +1 αν  $d(\mathbf{x}^A) \geq 0$  και το -1 ή το 0 αν  $d(\mathbf{x}^B) < 0$ . Η συναρτήσεις αυτές λέγονται και συναρτήσεις ενεργοποίησης (activation functions). Κατόπιν τούτων μπορούμε να παραστήσουμε διαγραμματικά τον ταξινομητή όπως δείχνεται στο Σχ.3.3.2.



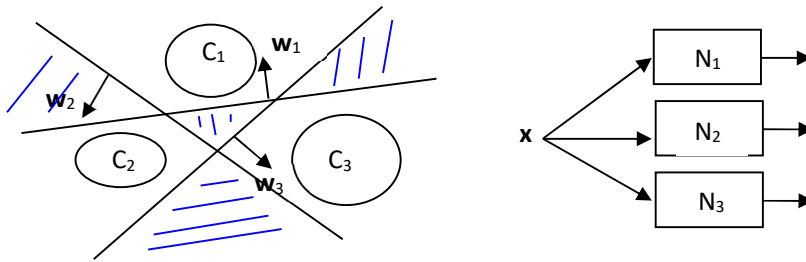
Σχ.3.3.2

Ο γραμμικός ταξινομητής είναι ένα απλό μοντέλο του βιολογικού νευρώνα και ονομάστηκε νευρώνας perceptron. Προτάθηκε αρχικά από τους McCulloch και Pitts για δυαδικές τιμές (0 ή 1) στην είσοδο και την έξοδο. Ολοκληρώθηκε από τον ψυχολόγο F. Rosenblatt, και αποτελεί το δομικό στοιχείο τεχνητών νευρωνικών δικτύων (ΤΝΔ) που μπορεί να αποτελούνται από εκατομμύρια νευρώνες.

- [1] McCulloch, W; Pitts, W (1943). "[A Logical Calculus of Ideas Immanent in Nervous Activity](#)". *Bulletin of Mathematical Biophysics*. **5** (4): 115-133.
- [2] Rosenblatt, Frank (1957). "The Perceptron—a perceiving and recognizing automaton". *Report 85-460-1*. Cornell Aeronautical Laboratory.

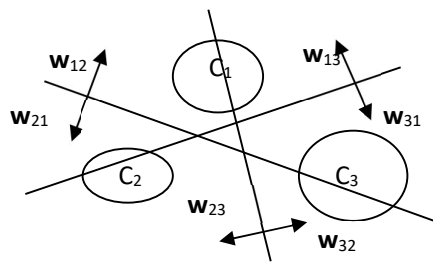
### 3.4.1. Η περίπτωση πολλών κλάσεων.

Στην περίπτωση που το πλήθος των κλάσεων είναι  $M > 2$  τότε θα πρέπει να χρησιμοποιήσουμε περισσότερους του ενός γραμμικούς ταξινομητές. Μία περίπτωση είναι να διαχωρίσουμε γραμμικά κάθε τάξη από τις υπόλοιπες (εάν αυτό είναι εφικτό) με  $M$  το πλήθος γραμμικές διακριτικές συναρτήσεις. Αν  $d_\mu(\mathbf{x}) = \mathbf{w}_\mu^T \cdot \mathbf{x} + c > 0$ ,  $\mu = 1 \dots M$ , όταν το  $\mathbf{x}$  ανήκει στη κλάση  $C_\mu$  (Σχ. 3.4.1.1), τότε ένα πρότυπο ανήκει στην κλάση της οποίας η γραμμική διακριτική συνάρτηση  $d_\mu(\mathbf{x})$  είναι θετική. Στην περίπτωση υλοποίησης με ΝΔ και βηματική έξοδο 0/1 το άγνωστο πρότυπο ταξινομείται στην κλάση που ο αντίστοιχος νευρώνας έχει έξοδο 1. Μια τέτοια προσέγγιση όπως δείχνεται στο σχήμα. 3.4.1 είναι δυνατόν να αποδώσει ένα πρότυπο σε καμία ή σε πολλές τάξεις.

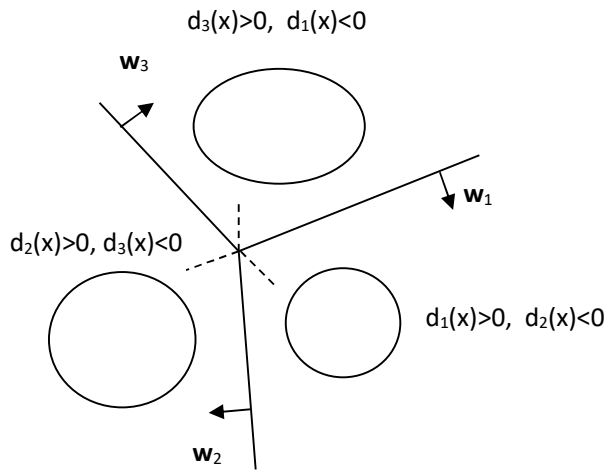


Σχ. 3.4.1.1

Ένας άλλος τρόπος αντιμετώπισης του προβλήματος είναι ο διαχωρισμός των κλάσεων ανά δύο αδιαφορώντας για τις υπόλοιπες. Το πρότυπο αποδίδεται σε εκείνη την κλάση που έχει  $M-1$  θετικές τιμές στις διακριτικές συναρτήσεις που την χωρίζουν από τις υπόλοιπες. Η μέθοδος αυτή είναι υπολογιστικά πολύπλοκη στην περίπτωση των πολλών κλάσεων (Σχ.3.4.1.2). Τέλος αν οι διαχωριστικές επιφάνειες είναι όπως στο σχήμα 3.4.1.3 αποφεύγεται η ύπαρξη διφορούμενων περιοχών στο χώρο των προτύπων. Συνδυασμοί γραμμικών ή κατά τμήματα γραμμικών ταξινομητών με την χρήση διαφόρων συναρτήσεων στην έξοδο έχουν προταθεί όπως είναι τα ΝΔ ADALINE, MADALINE οι μηχανές επιτροπής (COMMITTEE MACHINES) κ.α. Στα συστήματα αυτά δεν θα αναφερθούμε ενδελεχώς δεδομένου ότι η κεντρική ιδέα της λειτουργίας του καλύπτεται από όσα αναφέρθηκαν μέχρι τώρα. Σε επόμενο κεφάλαιο θα περιγράψουμε συστήματα που αποτελούνται από διαδοχικά επίπεδα ταξινομητών (MLP: multilayer perceptron) με παραγωγίσιμες συναρτήσεις ενεργοποίησης για την επίλυση μη γραμμικών προβλημάτων.



Σχ. 3.4.1.2



Σχ. 3.4.1.3

### 3.4.2 Η περίπτωση της XOR, ταξινομητές πολλών επιπέδων

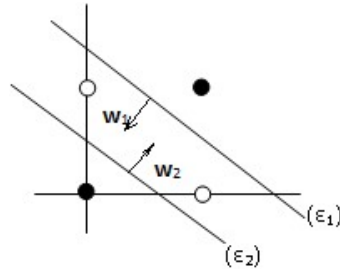
Όταν το πρόβλημα της ταξινόμησης αφορά περισσότερες από δύο κλάσεις ή όταν οι κλάσεις δεν είναι γραμμικά διαχωρίσιμες, είναι δυνατόν να επιτευχθούν λύσεις με κατάλληλους συνδυασμούς γραμμικών ταξινομητών. Μια τέτοια χαρακτηριστική περίπτωση είναι αυτή της λογικής πύλης XOR της οποίας ο πίνακας αληθείας δείχνεται στον Πιν.3.4.2.1.

A	B	A XOR B
0	0	0
0	1	1
1	0	1
1	1	0

Πίνακας 3.4.2.1.

Σύμφωνα με αυτόν οι συνδυασμοί των τιμών των λογικών μεταβλητών  $\alpha$ ,  $\beta$  αποτελούν τέσσερα πρότυπα που περιγράφονται από τα διανύσματα στοιχείων συνόλου  $\Omega = \{(0,0), (0,1), (1,0), (1,1)\}$  και η πράξη  $\alpha \text{ XOR } \beta$  ορίζει τις κλάσεις  $C_0 = \{(0,0), (1,1)\}$  και  $C_1 = \{(0,1), (1,0)\}$ . Στο Σχ.3.4.2.1 φαίνονται τα άκρα των διανυσμάτων στο  $E^2$ . Είναι προφανές ότι οι κλάσεις  $C_0, C_1$  δεν διαχωρίζονται με μια ευθεία. Ο διαχωρισμός των κλάσεων μπορεί να γίνει με δύο ευθείες

(Σχ.3.4.2.1) που ορίζουν μία ζώνη στο εσωτερικό της οποίας βρίσκονται τα πρότυπα της κλάσης  $C_1$ . Η ευθεία  $(\epsilon_1)$  μπορεί να προσδιοριστεί από έναν γραμμικό ταξινομητή  $T_1$  που θα διαχωρίζει το πρότυπο  $(1,1)$  από τα υπόλοιπα. Η ευθεία  $(\epsilon_2)$  μπορεί να προσδιορισθεί από έναν γραμμικό ταξινομητή  $T_2$  που θα διαχωρίζει το πρότυπο  $(0,0)$  από τα υπόλοιπα.



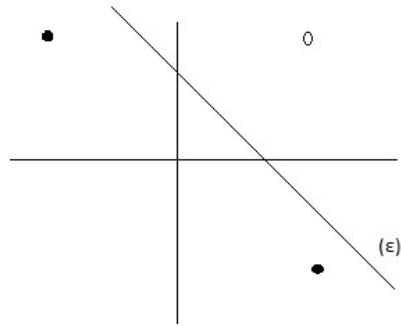
Σχ. 3.4.2.1

Οι έξοδοι των  $T_1, T_2$  θα είναι οι τιμές των συναρτήσεων  $\sigma_1=f(\mathbf{x}), \sigma_2=f(\mathbf{x})$  για  $\mathbf{x} \in \Omega$  και  $f$  την συνάρτηση προσήμου όπως φαίνονται στον ακόλουθο πίνακα 3.4.2.2 :

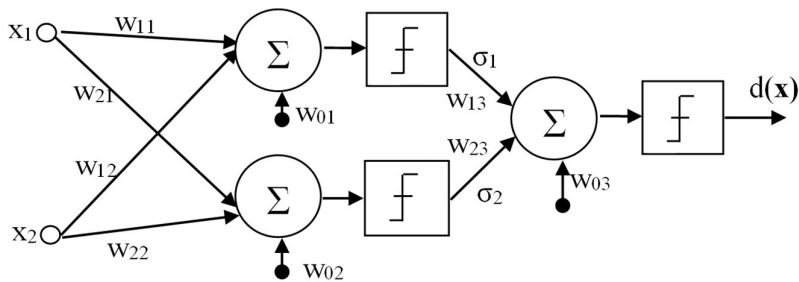
$x_1$	$x_2$	$\sigma_1$	$\sigma_2$	ΚΛΑΣΗ
0	0	+1	-1	$C_0$
0	1	+1	+1	$C_1$
1	0	+1	+1	$C_1$
1	1	-1	+1	$C_0$

Πίνακας 3.4.2.2

Οι τιμές των  $\sigma_1(\mathbf{x}), \sigma_2(\mathbf{x})$  αποτελούν ένα νέο σύνολο προτύπων  $\Phi=\{(+1,-1), (-1,+1), (+1,+1)\}$ . Στο Σχ.3.4.2.3 φαίνεται ο χώρος του  $\Phi$ . Ένας γραμμικός ταξινομητής  $T$  μπορεί να προσδιορίσει την  $\epsilon$ . Το όλο σύστημα φαίνεται στο Σχ.3.4.2.3. Αφήνεται σαν άσκηση στον αναγνώστη ο υπολογισμός των συναπτικών βαρών του ΝΔ του σχήματος 3.4.2.3 με γραμμικές διακριτικές συναρτήσεις όπως των ευθειών των σχημάτων 3.4.2.1, 3.4.2.3. Προσέξτε παρατηρώντας τον πίνακα 3.4.2.2, ότι ο ένας νευρώνας του κρυφού επιπέδου υλοποιεί την λογική πράξη OR  $(\alpha+\beta)$  και ο άλλος την λογική πράξη AND  $(\alpha \cdot \beta)$ .



Σχ.3.4.2.3.



Σχήμα 3.4.2.3. Γραμμικός ταξινομητής-NΔ Perceptron δύο επιπέδων.

Με τη χρήση περισσότερων νευρώνων παρατεταγμένων σε επίπεδα όπως στο Σχ.3.4.2.3 δημιουργούμε πολυεπίπεδα δίκτυα perceptron (*MLP: Multi Layer Perceptron*) που μπορούν να διαχωρίσουν κλάσεις μη γραμμικά διαχωρίσιμες. Τα επίπεδα ανάμεσα στο επίπεδο εισόδου των τιμών του ανύσματος  $\mathbf{x}$  και του επιπέδου των νευρώνων εξόδου, ονομάζονται κρυφά επίπεδα. Το νευρωνικό δίκτυο του Σχ.3.4.2.3 έχει ένα κρυφό επίπεδο δύο νευρώνων. Ο διαχωρισμός μη γραμμικά διαχωρίσιμων κλάσεων με την χρήση MLP όπως παρουσιάστηκε μέχρι τώρα, απαιτεί την παρέμβαση του ανθρώπου σχεδιαστή και δεν είναι μια γενική και αυτοματοποιημένη διαδικασία. Οι περισσότεροι αλγόριθμοι εκπαίδευσης βασίζονται στην εύρεση του ελαχίστου μιας κατάλληλης συνάρτησης κόστους των συναπτικών βαρών. Μια τέτοια προσέγγιση βασίζεται στην κατάβαση αντίθετα από την κλίση της συνάρτησης αυτής (*gradient descent*). Οι βηματικές συναρτήσεις που χρησιμοποιήθηκαν μέχρι τώρα δεν είναι παραγωγίσιμες στο μηδέν. Άλλες συναρτήσεις παραγωγίσιμες σε όλο το πεδίο ορισμού τους μπορούν να χρησιμοποιηθούν και να δώσουν ικανοποιητική λύση όπως θα δούμε ακολούθως.



### 3.5 Πολυεπίπεδοι ταξινομητές – Διόρθωση σφάλματος με οπισθοδιάδοση (Back Error Propagation)

Ο perceptron αν και εντυπωσίασε με την ομοιότητα που εμφάνιζε με τα βιολογικά νευρωνικά δίκτυα εν τούτοις δεν μπορούσε να επιλύσει αυτόματα ακόμη και απλά προβλήματα. Για περισσότερο από μια δεκαετία η εφαρμογή του και το έντονο ερευνητικό ενδιαφέρον ατόνησαν αλλά δεν εγκαταλείφθηκαν. Οι Rumelhart et al<sup>[3]</sup> πρότειναν μια νέα εκδοχή MLP ΤΝΔ που έδωσε μεγάλη ώθηση στον αντίστοιχο επιστημονικό χώρο.

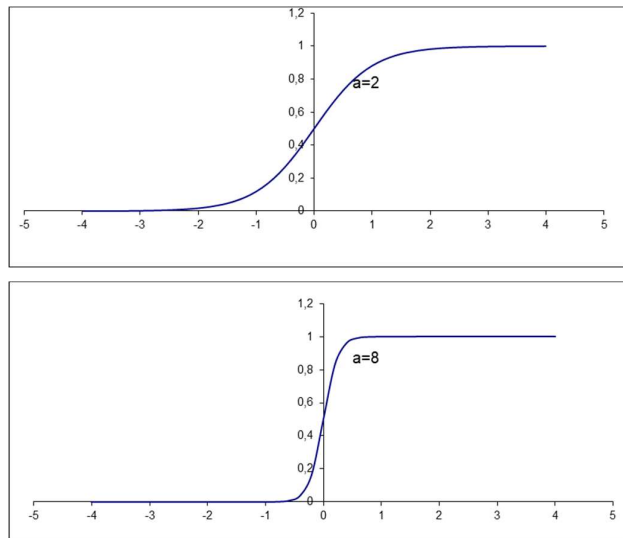
[3] [Rumelhart, David E.](#); [Hinton, Geoffrey E.](#); [Williams, Ronald J.](#) (1986a). "Learning representations by back-propagating errors". *Nature*. **323** (6088): 533–536.

Στην νέα αυτή εκδοχή η βηματική συνάρτηση εξόδου αντικαθίσταται από συνεχείς και παραγωγίσιμες συναρτήσεις που την προσεγγίζουν. Τέτοιες συναρτήσεις είναι οι ακόλουθες:

- Σιγμοειδής-λογιστική (sigmoid logistic):

$$f(\sigma) = \frac{1}{1+e^{-a\sigma}} \quad (3.5.1)$$

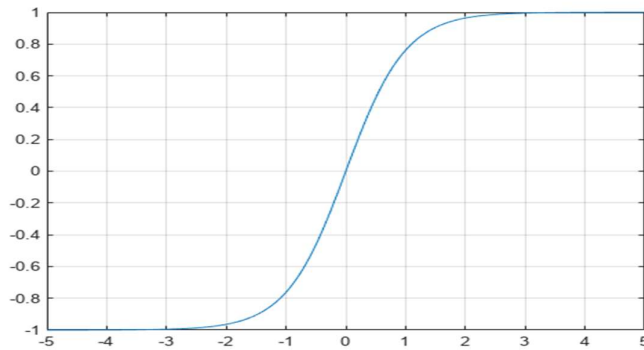
Η παράμετρος  $a$  λέγεται παράμετρος κλίσης (slope) ή λοξότητα, (Σχ.3.5.1)



Σχήμα.3.5.1

- Σιγμοειδής-Υπερβολική εφαπτομένη (hyperbolic tangent):

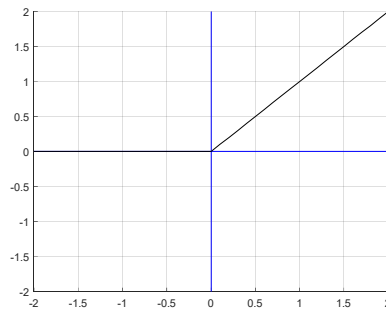
$$f(\sigma) = \tanh(\sigma) = \frac{1-e^{-2\sigma}}{1+e^{-2\sigma}}$$



Σχήμα. Η συνάρτηση της υπερβολικής εφαπτομένης

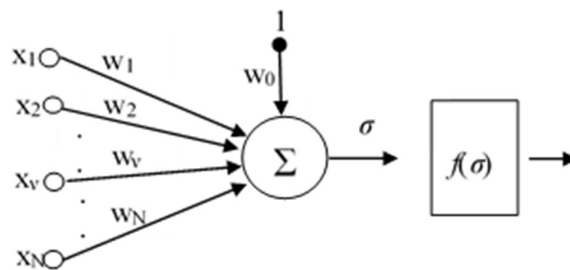
- Συνάρτηση ράμπας (ramp function):  $f(\sigma) = \begin{cases} 0 & \text{αν } \sigma < 0 \\ \sigma & \text{αν } \sigma \geq 0 \end{cases} = \max(0, \sigma)$ .

Ονομάζεται και (ReLU: Rectified Linear Unit). Η ReLU χρησιμοποιείται κατά κόρον στα συνελκτικικά νευρωνικά δίκτυα που θα παρουσιασθούν σε επόμενο κεφάλαιο.



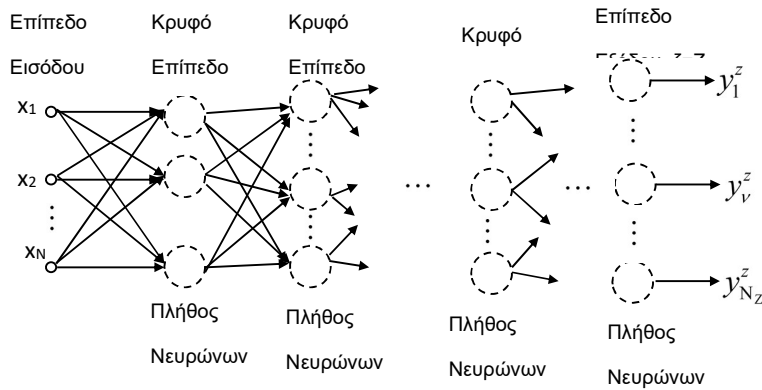
Σχήμα: Η συνάρτηση ReLU.

Η χρήση της  $f(x)$  διευκολύνει την εφαρμογή μιας μεθόδου ελαχιστοποίησης της κατάλληλης συνάρτησης κόστους. Ο τεχνητός νευρώνας έχει τώρα την μορφή του Σχ. 3.5.2



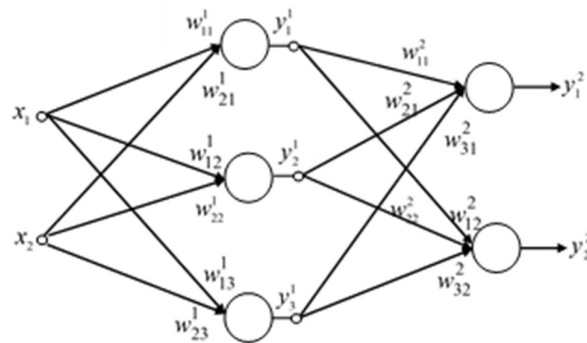
Σχ.3.5.2.

Το όλο ΝΔ όπως θα παρουσιασθεί και θα αναλυθεί ακολούθως βασίζεται στα παραπάνω. Είναι ένας πολυεπίπεδος perceptron και ονομάζεται λόγω της μεθόδου εκπαίδευσής του, ΝΔ Back error propagation ή για συντομία Back propagation. Στο νευρωνικό δίκτυο υπάρχουν επίπεδα δηλαδή ομάδες νευρώνων, πλήρως διασυνδεδεμένα μεταξύ τους Σχ.3.5.3. Η έξοδος κάθε νευρώνα ενός επιπέδου εισέρχεται σε κάθε νευρώνα του επόμενου επιπέδου, με αντίστοιχη σύναψη (βάρους) Σχ.3.5.2. Οι τιμές του πίνακα εισόδου αποτελούν το επίπεδο εισόδου. Εκτός από το επίπεδο εισόδου υπάρχουν Z το πλήθος επόμενα διαδοχικά επίπεδα, το τελευταίο των οποίων ονομάζεται επίπεδο εξόδου. Τα επίπεδα που ενδεχομένως υπάρχουν μεταξύ των επιπέδων εισόδου και εξόδου, λέγονται κρυφά επίπεδα.



Σχ.3.5.3.

Στο Σχ. 3.5.4. φαίνεται ένα νευρωνικό όπου το πρότυπο εισόδου έχει δύο χαρακτηριστικά  $[x_1, x_2]^T$  που εισέρχονται σε ένα κρυφό επίπεδο τριών νευρώνων οι έξοδοι των οποίων είναι εισοδοί ενός επόμενου επιπέδου με δύο νευρώνες και δύο τελικές εξόδους.



Σχ.3.5.4.

Για την αρίθμηση των επιπέδων θα χρησιμοποιούμε τον δείκτη  $\zeta$  με τιμές  $\zeta=0, \dots, Z$ , με  $\zeta=0$  για το επίπεδο εισόδου και  $\zeta=Z$  για το επίπεδο εξόδου Σχ. 3.5.4. Ονομάζουμε  $N_\zeta$  το πλήθος των νευρώνων ενός επιπέδου  $\zeta$ . Ως εκ τούτου το πλήθος των στοιχείων του πίνακα εισόδου  $\mathbf{x}$  είναι  $N_0$  και του πίνακα εξόδου  $\mathbf{y}$  είναι  $N_Z$ . Αν  $\nu$  είναι δείκτης για αρίθμηση των νευρώνων του επιπέδου  $\zeta$ , τότε οι συνάψεις του συνθέτουν έναν πίνακα γραμμής  $\mathbf{w}_\nu^\zeta$  και  $\mu$  δείκτης για την αρίθμηση των νευρώνων του προηγούμενου επιπέδου  $\zeta-1$  που αποτελείται από  $N_{\zeta-1}$  νευρώνες, τότε

$$\mathbf{w}_\nu^\zeta = [w_{1\nu}^\zeta, \dots, w_{\mu\nu}^\zeta, \dots, w_{N_{\zeta-1}\nu}^\zeta, w_{0\nu}^\zeta]^T$$

Η έξοδος του αθροιστή του νευρώνα θα είναι  $\sigma_\nu^\zeta$

$$\sigma_\nu^\zeta = \sum_{\mu=1}^{N_{\zeta-1}} w_{\mu\nu}^\zeta \cdot y_\mu^{\zeta-1} + w_{0\nu}^\zeta$$

ή ακόμη  $\sigma_\nu^\zeta = (\mathbf{y}^{\zeta-1})^T \cdot \mathbf{w}_\nu^\zeta$  με  $\mathbf{y}^{\zeta-1} = [y_1, \dots, y_{N_{\zeta-1}}, 1]^T$

Η τελική έξοδος του νευρώνα,  $y_\nu^\zeta$  προκύπτει από την σχέση

$$y_\nu^\zeta = f(\sigma_\nu^\zeta)$$

όπου  $f(\cdot)$  σιγμοειδής συνάρτηση, π.χ. η λογιστική. Στο Σχ. 3.5.5 φαίνεται αναλυτικά η δομή ενός νευρώνα όπως περιγράφηκε παραπάνω.

### Εκπαίδευση με οπισθοδρόμηση του σφάλματος (back-error propagation)

Εστώ  $I$  το πλήθος ζευγών από πίνακες εισόδου και των αντίστοιχων επιθυμητών πινάκων εξόδου με γνωστές τιμές, το σύνολο εκπαίδευσης  $S$  ορίζεται ως  $S = \{(\mathbf{x}_i, \mathbf{y}_i) / (\mathbf{x}_i, \mathbf{y}_i)\}$  ζεύγος με  $\mathbf{x}_i$  πίνακα στήλης εισόδου και  $\mathbf{y}_i$  τον αντίστοιχο επιθυμητό πίνακα στήλης εξόδου,  $i=1, \dots, I$ .

Αν ο δείκτης  $\nu$  αριθμεί τους νευρώνες του επιπέδου εξόδου  $Z$ ,  $\nu=1, \dots, N_Z$  και ο πίνακας εξόδου του ΝΔ είναι  $\mathbf{y}^Z = [y_1^Z, \dots, y_\nu^Z, \dots, y_{N_Z}^Z]^T$  για συγκεκριμένο ζεύγος  $(\mathbf{x}_i, \mathbf{y}_i)$  ορίζουμε μία συνάρτηση  $\Delta(i)$  με χρήση της Ευκλείδειας απόστασης σύμφωνα με την σχέση:

$$\Delta(i) = \frac{1}{2} D_E(\mathbf{y}^Z, \mathbf{y}_i)^2 = \frac{1}{2} \sum_{\nu=1}^{N_Z} (y_\nu^Z - y_{\nu i})^2$$

Η σχέση  $\Delta(i)$  είναι ένα άθροισμα τετραγωνικών σφαλμάτων μεταξύ της παραγόμενης εξόδου  $\mathbf{y}^Z$  του ΝΔ όταν η είσοδος είναι  $\mathbf{x}_i$  και της επιθυμητής εξόδου  $\mathbf{y}_i$ . Μπορούμε να ορίσουμε τώρα την συνάρτηση κόστους  $K(\cdot)$  που θα έχει ανεξάρτητες

μεταβλητές όλα τα  $\mathbf{w}_v^\zeta$  για το συγκεκριμένο σύνολο εκπαίδευσης  $S$  σύμφωνα με την σχέση

$$K(\mathbf{w}_v^\zeta) = \sum_{i=1}^I \Delta(i)$$

Αντιμετωπίζοντας την εύρεση ελάχιστης τιμής της  $K$  ως πρόβλημα βελτιστοποίησης (*optimization*) οι τιμές των  $\mathbf{w}_v^\zeta$  μπορούν να εκτιμηθούν επαναληπτικά με την μέθοδο **καθόδου κατά την κλίση (gradient descent, GD)** σύμφωνα με σχέση

$$\mathbf{w}_v^\zeta(t+1) = \mathbf{w}_v^\zeta(t) - \rho \left. \frac{\partial K}{\partial \mathbf{w}_v^\zeta} \right|_{\mathbf{w}_v^\zeta(t)} = \mathbf{w}_v^\zeta(t) - \rho \sum_{i=1}^I \left. \frac{\partial \Delta(i)}{\partial \mathbf{w}_v^\zeta} \right|_{\mathbf{w}_v^\zeta(t)} \quad (3.5.1)$$

Σύμφωνα με τον κανόνα παραγωγίσης της αλυσίδας

$$\frac{\partial \Delta(i)}{\partial \mathbf{w}_v^\zeta} = \frac{\partial \Delta(i)}{\partial \sigma_v^\zeta} \cdot \frac{\partial \sigma_v^\zeta}{\partial \mathbf{w}_v^\zeta}$$

Ο παράγον  $\frac{\partial \sigma_v^\zeta}{\partial \mathbf{w}_v^\zeta}$  για όλα τα  $\mathbf{w}_v^\zeta$  (για κάθε  $v$  και  $\zeta > 0$ ) είναι

$$\frac{\partial \sigma_v^\zeta}{\partial \mathbf{w}_v^\zeta} = \frac{\partial}{\partial \mathbf{w}_v^\zeta} \left( (\mathbf{y}^{\zeta-1})^T \cdot \mathbf{w}_v^\zeta \right) = \frac{\partial}{\partial \mathbf{w}_v^\zeta} \left( y_1^{\zeta-1} \cdot w_{1v}^\zeta + \dots + y_\mu^{\zeta-1} \cdot w_{\mu v}^\zeta + \dots + y_{N_{\zeta-1}}^{\zeta-1} \cdot w_{N_{\zeta-1}v}^\zeta + w_{0v}^\zeta \right) = (\mathbf{y}^{\zeta-1})^T$$

Απομένει τώρα ο υπολογισμός του πρώτου παράγοντα του γινομένου, τον οποίο ονομάζουμε

$$\delta_v^\zeta(i) = \delta_v^\zeta = \frac{\partial \Delta(i)}{\partial \sigma_v^\zeta} \quad \text{και} \quad \boldsymbol{\delta}^\zeta = \left[ \delta_1^\zeta, \dots, \delta_v^\zeta, \dots, \delta_{N_\zeta}^\zeta \right]^T$$

Θα υπολογίσουμε πρώτα το  $\delta_v^\zeta$  (παραλήφθηκε ο δείκτης  $i$  για απλοποίηση στην γραφή) για έναν νευρώνα  $v$  του επιπέδου εξόδου ( $\zeta=Z$ ,  $N_\zeta=N_Z$ ) και δείκτη αρίθμησης των νευρώνων  $\mu=1, \dots, N_Z$ .

$$\begin{aligned} \delta_v^Z &= \frac{\partial \Delta(i)}{\partial \sigma_v^Z} = \frac{\partial}{\partial \sigma_v^Z} \left( \frac{1}{2} \sum_{\mu=1}^{N_Z} (y_\mu^Z - y_{\mu i})^2 \right) \quad \text{με} \quad y_\mu^Z = f(\sigma_\mu^Z) \Rightarrow \\ \delta_v^Z &= \frac{1}{2} \sum_{\mu=1}^{N_Z} \frac{\partial}{\partial \sigma_v^Z} (f(\sigma_\mu^Z) - y_{\mu i})^2 = \frac{2}{2} (f(\sigma_v^Z) - y_{v i}) f'(\sigma_v^Z) \Rightarrow \\ \delta_v^Z &= (y_v^Z - y_{v i}) f'(\sigma_v^Z) \end{aligned} \quad (3.5.2)$$

Για τα κρυφά επίπεδα ( $0 < \zeta < Z$ ) ο υπολογισμός του  $\delta_v^\zeta$  στο  $\zeta$  επίπεδο είναι περιπλοκότερος. Θα βασισθεί στις τιμές διόρθωσης των  $\boldsymbol{\delta}^{\zeta+1} = \left[ \delta_1^{\zeta+1}, \dots, \delta_{N_{\zeta+1}}^{\zeta+1} \right]^T$  του επομένου επιπέδου του  $\zeta+1$ , αρχίζοντας από το επίπεδο που προηγείται του επιπέδου εξόδου ( $\zeta=Z-1$ ,  $Z=\zeta+1$ ) και οδεύοντας προοδευτικά προς τα πίσω (*back-*

*error propagation*). Συγκεκριμένα αν  $\zeta$  ένα κρυφό επίπεδο, το επόμενο του θα είναι το  $\zeta+1$ . Έστω ακόμη  $\mu$  ένας μετρητής αρίθμησης των νευρώνων του  $\zeta$  επιπέδου και  $\kappa$  ένας μετρητής αρίθμησης των νευρώνων του  $\zeta+1$  επιπέδου. Η συνάρτηση  $\Delta(i)$  εξαρτάται από τα  $\sigma_1^{\zeta+1}, \dots, \sigma_\mu^{\zeta+1}, \dots, \sigma_{N_{\zeta+1}}^{\zeta+1}$  και κάθε  $\sigma_\kappa^{\zeta+1}$  εξαρτάται από το  $\sigma_\nu^\zeta$  του  $\nu$ -οστού νευρώνα του  $\zeta$  επιπέδου. Σύμφωνα με τον κανόνα της αλυσιδωτής παραγώγισης

$$\begin{aligned} \delta_\nu^\zeta &= \frac{\partial \Delta_i}{\partial \sigma_\nu^\zeta} = \frac{\partial \Delta_i}{\partial \sigma_\mu^{\zeta+1}} \frac{\partial \sigma_\mu^{\zeta+1}}{\partial \sigma_\nu^\zeta} \\ \frac{\partial \Delta_i}{\partial \sigma_\mu^{\zeta+1}} &= \left[ \frac{\partial \Delta_i}{\partial \sigma_1^{\zeta+1}}, \dots, \frac{\partial \Delta_i}{\partial \sigma_\mu^{\zeta+1}}, \dots, \frac{\partial \Delta_i}{\partial \sigma_{N_{\zeta+1}}^{\zeta+1}} \right] = [\delta_1^{\zeta+1}, \dots, \delta_\mu^{\zeta+1}, \dots, \delta_{N_{\zeta+1}}^{\zeta+1}] \\ \sigma^{\zeta+1} &= (\mathbf{W}^{\zeta+1})^T \cdot \mathbf{y}^\zeta = \begin{bmatrix} (\mathbf{w}_1^{\zeta+1})^T \cdot \mathbf{y}^\zeta \\ \vdots \\ (\mathbf{w}_\mu^{\zeta+1})^T \cdot \mathbf{y}^\zeta \\ \vdots \\ (\mathbf{w}_{N_{\zeta+1}}^{\zeta+1})^T \cdot \mathbf{y}^\zeta \end{bmatrix} = \\ &= \begin{bmatrix} w_{11}^{\zeta+1} \cdot f(\sigma_1^\zeta) + \dots + w_{\nu 1}^{\zeta+1} \cdot f(\sigma_\nu^\zeta) + \dots + w_{N_{\zeta+1} 1}^{\zeta+1} \cdot f(\sigma_{N_\zeta}^\zeta) + w_{01}^{\zeta+1} \\ \vdots \\ w_{1\mu}^{\zeta+1} \cdot f(\sigma_1^\zeta) + \dots + w_{\nu\mu}^{\zeta+1} \cdot f(\sigma_\nu^\zeta) + \dots + w_{N_{\zeta+1}\mu}^{\zeta+1} \cdot f(\sigma_{N_\zeta}^\zeta) + w_{0\mu}^{\zeta+1} \\ \vdots \\ w_{1,N_{\zeta+1}}^{\zeta+1} \cdot f(\sigma_1^\zeta) + \dots + w_{\nu,N_{\zeta+1}}^{\zeta+1} \cdot f(\sigma_\nu^\zeta) + \dots + w_{N_{\zeta+1},N_{\zeta+1}}^{\zeta+1} \cdot f(\sigma_{N_\zeta}^\zeta) + w_{0,N_{\zeta+1}}^{\zeta+1} \end{bmatrix} \\ \frac{\partial \sigma_\mu^{\zeta+1}}{\partial \sigma_\nu^\zeta} &= \begin{bmatrix} w_{\nu 1}^{\zeta+1} \cdot f'(\sigma_\nu^\zeta) \\ \vdots \\ w_{\nu\mu}^{\zeta+1} \cdot f'(\sigma_\nu^\zeta) \\ \vdots \\ w_{\nu N_{\zeta+1}}^{\zeta+1} \cdot f'(\sigma_\nu^\zeta) \end{bmatrix} \\ \delta_\nu^\zeta &= \frac{\partial \Delta_i}{\partial \sigma_\nu^\zeta} = \frac{\partial \Delta_i}{\partial \sigma_\mu^{\zeta+1}} \frac{\partial \sigma_\mu^{\zeta+1}}{\partial \sigma_\nu^\zeta} = \\ &= [\delta_1^{\zeta+1}, \dots, \delta_\mu^{\zeta+1}, \dots, \delta_{N_{\zeta+1}}^{\zeta+1}] \cdot \begin{bmatrix} w_{\nu 1}^{\zeta+1} \cdot f'(\sigma_\nu^\zeta) \\ \vdots \\ w_{\nu\mu}^{\zeta+1} \cdot f'(\sigma_\nu^\zeta) \\ \vdots \\ w_{\nu N_{\zeta+1}}^{\zeta+1} \cdot f'(\sigma_\nu^\zeta) \end{bmatrix} \Rightarrow \\ \delta_\nu^\zeta &= \left[ \sum_{\mu=1}^{N_{\zeta+1}} \delta_\mu^{\zeta+1} \cdot w_{\nu\mu}^{\zeta+1} \right] \cdot f'(\sigma_\nu^\zeta) \quad (3.5.3) \end{aligned}$$

Τελικά από τις σχέσεις (3.5.1), (3.5.2), (3.5.3) συνεπάγεται ότι

$$\mathbf{w}_\nu^\zeta(t+1) = \mathbf{w}_\nu^\zeta(t) - \rho \sum_{\forall (x_i, y_i)} \delta_\nu^\zeta \cdot \mathbf{y}^{\zeta-1}(i) \quad (3.5.4)$$

όπου  $\delta_\nu^\zeta$  δίνεται από την σχέση (3.5.2)

$$\delta_v^Z = (y_v^Z - y_{v,i})f'(\sigma_v^Z)$$

όταν ο νευρώνας  $v$  βρίσκεται στο επίπεδο εξόδου και από την (3.5.3)

$$\delta_v^Z = \left[ \sum_{\mu=1}^{N_{\zeta+1}} \delta_{\mu}^{\zeta+1} \cdot w_{v\mu}^{\zeta+1} \right] \cdot f'(\sigma_v^Z)$$

όταν ο νευρώνας  $v$  βρίσκεται σε κρυφό επίπεδο ξεκινώντας από το τελευταίο και υποχωρώντας προοδευτικά μέχρι το πρώτο επίπεδο (*back-error propagation*). Ο δείκτης  $i$  του αθροίσματος της (3.5.4) ενυπάρχει στον ορισμό του  $\delta_v^Z = \frac{\partial \Delta(i)}{\partial \sigma_v^Z}$ .

Το ΝΔ Back error propagation είναι από τα πλέον χρησιμοποιούμενα ΝΔ και έχει εφαρμοσθεί σε πληθώρα εφαρμογών από διαφορετικές επιστημονικές περιοχές. Στα κρυφά επίπεδά του προσδιορίζονται ουσιαστικά αυτόματα τα αποτελεσματικά χαρακτηριστικά για την ταξινόμηση. Το πλήθος των κλάσεων μπορεί να είναι το ίδιο με αυτό των ανυσμάτων εισόδου του συνόλου εκπαίδευσης και έτσι το ΝΔ να προσεγγίζει ένα μετασχηματισμό των ανυσμάτων εισόδου στα ανύσματα εξόδου. Το σημαντικότερο μειονέκτημα του back propagation είναι ο χρόνος ολοκλήρωσης της εκπαίδευσης του ή χρόνος σύγκλισης, όπως λέγεται αλλιώς. Είναι δυνατόν να χρειασθούν εκατοντάδες χιλιάδες επαναλήψεις εωσότου συγκλίνει ακόμη και για σχετικά απλές εφαρμογές. Σε κάποιες εφαρμογές χρειάστηκαν μερικές μέρες για την σύγκλιση του συστήματος. Ο εγκλωβισμός της διαδικασίας σύγκλισης σε τοπικά ελάχιστα της συνάρτησης κόστους είναι ένα επιπρόσθετο πρόβλημα που θα αναλύσουμε σε ακόλουθη παράγραφο.

Παραλλαγή της παραπάνω διαδικασίας είναι η διόρθωση των βαρών για κάθε πρότυπο  $\mathbf{x}_i$  του συνόλου εκπαίδευσης με τυχαία επιλογή. Σ' αυτήν την περίπτωση ονομάζουμε την μέθοδο εύρεσης ελαχίστου της συνάρτησης **στοχαστική κάθοδο κατά την κλίση (stochastic gradient descent, SGD)**. Οι μέθοδοι αυτοί είναι αργές και απαιτούν πολλές επαναλήψεις. Μία βελτίωση τους εύκολα υλοποιήσιμη είναι η πρόσθεση ενός ακόμη όρου στην σχέση διόρθωσης που ονομάζουμε όρο **ορμής** ή **ροπής (momentum term)**. Συγκεκριμένα η κατάβαση κατά την κλίση για μια συνάρτηση κόστους  $K(\mathbf{w}_v^Z)$  δίνεται σύμφωνα με όσα προαναφέρθηκαν από τις σχέσεις:

$$\begin{aligned} \mathbf{w}_v^Z(t+1) &= \mathbf{w}_v^Z(t) + \Delta \mathbf{w}_v^Z(t+1) \\ \Delta \mathbf{w}_v^Z(t+1) &= -\rho \left. \frac{\partial K}{\partial \mathbf{w}_v^Z} \right|_{\mathbf{w}_v^Z(t)} = -\rho \nabla K_{\mathbf{w}_v^Z(t)} \end{aligned}$$

Στην διόρθωση με ορμή

$$\Delta \mathbf{w}_v^Z(t+1) = -\rho \left. \frac{\partial K}{\partial \mathbf{w}_v^Z} \right|_{\mathbf{w}_v^Z(t)} + \alpha \Delta \mathbf{w}_v^Z(t)$$

$$\Delta \mathbf{w}_v^Z(t) = \Delta \mathbf{w}_v^Z(t) - \Delta \mathbf{w}_v^Z(t-1)$$

Και  $\alpha \in (0, 1)$ , με συνήθεις τιμές από 0.2 έως 0.8.

Προστίθεται δηλαδή και ένα μέρος της προηγούμενης μεταβολής των βαρών του νευρώνα. Η ποσότητα  $\alpha \Delta \mathbf{w}_v^\zeta(t)$  είναι ο όρος της ορμής (*momentum term*). Η χρήση του μειώνει το πλήθος των επαναλήψεων αυξάνοντας το βήμα διόρθωσης των βαρών κατά την κάθοδο όταν η συνάρτηση κόστους μεταβάλλεται αργά, δηλαδή πλάτος της βάρθρωσης ( $|\nabla K_{\mathbf{w}_v^\zeta(t)}|$ ) είναι πολύ μικρό. Επίσης εξομαλύνει τις απότομες αλλαγές της διεύθυνσης της διόρθωσης. Αυτές οι ιδιότητες της ορμής φαίνονται αλγεβρικά αν υπολογίσουμε της διόρθωση μετά από  $T$  επαναλήψεις από μία αρχική τιμή  $\mathbf{w}(0)$  των βαρών κάποιου νευρώνα (τα σύμβολα  $v, \zeta$  έχουν παραληφθεί για απλότητα). Θέτοντας  $\mathbf{G}(t) = -\rho \nabla K_{\mathbf{w}_v^\zeta(t)}$  θα ισχύει:

$$\Delta \mathbf{w}(T) = -\rho \mathbf{G}(T) + \alpha \cdot \Delta \mathbf{w}(T-1)$$

$$\Delta \mathbf{w}(T-1) = -\rho \mathbf{G}(T-1) + \alpha \cdot \Delta \mathbf{w}(T-2) \Rightarrow \alpha \Delta \mathbf{w}(T-1) = -\rho \alpha \mathbf{G}(T-1) + \alpha^2 \cdot \Delta \mathbf{w}(T-2)$$

$$\Delta \mathbf{w}(T-2) = -\rho \mathbf{G}(T-2) + \alpha \cdot \Delta \mathbf{w}(T-3) \Rightarrow \alpha^2 \Delta \mathbf{w}(T-2) = -\rho \alpha^2 \mathbf{G}(T-2) + \alpha^3 \cdot \Delta \mathbf{w}(T-3)$$

⋮

$$\Delta \mathbf{w}(2) = -\rho \mathbf{G}(2) + \alpha \cdot \Delta \mathbf{w}(1) \Rightarrow \alpha^{T-2} \Delta \mathbf{w}(2) = -\rho \alpha^{T-2} \mathbf{G}(2) + \alpha^{T-1} \cdot \Delta \mathbf{w}(1)$$

$$\Delta \mathbf{w}(1) = -\rho \mathbf{G}(1) + \alpha \cdot \Delta \mathbf{w}(0) \Rightarrow \alpha^{T-1} \Delta \mathbf{w}(1) = -\rho \alpha^{T-1} \mathbf{G}(1) + \alpha^T \cdot \Delta \mathbf{w}(0)$$

Προσθέτοντας τις εξισώσεις μετά της συνεπαγωγές προκύπτει:

$$\Delta \mathbf{w}(T) = -\rho \sum_{t=0}^{T-1} \alpha^t \mathbf{G}(T-t) + \alpha^T \cdot \Delta \mathbf{w}(0)$$

Μετά από μερικές επαναλήψεις, επειδή  $0 < \alpha < 1$ , η μεταβολή  $\Delta \mathbf{w}(T)$  επηρεάζεται από τις προηγούμενες τιμές της  $\mathbf{G}$  που φθίνουν κατά  $\alpha^t$  και ο όρος  $\alpha^T \cdot \Delta \mathbf{w}(0)$  γίνεται αμελητέος. Αν η  $\mathbf{G}(\cdot)$  στο διάστημα αυτό των επαναλήψεων έχει σταθερή τιμή  $\mathbf{G}$  τότε

$$\Delta \mathbf{w}(T) = -\rho(1 + \alpha + \alpha^2 + \dots + \alpha^{T-1})\mathbf{G} = -\rho \frac{1 - \alpha^T}{1 - \alpha} \mathbf{G} \approx -\frac{\rho}{1 - \alpha} \mathbf{G}$$

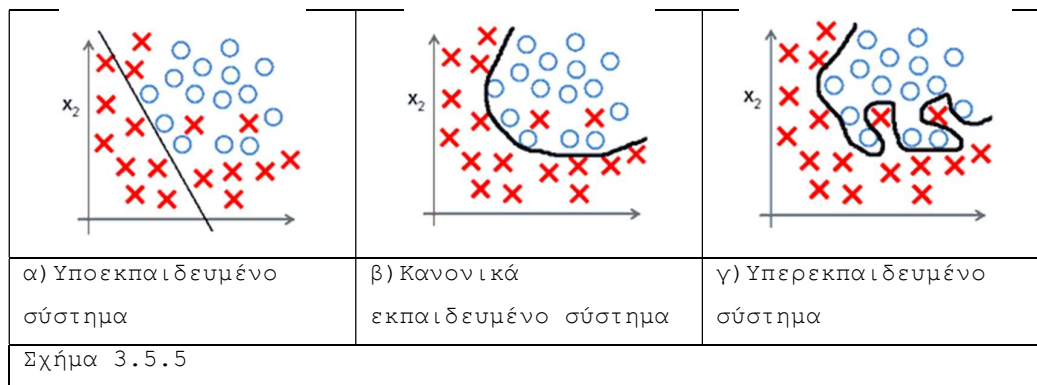
Για  $\rho=0.2$  και  $\alpha=0.5$  η αλλαγή των τιμών των βαρών θα είναι  $0.4\mathbf{G}$  ενώ χωρίς τον όρο ορμής θα ήταν  $0.2\mathbf{G}$ . Η μέθοδος κατάβασης της κλίσης με όρο ορμής (Gradient Descent with Momentum, **GDM**), εφαρμόζεται με διόρθωση ανά πρότυπο με τυχαία επιλογή και τιτλοφορείται ως Stochastic Gradient Descent with Momentum, **SGDM**.

Το πλήθος των επιπέδων και των νευρώνων των είναι ένα ζήτημα που μας απασχολεί. Δεν υπάρχει ένας γενικός τύπος-αλγόριθμος που να μας οδηγεί σε ένα βέλτιστο πλήθος παραμέτρων-βαρών. Ένας μεγάλος αριθμός βαρών αυξάνει το υπολογιστικό κόστος και μπορεί να οδηγήσει σε λεπτομερή εκμάθηση του συνόλου εκπαίδευσης ακόμη και του θορύβου των μετρήσεων, εις βάρος της γενικότητας του, δηλαδή της επιτυχούς ταξινόμησης νέων αγνώστων προτύπων. Αυτό το φαινόμενο το ονομάζουμε *υπερεκπαίδευση*. Μικρός αριθμός παραμέτρων μπορεί να οδηγήσει σε αδυναμία εκπαίδευσης ή χαμηλή απόδοση-ακρίβεια ταξινόμησης. Την



περίπτωση αυτή αναφέρουμε ως *υποεκπαίδευση*. Στο Σχήμα 3.5.5 φαίνεται η διαχωριστική καμπύλη τριών συστημάτων ενός προβλήματος δύο διαστάσεων για τις περιπτώσεις που αναφέραμε.

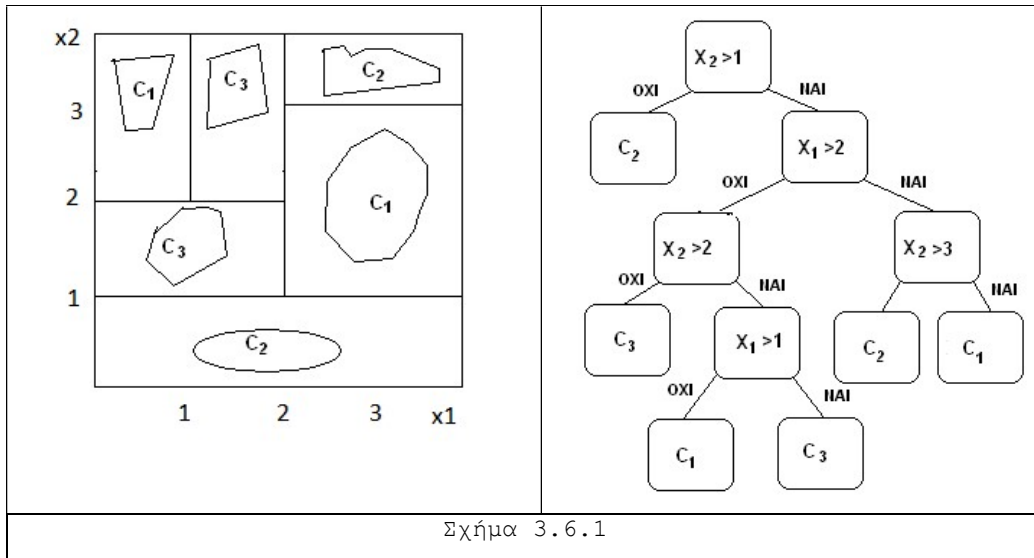
Μία απλή μέθοδος είναι να ξεκινάμε με ένα μικρό πλήθος παραμέτρων και να το αυξάνουμε προοδευτικά μέχρι να έχουμε ικανοποιητικό αποτέλεσμα. Εκείνο που συχνά εμφανίζεται όταν έχουμε υψηλή απόδοση εκπαίδευσης είναι η υπερεκπαίδευση που διαπιστώνεται όταν το σύστημά μας ταξινομεί πρότυπα για τα οποία δεν εκπαιδεύτηκε και παρουσιάζει χαμηλή απόδοση. Για να εντοπίσουμε νωρίς το πρόβλημα χωρίζουμε το αρχικό σύνολο  $S$  των διαθέσιμων δεδομένων σε δύο μέρη. Το ένα χρησιμοποιείται ως σύνολο εκπαίδευσης (Training set,  $T$ ) και το άλλο ως σύνολο επικύρωσης του αποτελέσματος (Validation set,  $V$ ). Θα πρέπει το ποσοστό επιτυχίας επιτυχούς ταξινόμησης (accuracy) να είναι υψηλό με παραπλήσιες τιμές.



### 3.6 Δένδρα απόφασης

Τα δένδρα απόφασης (decision trees) αποτελούν μια ευρεία κατηγορία τεχνικών μη γραμμικών ταξινομητών στην οποία ο χώρος των χαρακτηριστικών διαιρείται σε περιοχές που αντιστοιχούν στις επιθυμητές κλάσεις. Η απόφαση της ταξινόμησης ενός προτύπου προκύπτει από τις απαντήσεις σε ερωτήματα που υποβάλλονται σύμφωνα με μία δενδρική δομή. Θα παρουσιάσουμε ακολούθως μια αντιπροσωπευτική και δημοφιλή τεχνική των δένδρων απόφασης κατά την οποία ο χώρος χωρίζεται σε υπερ-παραλληλόγραμμα. Για ένα υπό ταξινόμηση πρότυπο κάθε ερώτημα αφορά ένα χαρακτηριστικό  $x_n$  και είναι της μορφής «ισχύει  $x_n < t$ » όπου  $t$  κατάλληλη τιμή κατωφλίου. Τα δένδρα απόφασης αυτής της μορφής ονομάζονται συνήθη δυαδικά δένδρα ταξινόμησης - ΣΔΔΤ (ordinary binary classification trees - OBCT). Η λειτουργία ενός ΣΔΔΤ γίνεται εύκολα αντιληπτή με ένα παράδειγμα δύο χαρακτηριστικών και κλάσεις κατανεμημένες όπως στο

Σχ. [ 3.6.1]. Με απλή γεωμετρική παρατήρηση μπορούμε να οδηγηθούμε στο δένδρο απόφασης του που φαίνεται στο ίδιο σχήμα.



Ο οπτικός διαχωρισμός βέβαια δεν είναι δυνατός σε χώρους υψηλότερης του τρία διάστασης και σε κάθε περίπτωση η εκτέλεση από Υπολογιστές απαιτεί την αλγοριθμική και μαθηματικολογική διατύπωση. Παρατηρώντας τον χώρο των προτύπων του παραδείγματος και τις ορθογώνιες περιοχές που περικλείουν τις κλάσεις, σε σχέση με το αντίστοιχο δένδρο απόφασης διαπιστώνουμε ότι σε κάθε κόμβο  $t$  τίθεται ως ερώτημα μια συνθήκη η ισχύς της οποίας χωρίζει ένα υποσύνολο  $X_t$  του χώρου των προτύπων σε δύο μέρη. Ο ριζικός κόμβος αφορά όλο σύνολο εκπαίδευσης. Η απάντηση του ερωτήματος μπορεί να είναι θετική (Ναι) ή αρνητική (Όχι) και από αυτή καθορίζονται τα δύο υποσύνολα  $X_{tN}$ ,  $X_{tO}$  που ικανοποιούν τις σχέσεις

$$X_{tN} \cup X_{tO} = X_t \text{ και } X_{tN} \cap X_{tO} = \emptyset$$

Πρέπει ακόμη να καθορισθούν κριτήρια ώστε:

- ο διαχωρισμός να είναι βέλτιστος
- να τερματίζεται ο διαχωρισμός ενός κόμβου
- να αντιστοιχίζεται σε ένα κόμβο-φύλλο μία κλάση

Μια ερώτηση κόμβου είναι της μορφής  $x_n > \alpha$ , όπου  $x_n$  η τιμή του  $n$  χαρακτηριστικού ενός προτύπου και  $\alpha$  μια τιμή κατωφλίου. Το  $\alpha$  μπορεί να πάρει απεριόριστες τιμές στο πεδίο μεταβολής του, εν τούτοις αν  $L_t$  είναι το πλήθος των προτύπων που ανήκουν στον  $X_t$ , αρκεί ένα πεπερασμένο πλήθος τιμών  $\alpha_{νλ}$ ,

$\lambda=1,2,\dots, \Lambda_t$  , για να διαχωριστεί ο χώρος. Η τιμή  $\alpha_{\lambda}$  θα επιλεγεί ώστε να ικανοποιείται το κριτήριο του βέλτιστου διαχωρισμού.

#### Καθορισμός του κριτηρίου διαχωρισμού

Ο διαχωρισμός του  $X_t$  πρέπει να είναι τέτοιος ώστε σε ένα τουλάχιστον από τα δύο μέρη του να κυριαρχούν πληθυσμιακά τα πρότυπα μιας μόνο κλάσης, να είναι δηλαδή «καθαρό» σύμφωνα με την ορολογία των δυαδικών δένδρων απόφασης. Ως μέτρο της καθαρότητας ενός κόμβου μπορεί να χρησιμοποιηθεί η εντροπία (μέση πληροφορία) γνωστή από την θεωρία της πληροφορίας. Για έναν κόμβο  $t$  η εντροπία  $I(t)$  δίνεται από την σχέση

$$I(t) = - \sum_{i=1}^M P(C_i|t) \cdot \log_2 P(C_i|t)$$

Όπου  $P(C_i|t)$  η πιθανότητα ένα πρότυπο του χώρου  $X(t)$  να ανήκει στην κλάση  $C_i$ . Υπενθυμίζεται ότι  $0 \cdot \log_2 0 = 0$  και ότι  $I(t)$  μεγιστοποιείται αν οι τιμές  $P(C_i|t)$  είναι ίσες μεταξύ τους. Οι πιθανότητες  $P(C_i|t)$  εκτιμώνται με βάση τα ποσοστά  $\Lambda_t^i / \Lambda_t$  όπου  $\Lambda_t$  το πλήθος των προτύπων του κόμβου και  $\Lambda_t^i$  το πλήθος όσων εξ αυτών ανήκουν στην κλάση  $C_i$ . Με τον διαχωρισμό του κόμβου η καθαρότητά του μεταβάλλεται κατά  $\Delta I(t)$  σύμφωνα με τη σχέση

$$\Delta I(t) = I(t) - \frac{\Lambda_{tN}}{\Lambda_t} \cdot I(t_N) - \frac{\Lambda_{tO}}{\Lambda_t} \cdot I(t_O)$$

Η ποσότητα  $\Delta I(t)$  υπολογίζεται για όλα τα χαρακτηριστικά και τις κατάλληλες τιμές κατώφλιου των προτύπων που ανήκουν στους χώρους  $X_t, X_{tN}, X_{tO}$ . Ακολουθώς επιλέγεται το χαρακτηριστικό και το κατώφλι που μεγιστοποιούν την  $\Delta I(t)$  και προκύπτει το ερώτημα του κόμβου  $t$ .

#### Καθορισμός του κριτηρίου τερματισμού του διαχωρισμού

Ο διαχωρισμός του χώρου  $X_t$ , είναι δυνατόν να τερματίσει όταν η μέγιστη τιμή  $\Delta I(t)$  είναι μικρότερη από ένα προκαθορισμένο όριο. Εναλλακτικά μπορεί να τερματιστεί αν ο πληθυσμός  $\Lambda_t$  είναι αρκετά μικρός ή όταν τα πρότυπα ανήκουν σε μία κλάση ( $\Delta I(t)=0$ ).

#### Καθορισμός κριτηρίου αντιστοίχισης μιας κλάσης στον κόμβο-φύλλο

Σύμφωνα με τα παραπάνω μετά τον τερματισμό στο κόμβο  $t$  αποδίδεται ως όνομα το όνομα της κλάσης των προτύπων που υπερτερούν πληθυσμιακά στον χώρο  $X_t$ .

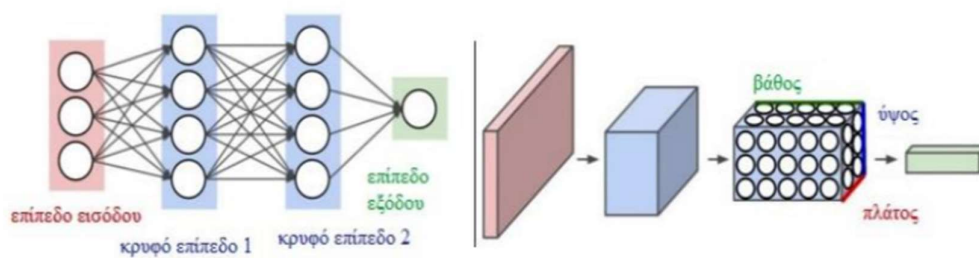
### **3.7 Συνελικτικά Νευρωνικά Δίκτυα**

Τα Συνελικτικά Νευρωνικά Δίκτυα (ΣΝΔ) (CNN:Convolutional Neural Networks) έχουν χρησιμοποιηθεί ευρέως στην αναγνώριση εικόνων και μοιάζουν με τα Τεχνητά Νευρωνικά Δίκτυα πρόσθιας τροφοδότησης (feed-forward networks, FFN).

Υπάρχουν διαφορές ανάμεσα σε αυτά τα δύο είδη δικτύων, που καθιστούν τα ΣΝΔ πιο ελκυστικά για χρήση. Πιο συγκεκριμένα στα κλασικά FFN η κλιμάκωση μεγάλων εικόνων οδηγεί σε μεγάλο μη διαχειρίσιμο πλήθος βαρών. Για παράδειγμα στο σύνολο δεδομένων CIFAR-10 οι εικόνες είναι μεγέθους, μόνο  $32 \times 32 \times 3$  (πλάτος  $\times$  ύψος  $\times$  χρωματικά κανάλια), και επομένως ένας πλήρως συνδεδεμένος νευρώνας στο πρώτο κρυφό επίπεδο ενός ΤΝΔ θα είχε  $32 \times 32 \times 3 = 3072$  βάρη. Παρά το γεγονός ότι αυτός ο αριθμός δείχνει διαχειρίσιμος, για είσοδο εικόνας μεγαλύτερων διαστάσεων, π.χ  $200 \times 200 \times 3$  θα είχαμε νευρώνες με  $200 \times 200 \times 3 = 120.000$  βάρη ο καθένας.

Για περισσότερους νευρώνες ο αριθμός των παραμέτρων θα μεγάλωνε ραγδαία. Συνεπώς η πλήρης συνδεσιμότητα είναι σπάταλη, και μπορεί λόγω των πολλών παραμέτρων να οδηγήσει εύκολα σε *υπερπροσαρμογή (overfitting)* του δικτύου. Από την άλλη πλευρά, η χρήση των ΣΝΔ εκμεταλλεύεται το γεγονός ότι η είσοδος αποτελείται από εικόνες, και περιορίζει την αρχιτεκτονική του με πιο έξυπνο τρόπο. Γενικότερα τα ΣΝΔ, διαθέτουν νευρώνες οι οποίοι έχουν 3 διαστάσεις (πλάτος-ύψος-βάθος), και έχουν την ιδιαιτερότητα να συνδέονται διαδοχικά με μία μικρή περιοχή του προηγούμενου επιπέδου αντίθετα με ότι γινόταν σε μία πλήρη σύνδεση. Κατ' ουσίαν ο νευρώνας είναι ένα φίλτρο που *συνελίσσεται* στα δεδομένα εισόδου του. Θα πρέπει να διευκρινίσουμε εδώ ότι λέγοντας «συνελίσσεται» εννοούμε την εφαρμογή της δισδιάστατης συνέλιξης με κάποια διευρυμένη χρήση. Οι όροι νευρώνας, κόμβος (node), φίλτρο και πυρήνας (kernel) χρησιμοποιούνται εναλλακτικά στην βιβλιογραφία των συνελικτικών δικτύων για το ίδιο πράγμα.

Στο σχήμα 3.7.1 απεικονίζεται αριστερά μία αρχιτεκτονική ενός κλασικού νευρωνικού δικτύου και δεξιά μία αρχιτεκτονική ενός συνελικτικού νευρωνικού δικτύου. Κάθε επίπεδο ενός ΣΝΔ μετασχηματίζει τον τρισδιάστατο όγκο εισόδου (input volume), σε έναν τρισδιάστατο όγκο εξόδου (output volume). Στο συγκεκριμένο παράδειγμα η έγχρωμη εικόνα του παραπάνω σχήματος αποτελεί το επίπεδο εισόδου επομένως θα έχει βάθος 3, όσα και τα κανάλια μιας RGB εικόνας.



Σχήμα 3.7.1. (Αρχιτεκτονική CNN)

Στα ΣΝΔ εκτελούνται τέσσερις βασικές λειτουργίες:

1. Συνέλιξη (Convolution)
2. Εφαρμογή μη γραμμικότητας (Non Linearity) – ReLU
3. Συγκέντρωση ή Υπο-Δειγματοληψία (Pooling / Sub sampling)
4. Κατηγοριοποίηση με πλήρως συνδεδεμένο επίπεδο (Fully Connected Layer for Classification)

## Επίπεδο Συνέλιξης (Convolution)

Θεμελιώδες δομικό επίπεδο ενός Συνελικτικού Νευρωνικού Δικτύου, όπως έχει προαναφερθεί, είναι το επίπεδο στο οποίο συνελίσεται η είσοδος με  $N$  το πλήθος φίλτρα (νευρώνες βαρών). Το αποτέλεσμα είναι πίνακες που αποτελούν χάρτες χαρακτηριστικών. Το επίπεδο αυτό ονομάζεται συνελικτικό επίπεδο. Τα φίλτρα αυτά στην περίπτωση των έγχρωμων εικόνων είναι τρισδιάστατα, καθώς το μέγεθός τους καθορίζεται από το ύψος  $R$  (πλήθος γραμμών), το πλάτος (πλήθος στηλών)  $C$  και το βάθος  $D$  (πλήθος πινάκων  $R \times C$ ). Στην περίπτωση που η είσοδος είναι μία έγχρωμη εικόνα το βάθος είναι τρία (3) για τα τρία βασικά χρώματα (Red Green Blue: RGB). Για την παραγωγή αυτών των χαρακτηριστικών, σαρώνεται ολόκληρη η εικόνα, πραγματοποιούνται πράξεις εσωτερικού γινομένου μεταξύ των τιμών του φίλτρου και της υποκειμένης περιοχής του πίνακα και στο τέλος εξάγεται το αποτέλεσμα το οποίο τοποθετείται στον χάρτη χαρακτηριστικών.

Η πράξη της *δισδιάστατης συνέλιξης* για δύο διακριτά σήματα  $x(r,c)$ ,  $w(r,c)$ , (δισδιάστατες ακολουθίες) ορίζεται τυπικά από την σχέση:

$$x * w = \{x * w\}(r,c) = \sum_k \sum_l x(k,l) \cdot w(r-k,c-l) \text{ με } r,c,k,l \in \mathbb{Z}$$

Η πράξη της *δισδιάστατης συσχέτισης* για δύο διακριτά σήματα  $x(m,n)$ ,  $y(m,n)$ , (δισδιάστατες ακολουθίες) ορίζεται τυπικά από την σχέση:

$$x \circ w = \{x \circ w\}(r,c) = \sum_k \sum_l x(k,l) \cdot w(k-r,l-c) \text{ με } r,c,k,l \in \mathbb{Z}$$

Αν  $w(k,l) = w(-k,-l)$  οι δύο πράξεις δίνουν το ίδιο αποτέλεσμα. Για ακολουθίες πεπερασμένου μήκους τα όρια των δεικτών  $k,l$  καθορίζονται από τις θέσεις που οι ακολουθίες έχουν όλες τις τιμές τους μηδενικές.

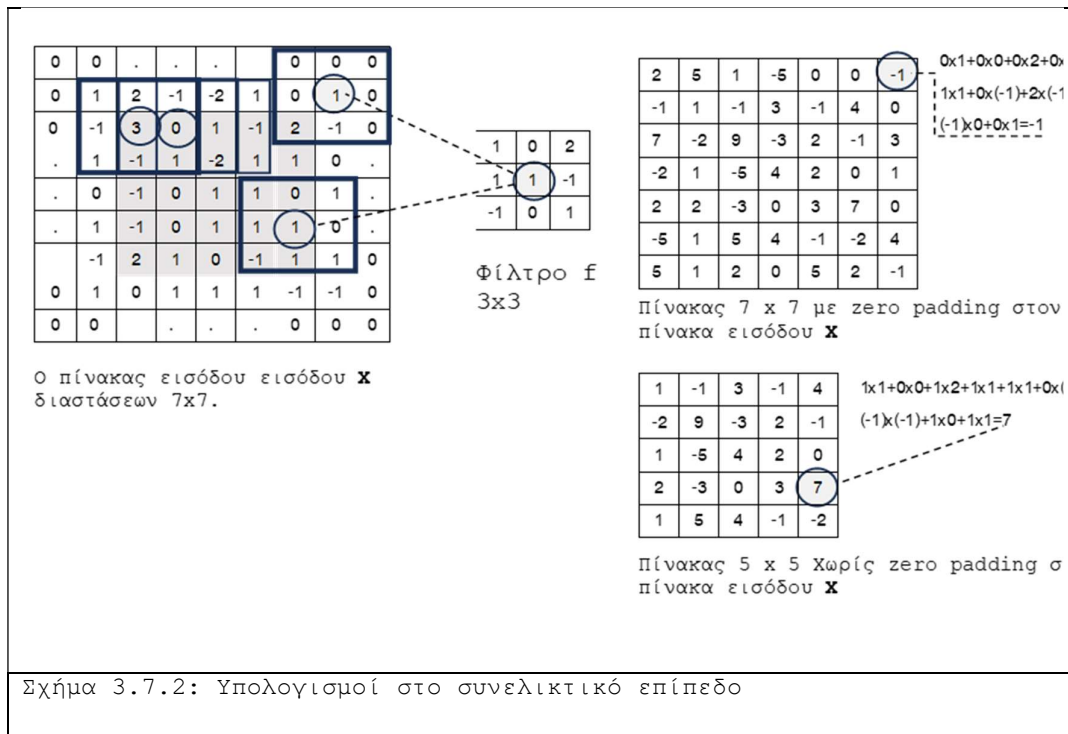
Στα συνελικτικά δίκτυα εφαρμόζεται η πράξη της συσχέτισης. Εν τούτοις ιστορικά επικράτησε να χρησιμοποιείται ο όρος συνέλιξη από την ευρεία χρήση

του στην επεξεργασία σημάτων. Θα διατηρήσουμε αυτήν την ορολογία στην συνέχεια.

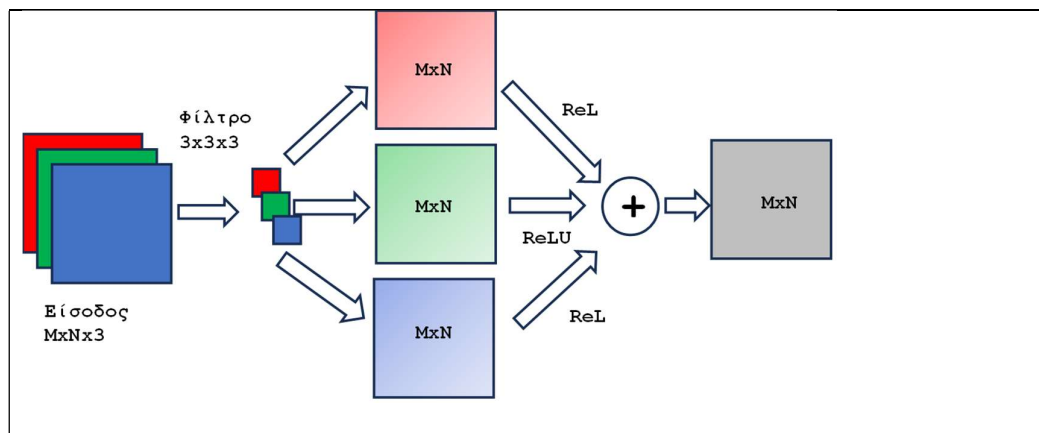
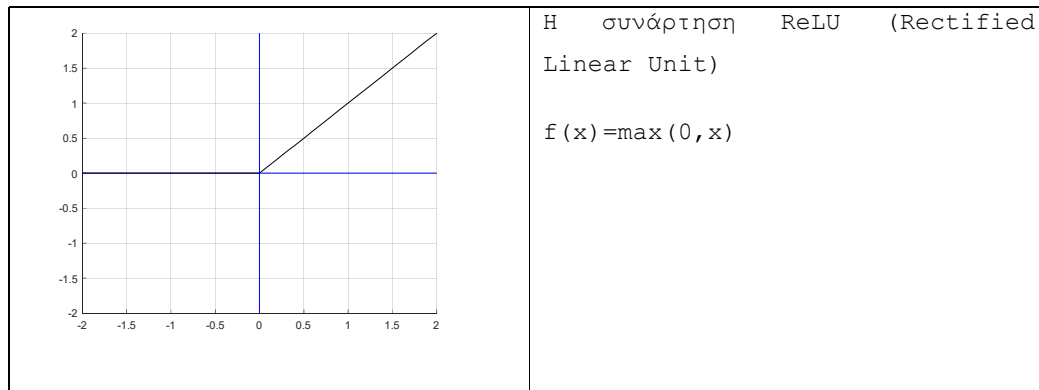
Ανεξαρτήτως των ορίων μέσα στα οποία ορίζεται η πράξη της συνέλιξης στο  $\mathbb{Z} \times \mathbb{Z}$ , εκείνο που μας ενδιαφέρει είναι οι τιμές των ακολουθιών και τις θεωρούμε ως στοιχεία πινάκων α) της εισόδου  $\mathbf{X}$  με διαστάσεις  $(R_x \times C_x)$  και β) του φίλτρου  $\mathbf{W}$  με διαστάσεις  $(R_w \times C_w)$ .

Στο ακόλουθο παράδειγμα (Σχ.3.7.2) φαίνεται η εφαρμογή της πράξης της συνέλιξης (συσχέτισης). Οι τιμές του του φίλτρου  $w(\cdot, \cdot)$  ολισθαίνουν επάνω στις τιμές των γραμμών και των στηλών της  $x(\cdot, \cdot)$  και υπολογίζεται το άθροισμα των γινομένων των ομοίωθων τιμών. Το αποτέλεσμα θα είναι μια ακολουθία πεπερασμένου μήκους με διαστάσεις  $(R_x + R_w - 1) \times (C_x + C_w - 1)$ . Αν θέλουμε οι διαστάσεις της προκύπτουσας ακολουθίας να είναι ίσες με αυτές τις εισόδου, αφαιρούμε  $R_w - 1$  γραμμές και  $C_w - 1$  στήλες από τα όριά της. Απλούστερα μπορούμε να ολισθαίνουμε το φίλτρο έτσι ώστε η κεντρική τιμή του φίλτρου  $w$  να φτάνει στα όρια της  $x$  και αποδίδοντας μηδενικές τιμές στη  $x$  εκτός των ορίων της (zero padding).

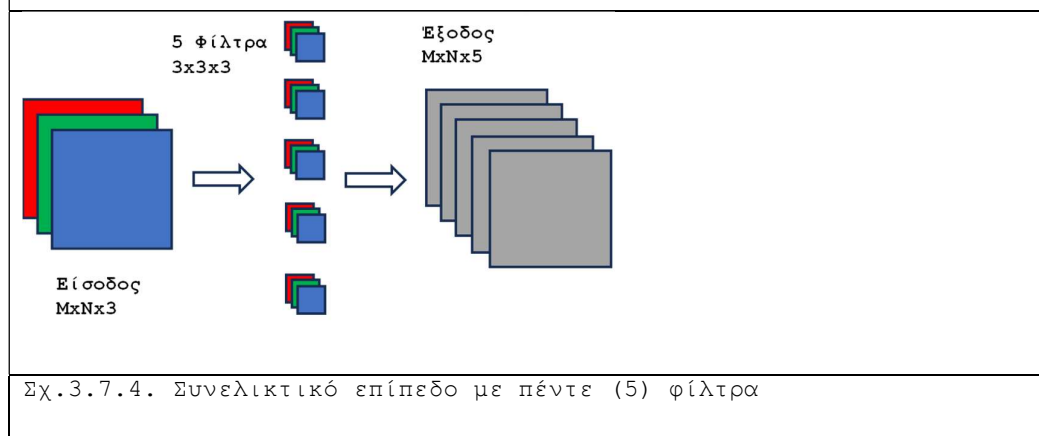
Είναι βολικό το φίλτρο να έχει περιττό αριθμό γραμμών και στηλών π.χ.  $3 \times 3$ ,  $5 \times 5$  ώστε να υπάρχει κεντρική τιμή. Στο τέλος της πράξης της συνέλιξης σε κάθε τιμή που θα προκύψει προσθέτουμε και ένα σταθερό όρο  $w_0$  (bias) η σημασία του οποίου είναι ίδια με αυτήν που αναφέρθηκε στον γραμμικό ταξινομητή.



Η τιμή κάθε στοιχείου του παραγόμενου γίνεται όρισμα μιας συνάρτησης ενεργοποίησης κυρίως της ReLU.



Σχ.3.7.3. Η εφαρμογή ενός φίλτρου σε είσοδο τριών καναλιών στο συνελικτικό επίπεδο



Σχ.3.7.4. Συνελικτικό επίπεδο με πέντε (5) φίλτρα

Αν η είσοδος έχει βάθος  $D$  δηλαδή αποτελείται από  $D$  το πλήθος πίνακες, ονομάζονται και κανάλια (channels), (π.χ.  $D=3$ , για είσοδο μία RGB έγχρωμη

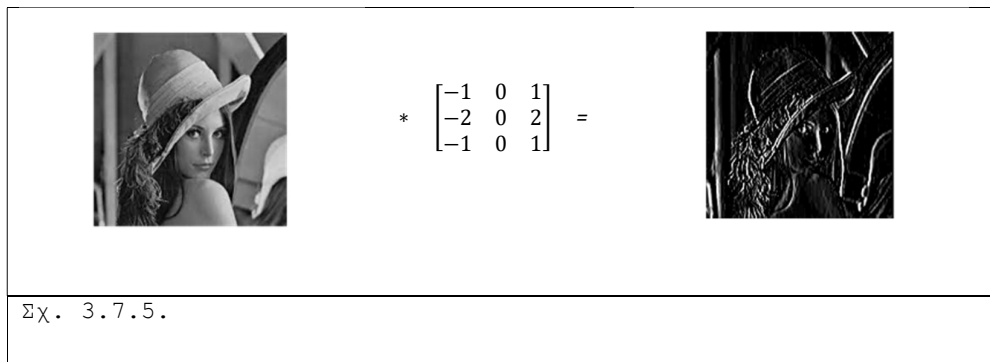
εικόνα), τότε και κάθε φίλτρο πρέπει να έχει βάθος  $D$ , δηλαδή να αποτελείται από  $D$  πίνακες. Για παράδειγμα αν η είσοδος είναι  $R_x \times C_x \times 3$  το φίλτρο πρέπει να είναι  $R_w \times C_w \times 3$ . Στην περίπτωση αυτή μετά την εφαρμογή του φίλτρου θα προκύψουν και  $D$  πίνακες ως αποτελέσματα. Οι πίνακες αυτοί είναι ίσων διαστάσεων και προστίθενται.

Τελικά μετά την άθροιση προκύπτει ένας πίνακας για κάθε φίλτρο. Αν εφαρμοσθούν  $N$  φίλτρα θα προκύψουν  $N$  πίνακες, όσοι και τα φίλτρα. Στην ορολογία των ΣΝΔ συχνά ένας τρισδιάστατος πίνακας καλείται και όγκος (*volume*). Άρα η είσοδος μας είναι ένας όγκος  $R_x \times C_x \times 3$  και με την εφαρμογή  $N$  φίλτρων παράγεται ως αποτέλεσμα ένας όγκος  $R_x \times C_x \times N$ . Τα παραπάνω απεικονίζονται στα ακόλουθα σχήματα (Σχ.3.7.3, Σχ.3.7.4)

Το φίλτρο κατά την ολίσθηση του είναι αποδεκτό να κινείται και με βήμα (δρασκειλιά, *stride*) μεγαλύτερο του ένα. Αυτό είναι μια παραλλαγή συνέλιξης που οδηγεί σε αποτελέσματα μικρότερων διαστάσεων και ακρίβειας. Έχει χρήση στην περίπτωση μεγάλων εικόνων και φίλτρων (π.χ για φίλτρο  $15 \times 15$ , *stride*=2).

Τα παραπάνω, δηλαδή η συνέλιξη της εισόδου με τα φίλτρα και η εφαρμογή της μη γραμμικότητας αποτελούν ένα *Επίπεδο Συνέλιξης (Convolutional Layer)*.

Οι πίνακες που εξάγονται από το επίπεδο συνέλιξης καλούνται **χάρτες χαρακτηριστικών**. Στο ακόλουθο Σχήμα 3.7.5 φαίνεται η εφαρμογή ενός φίλτρου *Sobel* στην εικόνα «lena».

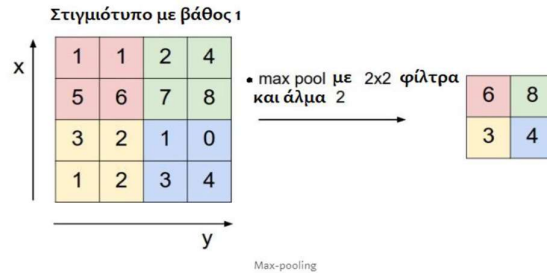


### Επίπεδο δειγματοληψίας (Pooling Layer)

Τα επίπεδα δειγματοληψίας συναντώνται ανάμεσα στα επίπεδα συνέλιξης των ΣΝΔ. Αυτό που κάνει αυτό το επίπεδο είναι να μειώνει το μέγεθος των πινάκων εισόδων και κατά συνέπεια και τους υπολογισμούς ενός δικτύου και θέτοντας υπό έλεγχο προβλήματα υπερπροσαρμογής (*over-fitting*). Λαμβάνοντας τις τιμές μιας μικρής περιοχής π.χ.  $(M \times M)$  ενός πίνακα εισόδου και με βήμα  $M$  παίρνουμε

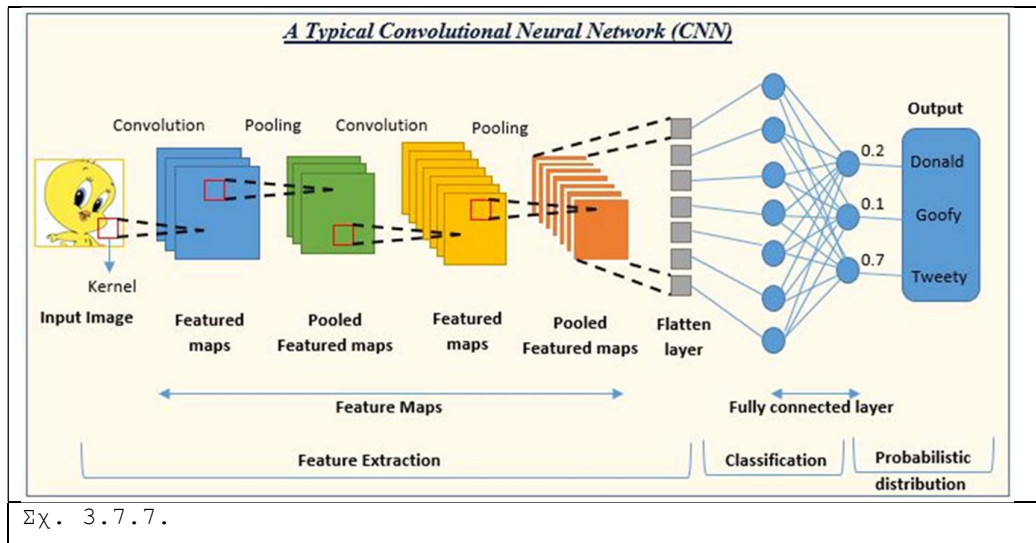


την μεγαλύτερη από αυτές (Max Pooling) ή τον μέσο όρο τους (Average Pooling) και δημιουργούμε έναν νέο πίνακα με μικρότερες διαστάσεις. Αν για παράδειγμα ο πίνακας είναι  $224 \times 224$  για  $M=2$  θα έχουμε ως αποτέλεσμα ένα πίνακα  $112 \times 112$ . Ένα παράδειγμα φαίνεται στο παρακάτω σχήμα 3.7.6.



Σχήμα 3.7.6. (Επίπεδο Υποδειγματοληψίας)

Επίπεδα συνέλιξης και υποδειγματοληψίας μπορούν να τοποθετούνται διαδοχικά καταλήγοντας σε πίνακες μικρότερων διαστάσεων κάθε φορά. Μετά από αυτήν την διαδοχή οι τιμές των τελικών πινάκων τοποθετούνται σε ένα άνυσμα στήλης (*flattening*). Το άνυσμα αυτό θα είναι η είσοδος ενός τελικού επιπέδου ταξινόμησης με πλήρη διασύνδεση, Σχήμα 3.7.7.

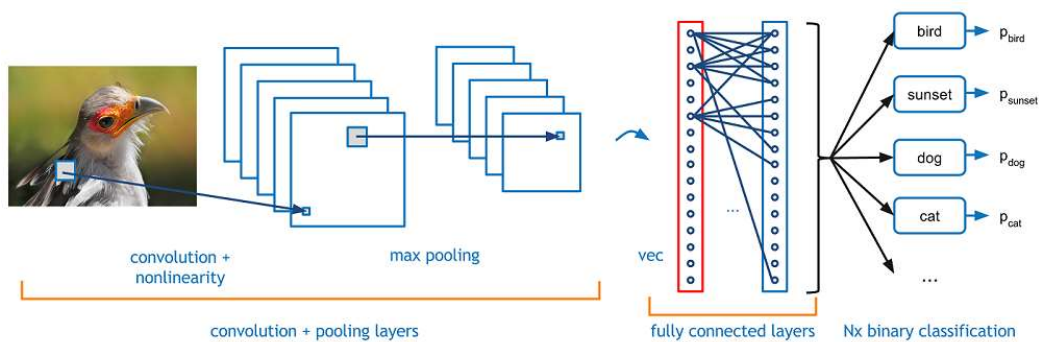


Σχ. 3.7.7.

### Επίπεδο Ταξινόμησης

Σε αυτό το επίπεδο οι νευρώνες έχουν πλήρεις συνδέσεις με όλες τις εξόδους των νευρώνων του προηγούμενου επιπέδου. Το πλήρως συνδεδεμένο επίπεδο (Fully Connected Layer) είναι ένας Perceptron πολλών επιπέδων (MLP), που

χρησιμοποιεί ως συνάρτηση κόστους την πολυωνυμική λογιστική (Multinomial Logistic) στο επίπεδο εξόδου. Για την είσοδο των δεδομένων στο επίπεδο αυτό θα πρέπει αυτά να δομήσουν ένα μονοδιάστατο πίνακα (άνυσμα στήλης). Στο ακόλουθο Σχήμα 3.7.8 παρουσιάζεται μια εκδοχή ενός CNN τεσσάρων κλάσεων.



Σχήμα 3.7.8.

### Υπολογισμός του πλήθους βαρών και εκπαίδευση στα ΣΝΔ

**Επίπεδο συνέλιξης :** Το συγκεκριμένο επίπεδο είναι αυτό που δημιουργεί τον πίνακα χαρακτηριστικών, οπότε έχουμε πίνακες βαρών. Οι αρχικές τιμές των βαρών είναι τυχαίες. Ο αριθμός των βαρών (παραμέτρων) σε ένα επίπεδο συνέλιξης με  $N$  κανάλια και  $M$  φίλτρα διαστάσεων  $r \times c$  θα είναι:  $r \times c \times N \times M$

**Επίπεδο δειγματοληψίας :** Στο επίπεδο δειγματοληψίας μειώνουμε το μέγεθος των διαστάσεων. Για παράδειγμα αν σε έναν χάρτη χαρακτηριστικών  $224 \times 224$  στοιχείων εφαρμοστεί δειγματοληψία σε  $2 \times 2$  υπό-περιοχές τότε αριθμός των στοιχείων του χάρτη χαρακτηριστικών είναι  $112 \times 112$ . Στο επίπεδο δειγματοληψία δεν υπάρχουν βάρη.

**Πλήρως συνδεδεμένο επίπεδο (FC) :** Το συγκεκριμένο επίπεδο έχει τον υψηλότερο αριθμό βαρών από κάθε άλλο επίπεδο. Για τον υπολογισμό του πλήθους των βαρών υπολογίζεται το γινόμενο του πλήθους των νευρώνων κάθε επιπέδου επί το πλήθος των εξόδων του προηγούμενου επιπέδου. Έτσι το πλήθος των βαρών είναι:  $\text{πλήθος νευρώνων στο τρέχον επίπεδο} \times (\text{πλήθος εξόδων στο προηγούμενο επίπεδο} + 1)$ , το +1 αφορά τον σταθερό όρο. Για παράδειγμα στο **CNN VGGNET** το πλήθος των βαρών φαίνεται στον ακόλουθο Πίνακα.

```

INPUT: [224x224x3]      memory: 224*224*3=150K  weights: 0
CONV3-64: [224x224x64] memory: 224*224*64=3.2M  weights: (3*3*3)*64 = 1,728
CONV3-64: [224x224x64] memory: 224*224*64=3.2M  weights: (3*3*64)*64 = 36,864
POOL2: [112x112x64]    memory: 112*112*64=800K  weights: 0
CONV3-128: [112x112x128] memory: 112*112*128=1.6M  weights: (3*3*64)*128 = 73,728
CONV3-128: [112x112x128] memory: 112*112*128=1.6M  weights: (3*3*128)*128 = 147,456
POOL2: [56x56x128]    memory: 56*56*128=400K  weights: 0
CONV3-256: [56x56x256] memory: 56*56*256=800K  weights: (3*3*128)*256 = 294,912
CONV3-256: [56x56x256] memory: 56*56*256=800K  weights: (3*3*256)*256 = 589,824
CONV3-256: [56x56x256] memory: 56*56*256=800K  weights: (3*3*256)*256 = 589,824
POOL2: [28x28x256]    memory: 28*28*256=200K  weights: 0
CONV3-512: [28x28x512] memory: 28*28*512=400K  weights: (3*3*256)*512 = 1,179,648
CONV3-512: [28x28x512] memory: 28*28*512=400K  weights: (3*3*512)*512 = 2,359,296
CONV3-512: [28x28x512] memory: 28*28*512=400K  weights: (3*3*512)*512 = 2,359,296
POOL2: [14x14x512]    memory: 14*14*512=100K  weights: 0
CONV3-512: [14x14x512] memory: 14*14*512=100K  weights: (3*3*512)*512 = 2,359,296
CONV3-512: [14x14x512] memory: 14*14*512=100K  weights: (3*3*512)*512 = 2,359,296
CONV3-512: [14x14x512] memory: 14*14*512=100K  weights: (3*3*512)*512 = 2,359,296
POOL2: [7x7x512]     memory: 7*7*512=25K   weights: 0
FC: [1x1x4096]        memory: 4096          weights: 7*7*512*4096 = 102,760,448
FC: [1x1x4096]        memory: 4096          weights: 4096*4096 = 16,777,216
FC: [1x1x1000]        memory: 1000          weights: 4096*1000 = 4,096,000

TOTAL memory: 24M * 4 bytes ~ = 93MB / image (only forward! ~*2 for bwd)
TOTAL params: 138M parameters

```

Πίνακας 3.7.1: Πλήθος των βαρών στο CNN VGGNet 16

## Εκπαίδευση των Συνελικτικών Νευρωνικών Δικτύων

Στα CNN η συνάρτηση κόστους βασίζεται στην διασταυρούμενη εντροπία (cross entropy) και όχι στο μέσο τετραγωνικό σφάλμα. Η προσέγγιση ακρότατου σημείου (ελαχίστου) επιδιώκεται κυρίως με μεθόδους όπως η SGDM (Stochastic Gradient Descent with Momentum) και η ADAM (Adaptive Moment Estimation). Στο τελικό επίπεδο εξόδου του ταξινομητή, συνάρτηση ενεργοποίησης είναι η *softmax* που μετατρέπει τις τιμές των αθροιστών σε πιθανότητες. Ακολούθως θα παρουσιάσουμε την διαδικασία εκπαίδευσης και διόρθωσης σφάλματος στα επίπεδα ενός CNN.

## Εκπαίδευση στο επίπεδο πλήρους διασύνδεσης

Στο τελικό επίπεδο εξόδου το πλήθος των νευρώνων ορίζεται να είναι ίσο με το πλήθος των κλάσεων, έστω  $M$ . Κάθε κλάση αντιστοιχίζεται σε έναν νευρώνα (τον ίδιο πάντοτε) και ονοματίζεται με την θέση του  $m=1..M$  στο επίπεδο εξόδου. Έστω το  $\mathbf{x} = [x_1 \dots x_n \dots x_N, 1]$  επαυξημένο διάνυσμα που εισέρχεται στο επίπεδο εξόδου και  $\mathbf{w}_m = [w_{1m}, \dots, w_{nm}, \dots, w_{Nm}, w_{0m}]^T$  τα επαυξημένα διανύσματα των βαρών των νευρώνων σ' αυτό.

Αρχικά υπολογίζονται οι ποσότητες  $\sigma_m = \mathbf{w}_m^T \cdot \mathbf{x}$ . Οι τιμές αυτές μετασχηματίζονται σύμφωνα την συνάρτηση softmax σε

$$y_m = \frac{e^{\mathbf{w}_m^T \cdot \mathbf{x}}}{\sum_k e^{\mathbf{w}_k^T \cdot \mathbf{x}}}$$

Οι τιμές  $y_m$  είναι θετικές και το άθροισμα τους  $\sum_{m=1}^M y_m = 1$ . Επιθυμούμε οι τιμές εξόδου  $y_m$  να είναι τέτοιες ώστε αν το άνυσμα εισόδου ανήκει στην κλάση  $c$ ,  $\mathbf{x}^c$ ,  $c \in \{1, \dots, M\}$  τότε η έξοδος  $y_c$  του  $c$  νευρώνα να είναι μεγαλύτερη των άλλων εξόδων και ει δυνατόν ίση με την μονάδα εκφράζοντας έτσι υψηλή πιθανότητα να ανήκει η είσοδος στην επιθυμητή κλάση. Με όρους πιθανοτήτων αν  $1, \dots, M$  είναι οι τιμές μιας τυχαίας διακριτής μεταβλητής  $\omega$ , τότε  $y_m = p(\omega = m | \mathbf{x}^c)$  είναι οι πιθανότητες που αντιστοιχούν στις τιμές τις  $\omega$  δηλαδή  $y_1 = p(1 | \mathbf{x}^c), \dots, y_m = p(m | \mathbf{x}^c), \dots, y_M = p(M | \mathbf{x}^c)$ . Το ιδανικό θα ήταν οι πιθανότητες  $p(\omega = m | \mathbf{x}^c)$  να ήταν ή να συνέκλιναν με πιθανότητες

$$q(\omega = m | \mathbf{x}^c) = \begin{cases} 1 & \text{αν } m = c \\ 0 & \text{αν } m \neq c \end{cases}$$

Ένα μέτρο σύγκρισης των  $p(\omega)$  και  $q(\omega)$  δίνεται από την σχέση της διασταυρούμενης εντροπίας (cross entropy):

$$H(q, p) = - \sum_{\omega} q(\omega) \cdot \log(p(\omega)) = - \sum_{m=1}^M q(\omega = m | \mathbf{x}^c) \cdot \log(p(\omega = m | \mathbf{x}^c)) = - \sum_{m=1}^M q(\omega = m | \mathbf{x}^c) \cdot \log(y_m) = - \log(y_c)$$

Για παράδειγμα σε ένα σύστημα ταξινόμησης τεσσάρων κλάσεων εισέρχεται το άνυσμα  $\mathbf{x}^2$  και οι έξοδοι των τεσσάρων νευρώνων είναι

$$y_1=0.2, y_2=0.5, y_3=0.1, y_4=0.2$$

και προφανώς

$$q(1 | \mathbf{x}^2)=0, q(2 | \mathbf{x}^2)=1, q(3 | \mathbf{x}^2)=0, q(4 | \mathbf{x}^2)=0$$

$$H(q, p) = -(0 \cdot \log(0.2) + 1 \cdot \log(0.5) + 0 \cdot \log(0.1) + 0 \cdot \log(0.2)) = 2$$

Όσο η  $y_2$  τείνει στην μονάδα η  $H(q, p)$  θα τείνει στο μηδέν.

Κατόπιν αυτού μπορούμε να εισάγουμε ως συνάρτηση κόστους (cost function) ή απώλειας (loss) την σχέση

$$L(\mathbf{W}) = - \sum_i \sum_m q(\omega = m | \mathbf{x}_i^c) \cdot \log(y_{mi}) = - \sum_i \log(y_{c,i})$$

που αθροίζει τις τιμές της cross entropy που προκύπτουν από την εισαγωγή στο σύστημα κάθε προτύπου  $\mathbf{x}_i$ ,  $i=1\dots I$ ,  $I$  το πλήθος των προτύπων.

$$y_{c,i} = \frac{e^{\mathbf{w}_c^T \cdot \mathbf{x}_i}}{\sum_m e^{\mathbf{w}_m^T \cdot \mathbf{x}_i}}$$

Για τον νευρώνα  $j$  στο επίπεδο εξόδου  $j \in \{1 \dots M\}$  και  $\mathbf{w}_j$  το άνωσμο στήλης των βαρών του ισχύει

$$\frac{\partial L}{\partial \mathbf{w}_j} = \frac{\partial}{\partial \mathbf{w}_j} \left( - \sum_i \log(y_{ci}) \right) = - \sum_i \frac{\partial}{\partial \mathbf{w}_j} \log(y_{ci}) = - \sum_i \frac{\partial}{\partial \mathbf{w}_j} \log \left( \frac{e^{\mathbf{w}_c^T \cdot \mathbf{x}_i}}{\sum_m e^{\mathbf{w}_m^T \cdot \mathbf{x}_i}} \right)$$

Για τον υπολογισμό της παραγώγου γράφουμε τα παραπάνω υπό μορφή σύνθετης συνάρτησης. Θα παραλείψουμε για λόγους διευκόλυνσης της γραφής τον δείκτη  $i$  (θέτω  $\mathbf{x} = \mathbf{x}_i$ )

$$\sigma_j = \mathbf{w}_j^T \cdot \mathbf{x}, \quad y_c = \frac{e^{\sigma_c}}{\sum_m e^{\sigma_m}}, \quad l_c = \log(y_c)$$

$$\frac{\partial l}{\partial \mathbf{w}_j} = \frac{\partial l}{\partial y_c} \frac{\partial y_c}{\partial \sigma_j} \frac{\partial \sigma_j}{\partial \mathbf{w}_j}$$

$$\frac{\partial \sigma_j}{\partial \mathbf{w}_j} = \mathbf{x}_i^T, \quad \frac{\partial l_i}{\partial y_c} = \frac{1}{y_c}, \quad \frac{\partial y_c}{\partial \sigma_j} = \frac{\partial}{\partial z_j} \left( \frac{e^{\sigma_c}}{\sum_m e^{\sigma_m}} \right) = \begin{cases} \frac{e^{\sigma_c} \cdot \sum_m e^{\sigma_m} - e^{\sigma_c} \cdot e^{\sigma_c}}{(\sum_m e^{\sigma_m})^2} = y_c(1 - y_c) & \text{αν } j = c \\ \frac{-e^{\sigma_c} \cdot e^{\sigma_j}}{(\sum_m e^{\sigma_m})^2} = -y_c \cdot y_j & \text{αν } j \neq c \end{cases}$$

$$\text{Άρα } \frac{\partial l_i}{\partial \mathbf{w}_j} = \begin{cases} (1 - y_c) \cdot \mathbf{x}_i^T & \text{αν } j = c \\ -y_j \cdot \mathbf{x}_i^T & \text{αν } j \neq c \end{cases}, \quad \text{για κάθε συναπτικό βάρος του } j \text{ νευρώνα}$$

$$\frac{\partial l_i}{\partial w_{nj}} = \begin{cases} (1 - y_c) \cdot x_{ni} & \text{αν } j = c \\ -y_j \cdot x_{ni} & \text{αν } j \neq c \end{cases} \quad \text{και τελικά}$$

$$\frac{\partial L}{\partial w_{nj}} = - \sum_i \frac{\partial l_i}{\partial w_{nj}}$$

### Αλγόριθμος οπισθοδρόμησης σφάλματος στο επίπεδο δειγματοληψίας.

Τα επίπεδα δειγματοληψίας (pooling) δεν διαθέτουν παραμέτρους που θα αλλάξουν κατά την εκπαίδευση του ΣΝΔ. Αυτό που έχουμε μόνο να κάνουμε είναι να υπολογίσουμε τις κλίσεις σε σχέση με το προηγούμενο επίπεδο. Αν  $R$  είναι μια περιοχή του επιπέδου που υποδειγματοληπτείται με την διαδικασία του maxpooling η έξοδος θα είναι  $y = \max\{x_i \in R\}$  και  $\frac{\partial y}{\partial x_i} = \begin{cases} 1 & \text{αν } y = x_i \\ 0 & \text{αν } y \neq x_i \end{cases}$ . Για παράδειγμα αν η περιοχή είναι  $2 \times 2$  θα έχουμε:

$x_1=5$	$x_2=8$	->Max Pooling-> $y=8$
$x_3=3$	$x_4=1$	

$\frac{\partial y}{\partial x_1} = 0$	$\frac{\partial y}{\partial x_2} = 1$
$\frac{\partial y}{\partial x_3} = 0$	$\frac{\partial y}{\partial x_4} = 0$

Σύμφωνα με τα παραπάνω η διόρθωση των βαρών του φίλτρου θα γίνει μόνο για την θέση του επί του πίνακα εισόδου (κατά την συνελικτική διαδικασία) που παρήχθη η μέγιστη τιμή (το 8). Στην περίπτωση του average pooling, η κλίση θα είναι παντού  $1/n$  όπου  $n$  είναι το πλήθος των στοιχείων του παραθύρου (για το παράδειγμα μας  $1/4$ ,  $n=4$ ).

### Αλγόριθμος οπισθοδρόμησης σφάλματος στο επίπεδο συνέλιξης

Στο συνελικτικό επίπεδο ένα φίλτρο είναι ένας νευρώνας με βάρη που το πλήθος τους καθορίζεται από τις διαστάσεις του φίλτρου. Οι τιμές εισόδου προέρχονται από τις τιμές του καναλιού καθώς αυτές συνελίσσονται με τα βάρη του φίλτρου. Αυτό μπορεί να θεωρηθεί σε μια παράλληλη υλοποίηση όπως φαίνεται στο Σχήμα 3.7.9, για δεδομένα μιας διάστασης και φίλτρο εύρους τρία ( $3 \times 1$ ).

Έστω  $\delta_0 = \frac{\partial L}{\partial \sigma_0}$ ,  $\delta_1 = \frac{\partial L}{\partial \sigma_1}$ ,  $\delta_2 = \frac{\partial L}{\partial \sigma_2}$ ,  $\delta_3 = \frac{\partial L}{\partial \sigma_3}$  οι μεταβολές της συνάρτησης κόστους ως προς τα αθροίσματα (έξοδοι αθροιστών),  $\mathbf{w} = [w_0, w_1, w_2]$ ,  $\boldsymbol{\sigma} = [\sigma_0, \sigma_1, \sigma_2]^T$ . Θα ισχύουν οι σχέσεις

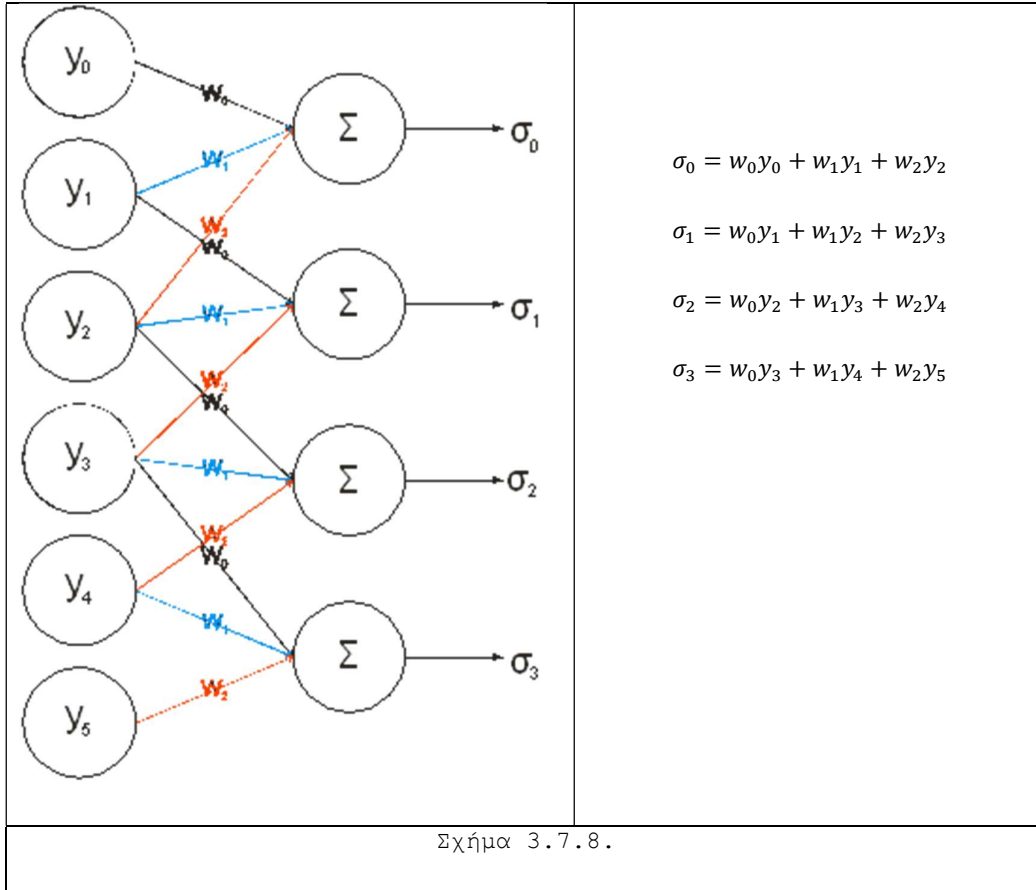
$$\sigma_0 = w_0 y_0 + w_1 y_1 + w_2 y_2, \quad \sigma_1 = w_0 y_1 + w_1 y_2 + w_2 y_3$$

$$\sigma_2 = w_0 y_2 + w_1 y_3 + w_2 y_4, \quad \sigma_3 = w_0 y_3 + w_1 y_4 + w_2 y_5$$

$$\frac{\partial L}{\partial w_0} = \frac{\partial L}{\partial \boldsymbol{\sigma}} \cdot \frac{\partial \boldsymbol{\sigma}}{\partial w_0} = [\delta_0, \delta_1, \delta_2, \delta_3] \cdot \begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ y_3 \end{bmatrix} = \delta_0 y_0 + \delta_1 y_1 + \delta_2 y_2 + \delta_3 y_3$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \boldsymbol{\sigma}} \cdot \frac{\partial \boldsymbol{\sigma}}{\partial w_1} = [\delta_0, \delta_1, \delta_2, \delta_3] \cdot \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} = \delta_0 y_1 + \delta_1 y_2 + \delta_2 y_3 + \delta_3 y_4$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial w_2} = [\delta_0, \delta_1, \delta_2, \delta_3] \cdot \begin{bmatrix} y_2 \\ y_3 \\ y_4 \\ y_5 \end{bmatrix} = \delta_0 y_2 + \delta_1 y_3 + \delta_2 y_4 + \delta_3 y_5$$



$$\sigma_0 = w_0 y_0 + w_1 y_1 + w_2 y_2$$

$$\sigma_1 = w_0 y_1 + w_1 y_2 + w_2 y_3$$

$$\sigma_2 = w_0 y_2 + w_1 y_3 + w_2 y_4$$

$$\sigma_3 = w_0 y_3 + w_1 y_4 + w_2 y_5$$

Για την οπισθοδιάδοση της διόρθωσης θα πρέπει να υπολογίσουμε τις ποσότητες:

$$\frac{\partial L}{\partial y_0}, \frac{\partial L}{\partial y_1}, \dots, \frac{\partial L}{\partial y_5}$$

$$\frac{\partial L}{\partial y_0} = \frac{\partial L}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial y_0} = [\delta_0, \delta_1, \delta_2, \delta_3] \cdot \begin{bmatrix} w_0 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \delta_0 w_0$$

$$\frac{\partial L}{\partial y_1} = \frac{\partial L}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial y_1} = [\delta_0, \delta_1, \delta_2, \delta_3] \cdot \begin{bmatrix} w_1 \\ w_0 \\ 0 \\ 0 \end{bmatrix} = \delta_0 w_1 + \delta_1 w_0$$

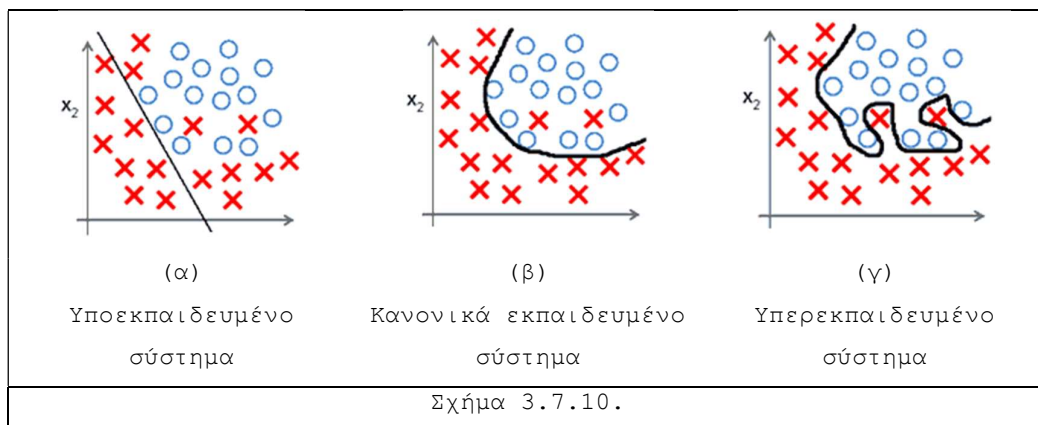
$$\frac{\partial L}{\partial y_2} = \delta_0 w_2 + \delta_1 w_1 + \delta_2 w_0, \quad \frac{\partial L}{\partial y_3} = \delta_1 w_2 + \delta_2 w_1 + \delta_3 w_0$$

$$\frac{\partial L}{\partial y_4} = \delta_2 w_2 + \delta_3 w_1, \quad \frac{\partial L}{\partial y_5} = \delta_3 w_2$$

Γενικότερα, σε ένα συνελικτικό επίπεδο, η κλίση του τρέχοντος επιπέδου προκύπτει με την συνέλιξη των  $\delta_i$  με τα συνελικτικά φίλτρα αντεστραμμένα.

Το πλήθος των επιπέδων και των νευρώνων των, είναι ένα ζήτημα που μας απασχολεί. Δεν υπάρχει ένας γενικός τύπος-αλγόριθμος που να μας οδηγεί σε ένα βέλτιστο πλήθος παραμέτρων-βαρών. Ένας μεγάλος αριθμός βαρών αυξάνει το υπολογιστικό κόστος και μπορεί να οδηγήσει σε λεπτομερή εκμάθηση του συνόλου εκπαίδευσης ακόμη και του θορύβου των μετρήσεων, εις βάρος της γενικότητας του, δηλαδή της επιτυχούς ταξινόμησης νέων αγνώστων προτύπων. Αυτό το φαινόμενο το ονομάζουμε *υπερεκπαίδευση*. Μικρός αριθμός παραμέτρων μπορεί να οδηγήσει σε αδυναμία εκπαίδευσης ή χαμηλή απόδοση-ακρίβεια ταξινόμησης. Την περίπτωση αυτή αναφέρουμε ως *υποεκπαίδευση*. Στο Σχήμα 3.7.10 φαίνεται η διαχωριστική καμπύλη τριών συστημάτων ενός προβλήματος δύο διαστάσεων για τις περιπτώσεις που αναφέραμε.

Μία απλή μέθοδος είναι να ξεκινάμε με ένα μικρό πλήθος παραμέτρων και να το αυξάνουμε προοδευτικά μέχρι να έχουμε ικανοποιητικό αποτέλεσμα. Εκείνο που συχνά εμφανίζεται όταν έχουμε υψηλή απόδοση εκπαίδευσης είναι η υπερεκπαίδευση που διαπιστώνεται όταν το σύστημά μας ταξινομεί πρότυπα για τα οποία δεν εκπαιδεύτηκε και παρουσιάζει χαμηλή απόδοση. Για να εντοπίσουμε νωρίς το πρόβλημα χωρίζουμε το αρχικό σύνολο  $S$  των διαθέσιμων δεδομένων σε δύο μέρη. Το ένα χρησιμοποιείται ως σύνολο εκπαίδευσης (Training set,  $T$ ) και το άλλο ως σύνολο επικύρωσης του αποτελέσματος (Validation set,  $V$ ). Θα πρέπει το ποσοστό επιτυχίας επιτυχούς ταξινόμησης (accuracy) να είναι υψηλό με παραπλήσιες τιμές.





## ΚΕΦΑΛΑΙΟ 4

### ΕΚΠΑΙΔΕΥΣΗ ΧΩΡΙΣ ΕΠΟΠΤΗ

Η εύρεση των συγκεντρώσεων των προτύπων όταν δεν είναι γνωστό το πλήθος τους και η μορφολογία τους είναι ένα πρόβλημα που απασχολεί την Υπολογιστική Νοημοσύνη, την ανάλυση και την εξόρυξη δεδομένων. Για την λύση του προτείνονται ενδιαφέρουσες τεχνικές με δομή ή όχι νευρωνικού δικτύου. Υπάρχουν ποικίλες ενδιαφέρουσες παραλλαγές του όπως η ιεραρχική συσταδοποίηση. Κρίσιμα ερωτήματα είναι το πλήθος των κλάσεων, οι αποστάσεις μεταξύ των προτύπων ή εσωτερική διασπορά και η μορφολογία των συγκεντρώσεων. Η μέτρηση πολλών χαρακτηριστικών και η ποικιλία των προτύπων είναι βασικοί παράγοντες που επιτείνουν την δυσκολία του προβλήματος.

Ακολούθως θα παρουσιάσουμε τρεις μεθόδους εκπαίδευσης χωρίς επόπτη. Οι δύο πρώτες είναι αλγόριθμοι που μπορούν να χρησιμοποιηθούν για την επίλυση σχετικών προβλημάτων ή να χρησιμοποιηθούν ως εργαλεία σε περισσότερο σύνθετες εργασίες. Η τρίτη είναι μία μέθοδος που βασίζεται στην λειτουργία ενός νευρωνικού δικτύου και μπορεί να αποτελέσει βάση για την σχεδίαση ταξινομητών. Η μέθοδος αυτή δίνει την δυνατότητα εποπτείας σε χώρων με περισσότερες των τριών διαστάσεων που η αναπαράστασή τους σε ένα σύστημα αξόνων είναι γεωμετρικά αδύνατη.

#### 4.1 Ο Αλγόριθμος ISODATA ή K-Μέσων (k-means ή c-means)

Ο Αλγόριθμος των K-μέσων τιτλοφορείται και ως ISODATA ή αλγόριθμος Lloyd's και είναι από τους δημοφιλέστερους αλγορίθμους συσταδοποίησης. Οι ομοιότητα των προτύπων περιγράφεται με την Ευκλείδεια Απόσταση. Κάθε συγκέντρωση καθορίζεται από ένα σημείο-αντιπρόσωπο στο χώρο των προτύπων και σ' αυτήν ανήκουν τα πρότυπα που απέχουν μικρότερη απόσταση από τον αντιπρόσωπο αυτής συγκριτικά με τους αντιπροσώπους των άλλων συγκεντρώσεων. Ο αντιπρόσωπος είναι ο μέσος όρος των προτύπων που αποτελούν την

συγκέντρωση και ονομάζεται κέντρο της κλάσης. Ακολούθως ο αλγόριθμος περιγράφεται για προκαθορισμένο πλήθος συγκεντρώσεων.

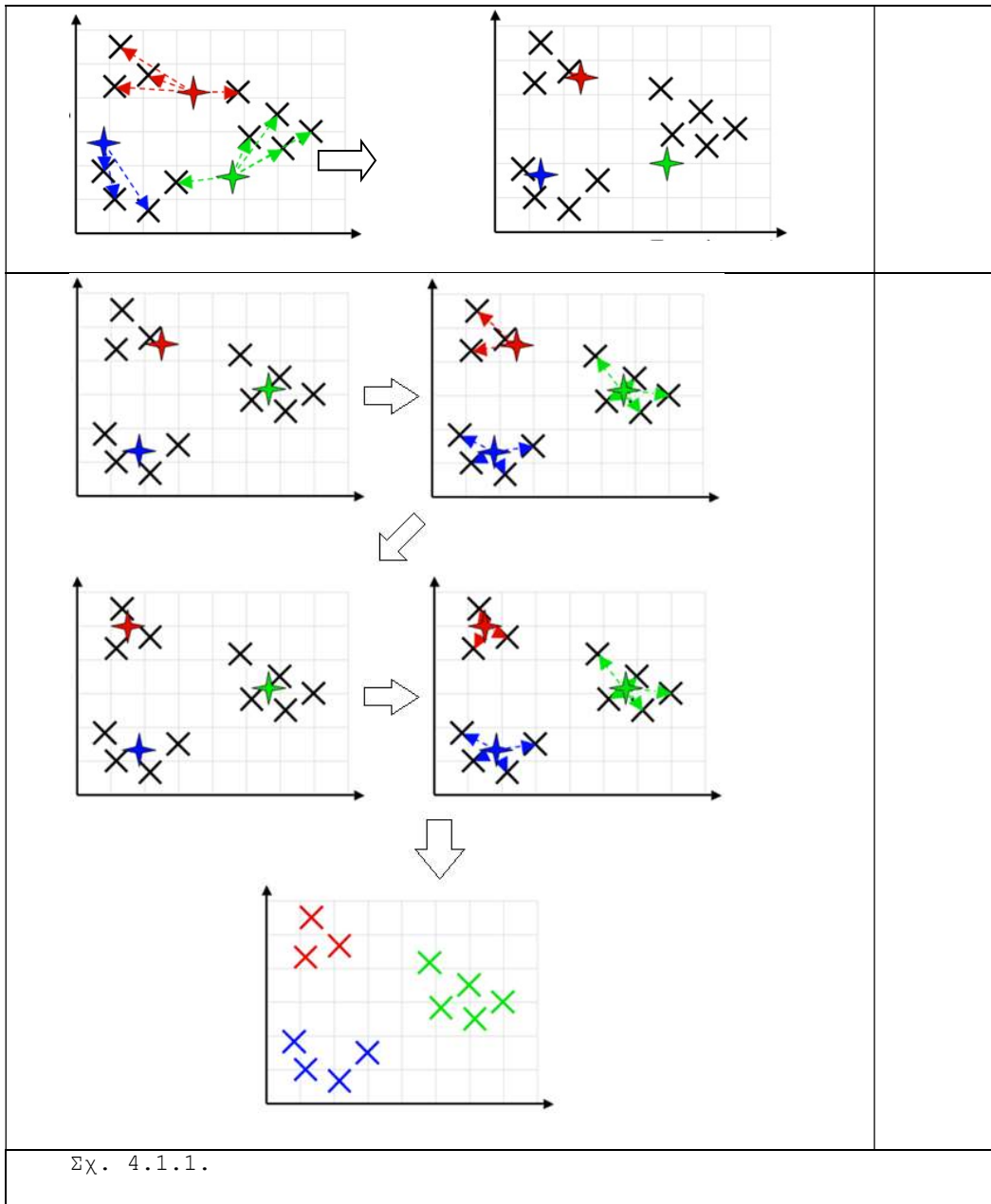
Έστω,

- I το πλήθος των προτύπων, K το πλήθος το συγκεντρώσεων,
- N το πλήθος των χαρακτηριστικών,
- $\mathbf{x}_i \in E^N$  πίνακας στήλης που περιγράφει το i πρότυπο,  $i=1..I$
- $\mathbf{c}_k \in E^N$  πίνακας στήλης που περιγράφει τον αντιπρόσωπο (κέντρο) της k κλάσης,  $k=1..K$
- B πίνακας  $I \times 1$  σε κάθε στοιχείο i του οποίου καταχωρείται η συγκέντρωση που ανήκει το i πρότυπο, (πχ  $b_i=k$ ),
- t μετρητής επανάληψης.

Για  $t=0$  αποδίδονται τυχαίες τιμές στα  $\mathbf{c}_k$ , για κάθε k.

<p>Για κάθε πρότυπο υπολογίζονται οι αποστάσεις του από τα κέντρα όλων των κλάσεων και καταχωρείται σ' αυτή που έχει κοντινότερο κέντρο στο πρότυπο.</p> <p>Για κάθε κλάση υπολογίζεται το πλήθος P των προτύπων που ανήκουν στην κλάση</p> <p>Υπολογίζεται το κέντρο της κλάσης</p> <p>Αν τα κέντρα παραμείνουν ίδια ο αλγόριθμος τερματίζεται</p>	<p>Επανάληψη</p> <p>Για i από 1 έως I</p> $b_i = \min_k \left\{ \sqrt{\sum_{v=1}^N (x_{vi} - c_{vk})^2} \right\}$ <p>Τέλος για.</p> <p>t = t+1</p> <p>Για k από 1 έως K</p> $P = \sum_{\forall b_i=k} 1$ $c_k = \frac{\sum_{\forall b_i=k} x_i}{P}$ <p>Τέλος για.</p> <p>Αν <math>c_k(t) = c_k(t+1)</math> για κάθε k, τέλος επανάληψης</p>
---	---

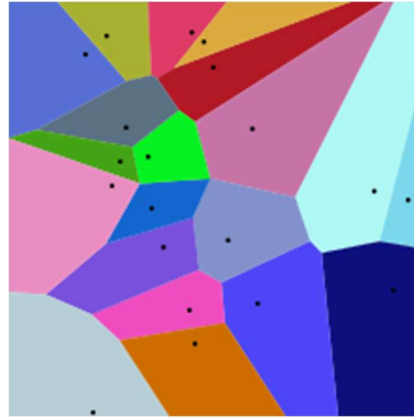
Στο Σχήμα 4.1.1 φαίνεται η μετακίνηση των κέντρων των συγκεντρώσεων σε διαδοχικές επαναλήψεις του αλγορίθμου.



Είναι δυνατόν επίσης κάποιο ή κάποια κέντρα να μην έχουν κανένα κοντινότερο πρότυπο και ως εκ τούτου να έχουμε λιγότερες από  $K$  συστάδες. Στην περίπτωση αυτή το κέντρο τοποθετείται κοντά σε πρότυπο που απέχει περισσότερο όλων από το κέντρο της συστάδας του.

Ο  $k$ -means θεωρείται ένας γρήγορος, αρκετά αποτελεσματικός αλγόριθμος με ευκολία στην υλοποίηση. Δίνει καλύτερα αποτελέσματα όταν το σύνολο των δεδομένων είναι δομημένο σε διαχώρισιμες συγκεντρώσεις. Μειονέκτημα είναι

ότι απαιτεί τον προκαθορισμό του αριθμού των κέντρων των συστάδων. Η τυχαία επιλογή της τιμής αυτών των κέντρων μπορεί να μη μας οδηγήσει σε καλά αποτελέσματα. Επίσης, ο αλγόριθμος παρουσιάζει αδυναμία στον χειρισμό δεδομένων με θόρυβο ή ακραίες τιμές και αποτυγχάνει για μη συμπαγές σύνολο δεδομένων (σφαιρικές συγκεντρώσεις). Οι μεσοκάθετοι των κέντρων αποτελούν τις γραμμικές διακριτικές συναρτήσεις των συγκεντρώσεων και συνθέτουν το λεγόμενο διάγραμμα Voronoi (Voronoi tessellation, Σχήμα 4.1.2.)



Σχ.4.1.2.

#### 4.2 Απεικόνιση αλυσίδας κοντινών γειτόνων.

Η απεικόνιση αλυσίδας (*chain mapping*) είναι μία μέθοδος απλή αλλά υπολογιστικά δαπανηρή που στοχεύει κυρίως στην εύρεση διαδοχικών γειτονικών σημείων που ονομάζουμε αλυσίδες. Η ιδιαίτερη μορφολογία των αλυσίδων δεν τις καθιστά ανιχνεύσιμες ως συγκεντρώσεις από αλγορίθμους όπως ο ISODATA. Χρησιμοποιείται συχνά στις μεθόδους ιεραρχικής συσταδοποίησης. Οδηγεί στην εποπτεία της κατανομής των προτύπων σε πολυδιάστατους χώρους και μπορεί να χρησιμοποιηθεί για την εκτίμηση του πλήθους και του περιεχομένου των συγκεντρώσεών τους. Σύμφωνα με αυτήν δημιουργούμε μία διαδρομή μεταβαίνοντας από κάθε πρότυπο στο γειτονικότερό του, ξεκινώντας από κάποιο τυχαίο. Καταγράφουμε διαδοχικά τις τιμές του δείκτη του προτύπου, του δείκτη του γείτονά του και της απόστασής τους, χωρίς να μεταβούμε σε πρότυπα από τα οποία έχουμε διέλθει. Δημιουργείται έτσι μια ακολουθία  $\mathbf{d}_n = [i, j, D_E(\mathbf{x}_i, \mathbf{x}_j)]$ . Θεωρούμε τον όρο  $m$  της ακολουθίας με την μεγαλύτερη τιμή απόστασης,  $\mathbf{d}_m = \max\{\mathbf{d}_n\}$ . Το σύνολο

$$C_1 = \{\mathbf{x}_i / \mathbf{x}_j \text{ έχει δείκτη την τιμή } i \text{ για } \mathbf{d}_n = [i, j, D_E(\mathbf{x}_i, \mathbf{x}_j)] \text{ με } n \leq m\}$$

αποτελεί την μία από δύο συγκεντρώσεις των του συνόλου των προτύπων. Η δεύτερη κλάση είναι το σύνολο

$$C_2 = \{\mathbf{x}_j / \mathbf{x}_j \text{ έχει δείκτη την τιμή } j \text{ για } \mathbf{d}_n = [i, j, D_E(\mathbf{x}_i, \mathbf{x}_j)] \text{ με } n \geq m\}$$

Θεωρώντας την επόμενη μεγαλύτερη τιμή  $m'$  που θα αντιστοιχεί σε στοιχεία μιας από τις δύο συγκεντρώσεις μπορούμε να διαχωρίσουμε αυτήν σε δύο μέρη όπως προηγουμένως κ.ο.κ.

Η μέθοδος απαιτεί τον υπολογισμό του πίνακα των αποστάσεων των προτύπων μεταξύ τους και είναι πολυπλοκότητας  $O(n^2)$ .

### ΠΑΡΑΔΕΙΓΜΑ

Έστω τα πρότυπα

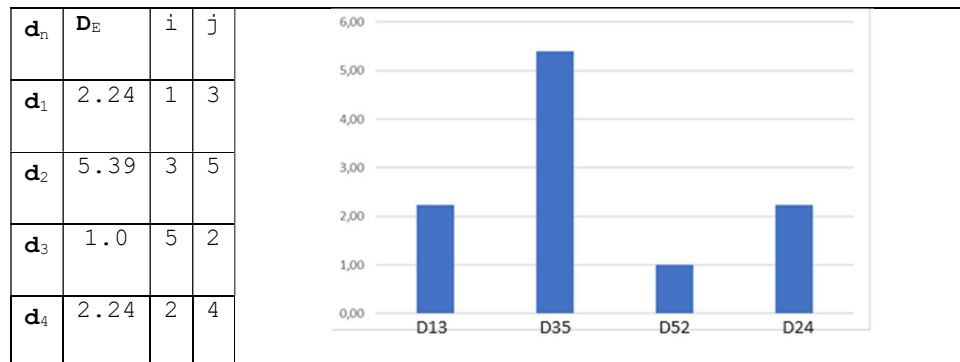
$$\mathbf{x}_1 = [10, 8]^T, \mathbf{x}_2 = [9, 7]^T, \mathbf{x}_3 = [1, 10]^T, \mathbf{x}_4 = [2, 8]^T, \mathbf{x}_5 = [4, 1]^T$$

Ο πίνακας των αποστάσεων τους είναι

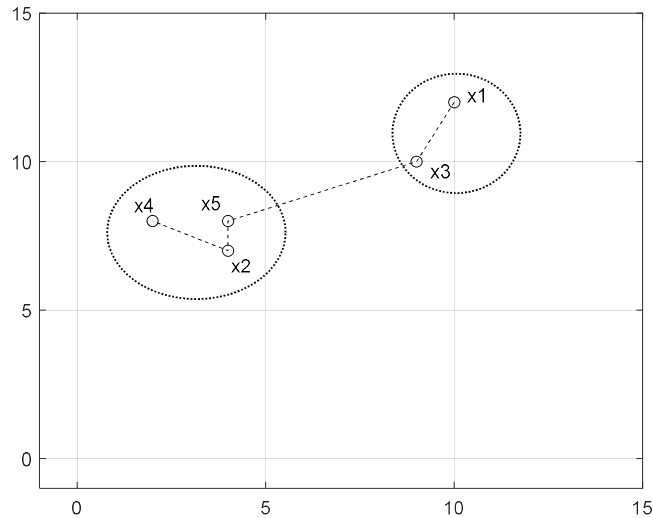
$D_{12}$ 7.81	$D_{13}$ 2.24	$D_{14}$ 8.94	$D_{15}$ 7.21
	$D_{23}$ 5.83	$D_{24}$ 2.24	$D_{25}$ 1.0
		$D_{34}$ 7.28	$D_{35}$ 5.39
			$D_{45}$ 2.0

Ξεκινώντας από το σημείο  $\mathbf{x}_1$  μεταβαίνουμε στο γειτονικότερό του κ.ο.κ.

Καταγράφεται η ακολουθία



Η τιμή 5.39 είναι η μεγαλύτερη στην στήλη των αποστάσεων και αφορά τον όρο  $m=2$ . Για  $n \leq 2$  ο δείκτης  $i$  με τιμές 1, 3 καθορίζει τα πρότυπα της συγκεντρώσης  $C_1 = \{\mathbf{x}_1, \mathbf{x}_3\}$  και για  $n \geq 2$  ο δείκτης  $j$  με τιμές 5, 2, 4 καθορίζει τα πρότυπα της συγκεντρώσης  $C_2 = \{\mathbf{x}_5, \mathbf{x}_4, \mathbf{x}_3\}$ , Σχήμα 4.2.1.



Σχήμα 4.2.1

### 4.3 Χάρτης απεικόνισης χαρακτηριστικών με αυτό-οργάνωση

Ο Χάρτης Απεικόνισης Χαρακτηριστικών με Αυτό-οργάνωση (SOFM: *Self organized Feature Map, SOM*) στοχεύει στην αναπαράσταση των σημείων ενός πολυδιάστατου χώρου που είναι ο αρχικός χώρος των προτύπων σε ένα χώρο λιγότερων διαστάσεων, συνήθως δύο διαστάσεων, διατηρώντας κατά το δυνατόν την τοπολογική-γεωμετρική δομή και τα στατιστικά χαρακτηριστικά των δεδομένων του αρχικού χώρου. Κατά την αναπαράσταση αυτή εκτός από την μείωση των διαστάσεων γίνεται και μείωση του πλήθους των σημείων στον τελικό χώρο που ονομάζουμε *χάρτη απεικόνισης των χαρακτηριστικών*. Βέβαια αφού τα σημεία του χάρτη είναι λιγότερα από αυτά του αρχικού είναι προφανώς ότι συχνά σε ένα σημείο του χάρτη θα απεικονίζονται πολλά σημεία του αρχικού χώρου. Εκείνο όμως που πρωτίστως πρέπει να διασφαλίζεται είναι ότι η γειτνίαση δύο σημείων του αρχικού χώρου θα συνεπάγεται και γειτνίαση των απεικονίσεών τους στο χάρτη (στο ίδιο ή σε γειτονικά σημεία). Αυτό μπορεί να διευκολύνει την οπτικοποίηση και την ανάλυση δεδομένων υψηλών διαστάσεων.

Ένας SOM είναι ένα τεχνητό νευρωνικό δίκτυο ενός επιπέδου που ονομάζεται και *επίπεδο ανταγωνισμού (competitive layer)*. Οι SOM προτάθηκαν από τον T.Kohonen<sup>[1]</sup> και συχνά καλούνται *Χάρτες Kohonen* ή *Δίκτυα Kohonen*.

[1] Kohonen, Teuvo (1982). "Self-Organized Formation of Topologically Correct Feature Maps". *Biological Cybernetics*. **43** (1): 59-69.

Τα σημεία του επιπέδου ανταγωνισμού (χάρτης) λέγονται *κόμβοι* και λειτουργούν ως νευρώνες, δηλαδή έχουν μεταβλητές-βάρος και είναι συνήθως

διατεταγμένοι σε τετραγωνικό ή εξαγωνικό πλέγμα δύο διαστάσεων που λέγεται και *κάναβος (grid)*, Σχήματα 4.3.1, 4.3.2. Ένα δίκτυο Kohonen είναι μια υπολογιστική προσέγγιση που βασίζεται σε ένα μηχανισμό ανταγωνισμού που παρατηρήθηκε σε βιολογικά συστήματα νευρώνων.

Κατά την φάση της εκπαίδευσης του δικτύου Kohonen δεν επιδιώκεται η ελαχιστοποίηση κάποιας συνάρτησης κόστους. Η διαδικασία εκπαίδευσης είναι επαναληπτική με προκαθορισμένο αριθμό επαναλήψεων που καθορίζεται από μια τιμή  $T$ . Χαρακτηριστικά του αλγορίθμου είναι η έννοια του *νικητή νευρώνα* και της *γειτονιάς* του. Ένας νευρώνας ονομάζεται νικητής κατά την φάση της εκπαίδευσης όπως θα αναλυθεί παρακάτω. Οι γειτονικοί του νευρώνες όπως ορίζονται είναι όσοι βρίσκονται σε μία περιοχή του κανάβου που έχει κέντρο τον νικητή. Αν η περιοχή είναι τετραγωνική σε ένα τετραγωνικό κανάβο η πλευρά της γειτονιάς μπορεί να είναι  $2 \times d + 1$  όπως φαίνεται στο σχήμα 4.3.2. Στον αλγόριθμο χρησιμοποιείται και η παράμετρος εκμάθησης  $\alpha$ . Το πλάτος της γειτονιάς και η παράμετρος εκμάθησης μεταβάλλονται κατά την διάρκεια της εκπαίδευσης.

Η εκπαίδευση του νευρωνικού δικτύου γίνεται ως εξής:

Καθορίζεται ο *κάναβος* του επιπέδου ανταγωνισμού, π.χ. όπως στο Σχήμα 4.3.1 όπου το επίπεδο ανταγωνισμού είναι τοποθετημένο σε έναν *κάναβο*  $8 \times 8$ . Αν και ο *κάναβος* είναι *δισδιάστατος*, οι νευρώνες έχουν αριθμηθεί με ένα δείκτη  $k$  που παίρνει ακέραιες τιμές από 0 έως 63. Αν  $\mathbf{x} \in R^N$  τότε κάθε νευρώνας  $k$  στον *κάναβο* έχει ένα διάνυσμα βαρών  $\mathbf{w}_k = [w_{1k} \dots w_{nk} \dots w_{Nk}]^T$ .

Εκτελούνται τα παρακάτω βήματα:

Ορίζουμε ένα μετρητή επανάληψης  $t=0,1,\dots,T$ . Αποδίδονται τυχαίες τιμές στα βάρη των νευρώνων  $w_{nk}(0)$  με τιμές που προκύπτουν από την πρόσθεση μικρών τυχαίων αριθμών στη μέση τιμή των ανυσμάτων  $\mathbf{x}$  όλου το συνόλου εκπαίδευσης.

Αρχικοποιείται η παράμετρος εκμάθησης  $\alpha(0)$  με μια μικρή θετική τιμή, συνήθως μεταξύ 0.2 και 0.5.

Αρχικοποιείται η παράμετρος  $d(0)$  με την τιμή 4, που είναι ίση με το μισό του εύρους του κανάβου.

Επιλέγεται ένα τυχαίο άνυσμα  $\mathbf{x}(t)$  από το σύνολο εκπαίδευσης.

Υπολογίζεται η έξοδος  $o_k(t)$  του νευρώνα  $k$  σύμφωνα με τη σχέση:

$$o_k(t) = \|\mathbf{x}(t) - \mathbf{w}_k(t)\| = \sqrt{\sum_{n=1}^N (x_n(t) - w_{nk}(t))^2}. \quad (4.3.1)$$

Είναι φανερό πως η  $o_k$  είναι η Ευκλείδεια απόσταση μεταξύ των  $\mathbf{x}(t)$  και  $\mathbf{w}_k(t)$ .

Ο νευρώνας  $c$  ανακηρύσσεται νικητής εάν ικανοποιείται η συνθήκη  $o_c(t) = \min\{o_k(t)\}$ . Εάν οι έξοδοι δύο νευρώνων είναι ίσες, τότε κατά σύμβαση επιλέγεται αυτός με το μικρότερο δείκτη.

Τα βάρη  $\mathbf{w}_k$  ανανεώνονται σύμφωνα με τις παρακάτω σχέσεις:

$$w_{nk}(t+1) = w_{nk}(t) + \Delta w_{nk}(t) \quad (4.3.3)$$

$$\Delta w_{nk} = \begin{cases} \alpha(t) \cdot (x_n(t) - w_{nk}(t)) & \text{αν } k \in N_c \\ 0 & \text{αν } k \notin N_c \end{cases} \quad (4.3.4)$$

όπου  $N_c$  το σύνολο των δεικτών των νευρώνων που βρίσκονται μέσα σε ένα τετράγωνο του δισδιάστατου κανάβου με κέντρο το νευρώνα νικητή και πλευρά  $2 \times d(t) + 1$  (γειτονιά του νικητή νευρώνα).

Αυξάνεται η μεταβλητή επανάληψης κατά ένα και αποδίδονται νέες τιμές στις μεταβλητές  $\alpha(t)$ ,  $d(t)$  σύμφωνα με τις σχέσεις:

$$a(t) = a(0) \cdot \left(1 - \frac{t}{T}\right) \quad (4.3.5)$$

$$d(t) = d(0) \cdot \left(1 - \frac{t}{T}\right) \quad (4.3.6)$$

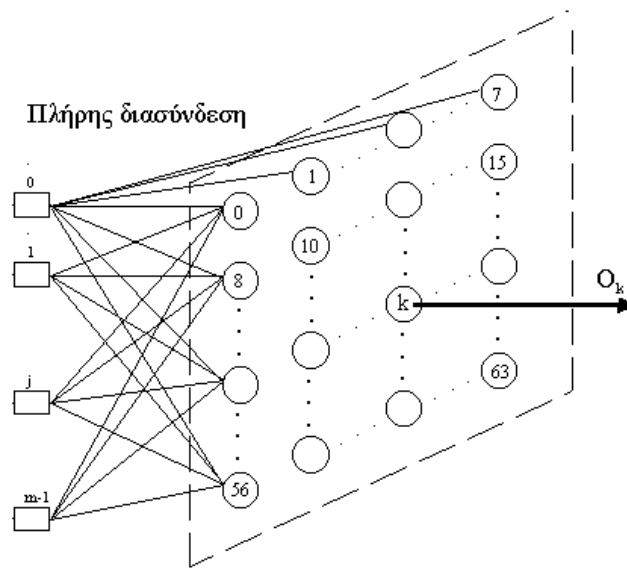
Τα βήματα 4 έως 8 επαναλαμβάνονται έως ότου η μεταβλητή  $t$  πάρει τη μέγιστη τελική τιμή  $T$ . Είναι φανερό πως οι μεταβλητές  $\alpha(t)$  και  $d(t)$  συγκλίνουν στο μηδέν καθώς η  $t$  τείνει στην τιμή  $T$ .

Μετά την εκπαίδευση κάθε άνυσμα εισόδου του SOM αντιστοιχίζεται- απεικονίζεται από τον νευρώνα που ανακηρύσσει νικητή. Κάθε νευρώνας του επιπέδου εξόδου αντιπροσωπεύει μία ομάδα προτύπων (cluster). Πρότυπα με μεγάλη ομοιότητα (γειτονικά στον αρχικό χώρο) αντιπροσωπεύονται από τον ίδιο νευρώνα ή γειτονικούς στο επίπεδο ανταγωνισμού. Νευρώνες που αντιπροσωπεύουν μεγάλο πλήθος σημείων του αρχικού χώρου είναι σημαντικοί ως κέντρα συγκεντρώσεων (κυρίαρχοι νευρώνες), Σχήμα 4.3.3. Νευρώνες που αντιπροσωπεύουν μικρό πλήθος σημείων μπορούν να απορριφθούν. Τα βάρη των νευρώνων έχουν τις διαστάσεις του αρχικού χώρου και ως εκ τούτου οι

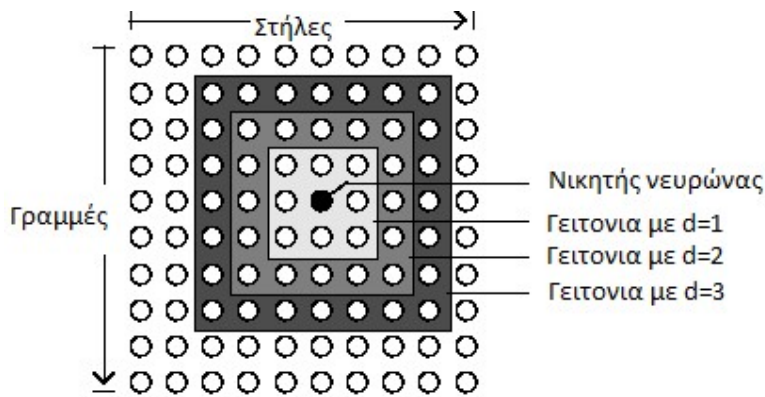


κυρίαρχοι νευρώνες αποτελούν μία αντιπροσωπευτική υπογραφή των αρχικών πολυπληθών δεδομένων.

Αν στους νευρώνες αποδοθούν ετικέτες (Labels) τότε το νευρωνικό δίκτυο Kohonen μπορεί να είναι προ-επεξεργαστής ενός συστήματος ταξινόμησης. Τα δίκτυα Kohonen αποτελούν εξαιρετικό εργαλείο ανάλυσης δεδομένων σε πολλά επιστημονικά πεδία. Εξελιγμένες επεκτάσεις τους είναι οι αναπτυσσόμενοι χάρτες αυτοοργάνωσης (*growing self-organizing map, GSOM*), οι σύμμορφοι χάρτες (*conformal maps*), οι χρονικά προσαρμοστικοί χάρτες αυτοοργάνωσης, (*time adaptive self-organizing map, TASOM*), οι προσανατολισμένοι και επεκτάσιμοι χάρτες (*oriented and scalable maps*).

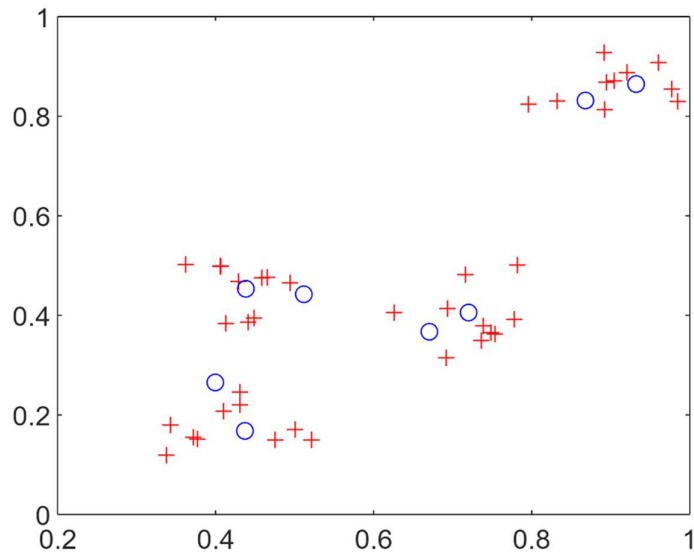


Σχήμα 4.3.1



Επίπεδο ανταγωνισμού σε τετραγωνικό κάναβο 10x10

Σχήμα 4.3.2



Σχήμα 4.3.3: Οκτώ νευρώνες (μπλε κύκλοι) περιγράφουν τα δεδομένα.

## ΚΕΦΑΛΑΙΟ 5

### ΑΝΑΛΥΣΗ ΧΑΡΑΚΤΗΡΙΣΤΙΚΩΝ

Μια σημαντική εργασία σε ένα σύστημα αναγνώρισης είναι η ανάλυση των μετρούμενων χαρακτηριστικών των προτύπων. Με την ανάλυση των χαρακτηριστικών πετυχαίνουμε την αξιολόγηση τους, την επιλογή των καταλληλότερων για την ταξινόμηση των προτύπων και την ελάττωση του πλήθους τους. Για παράδειγμα ένα χαρακτηριστικό με κοντινές τιμές μέσα σε κάθε μία κλάση και σαφώς διαφορετικές τιμές μεταξύ των κλάσεων είναι «καλό» χαρακτηριστικό ταξινόμησης. Με κατάλληλες μεθόδους ανάλυσης χαρακτηριστικών μπορούμε να οδηγηθούμε και σε νέα χαρακτηριστικά που προκύπτουν από τον γραμμικό συνδυασμό των αρχικών και είναι καταλληλότερα από αυτά. Ακολούθως θα δούμε μεθόδους ανάλυσης χαρακτηριστικών σε συστήματα εκμάθησης με και χωρίς επόπτη.

#### 5.1 Ανάλυση χαρακτηριστικών στην εκπαίδευση με επόπτη

Μια σημαντική εργασία σε ένα σύστημα αναγνώρισης είναι η επιλογή ενός μικρού συνόλου κατάλληλων χαρακτηριστικών από ένα πολύ μεγαλύτερο. Η επιλογή των ισχυρών χαρακτηριστικών είναι κρίσιμη για την αποτελεσματικότητα της ταξινόμησης. Για να ελαττωθεί το πλήθος των χαρακτηριστικών με την επιλογή ενός συνόλου «καλών χαρακτηριστικών» που βοηθούν την ταξινόμηση, εκτελείται μια διαδικασία αξιολόγησης με κριτήρια την *διαχωριστικότητα (separability)* και την *συσχέτιση (correlation)* των χαρακτηριστικών.

Για τον έλεγχο της ικανότητας διαχωρισμού καθορίζεται ο παράγοντας διαχωριστικότητας (*separability factor*)  $S_{avg}^l$  μεταξύ δύο ισοπίθανων κλάσεων A και B για κάθε χαρακτηριστικό  $v$  από τη σχέση:

$$S_v = \frac{|\mu_v^A - \mu_v^B|}{\sqrt{(\sigma_v^A)^2 + (\sigma_v^B)^2}} \quad (5.1)$$

όπου  $\mu_v^A$ ,  $\mu_v^B$ , είναι οι μέσες τιμές και  $\sigma_v^A$ ,  $\sigma_v^B$  οι τυπικές αποκλίσεις σε κάθε κλάση. Μεγάλη διαχωριστικότητα σημαίνει ότι το χαρακτηριστικό έχει μεγάλη ικανότητα να διαχωρίζει τις δύο κλάσεις μεταξύ τους. Ισοδύναμα χρησιμοποιείται η ποσότητα

$$S_v^2 = \frac{(\mu_v^A - \mu_v^B)^2}{(\sigma_v^A)^2 + (\sigma_v^B)^2}$$

γνωστή ως λόγος διάκρισης κατά Fisher.

Για την ανάλυση της συσχέτισης των χαρακτηριστικών υπολογίζεται ο παράγοντας συσχέτισης (*correlation factor*) για κάθε δύο χαρακτηριστικά  $v$  και  $\lambda$  που ανήκουν στην ίδια τάξη  $p$ , σύμφωνα με τη σχέση

$$C_{v\lambda}^p = \frac{\frac{1}{K_p} \sum_{k=1}^{K_p} (x_{vk} - \mu_v)(x_{\lambda k} - \mu_\lambda)}{\sigma_v \sigma_\lambda} = \frac{\sigma_{v\lambda}}{\sigma_v \sigma_\lambda} \quad (5.2)$$

όπου  $K_p$  είναι το πλήθος των στοιχείων της κλάσης  $p$ ,  $x_{vk}$  και  $x_{\lambda k}$  οι τιμές των χαρακτηριστικών,  $\mu_v, \mu_\lambda, \sigma_v$  και  $\sigma_\lambda$  οι μέσες τιμές και οι τυπικές αποκλίσεις των τιμών των χαρακτηριστικών στην τάξη  $p$ . Ο αριθμητής  $\sigma_{v\lambda}$  του κλάσματος λέγεται *ετεροσυσχέτιση* (*cross-correlation*) των τιμών των δύο χαρακτηριστικών. Ο παράγοντας συσχέτισης μετρά την ομοιότητα μεταξύ των δύο χαρακτηριστικών και παίρνει τιμές μεταξύ  $-1$  και  $+1$ . Μια τιμή κοντά στο  $+1$  ή στο  $-1$  σημαίνει ότι τα δύο χαρακτηριστικά συσχετίζονται πολύ ή συσχετίζονται αντίστροφα, αντίστοιχα. Μια τιμή κοντά στο μηδέν δείχνει ότι τα χαρακτηριστικά είναι κατά πολύ ασυσχέτιστα. Όταν δύο χαρακτηριστικά έχουν μεγάλη συσχέτιση μπορούμε να απορρίψουμε ένα εκ των δύο.

#### ΠΑΡΑΔΕΙΓΜΑ

Δίνονται οι πίνακες τεσσάρων προτύπων  $\Pi_1, \Pi_2, \Pi_3, \Pi_4$

$$x_1^A = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, x_2^A = \begin{bmatrix} 1 \\ 0.5 \\ 1 \end{bmatrix}, x_3^B = \begin{bmatrix} 0.5 \\ -1 \\ 0 \end{bmatrix}, x_4^B = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$$

Απορρίψτε το χειρότερο χαρακτηριστικό με το κριτήριο της διαχωριστικότητας.

$$\mu^A = \begin{bmatrix} (0+1)/2 \\ (1+0.5)/2 \\ (0+1)/2 \end{bmatrix} = \begin{bmatrix} 1/2 \\ 3/4 \\ 1/2 \end{bmatrix}, \quad \mu^B = \begin{bmatrix} (0.5+1)/2 \\ (-1+0)/2 \\ (0+1)/2 \end{bmatrix} = \begin{bmatrix} 3/4 \\ -1/2 \\ 1/2 \end{bmatrix}$$

$$\begin{aligned}
(\sigma_1^A)^2 &= \frac{(0 - 1/2)^2 + (1 - 1/2)^2}{2} = 1/4 & (\sigma_2^A)^2 &= \frac{(1 - 3/4)^2 + (1/2 - 3/4)^2}{2} = 1/16 \\
(\sigma_3^A)^2 &= \frac{(0 - 1/2)^2 + (1 - 1/2)^2}{2} = 1/4 & (\sigma_1^B)^2 &= \frac{(1/2 - 3/4)^2 + (1 - 3/4)^2}{2} = 1/16 \\
(\sigma_2^B)^2 &= \frac{(-1 + 1/2)^2 + (0 - 1/2)^2}{2} = 1/4 & (\sigma_3^B)^2 &= \frac{(0 - 1/2)^2 + (1 - 1/2)^2}{2} = 1/4
\end{aligned}$$

$$S_{AB}^1 = \frac{|1/2 - 3/4|}{\sqrt{1/4 + 1/16}}, \quad S_{AB}^2 = \frac{|3/4 + 1/2|}{\sqrt{1/16 + 1/4}}, \quad S_{AB}^3 = \frac{|1/2 - 1/2|}{\sqrt{1/4 + 1/4}} = 0$$

Ο συντελεστής διαχωριστικότητας για το 3<sup>ο</sup> χαρακτηριστικό είναι μηδέν και γι' αυτό μπορεί να απορριφθεί.

## 5.2 Ανάλυση χαρακτηριστικών στην εκπαίδευση χωρίς επόπτη

Για την ανάλυση χαρακτηριστικών στην εκπαίδευση χωρίς επόπτη, θα περιγράψουμε τον μετασχηματισμό Karhunen-Loeve (**KLT**) ή μέθοδο *ανάλυσης κύριων συνιστωσών* (ΑΚΣ), (*PCA: Principal Components Analysis*). Η ΑΚΣ είναι ένα μέσο με το οποίο «γεννώνται» νέα χαρακτηριστικά από τα υπάρχοντα, ασυσχέτιστα μεταξύ τους. Με τον μετασχηματισμό είναι δυνατή η αντιπροσώπευση του αρχικού συνόλου των χαρακτηριστικών με άλλα λιγότερα και αποτελεσματικότερα που διατηρούν σημαντικό μέρος της πληροφορίας των δεδομένων. Με τη χρήση της ΑΚΣ επιτυγχάνεται η αύξηση της διασποράς των τιμών των χαρακτηριστικών, η μείωση των διαστάσεων και η συμπίεση της πληροφορίας.

Για ένα σύνολο  $K$  ανυσμάτων σε  $N$ -διαστάσεων χώρο χαρακτηριστικών έστω ότι  $\mathbf{x}_k$  είναι το διάνυσμα που δίνεται από την σχέση.

$$\mathbf{x}_k = [x_{1k} \quad x_{2k} \quad \dots \quad x_{nk} \quad \dots \quad x_{Nk}]^T \quad (5.2.1)$$

$$\text{Με μέσο όρο } \boldsymbol{\mu}_k = E[\mathbf{x}_k] = \frac{1}{K} \sum_{k=1}^K \mathbf{x}_k = \mathbf{0} \quad (5.2.2)$$

όπου  $E[.]$  η μαθηματική προσδοκία. Αν ο μέσος όρος  $\boldsymbol{\mu}_k$  δεν είναι το μηδενικό διάνυσμα τότε τροποποιούμε τις τιμές αφαιρώντας από αυτές τον μέσο όρο τους ώστε οι τιμές που θα προκύψουν να έχουν μηδενικό μέσο όρο. Ο πίνακας συνδιασποράς  $\mathbf{C}_k$  όλων των ανυσμάτων του συνόλου εκπαίδευσης δίνεται από τη σχέση (με την προϋπόθεση ότι ο μέσος όρος είναι μηδέν).

$$\mathbf{C}_x = E[\mathbf{x}\mathbf{x}^T] = \frac{1}{K} \sum_{k=1}^K \mathbf{x}_k \mathbf{x}_k^T = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1v} & \cdots & \sigma_{1N} \\ \sigma_{21} & \sigma_2^2 & \cdots & \sigma_{2v} & \cdots & \sigma_{2N} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \sigma_{\lambda 1} & \sigma_{\lambda 2} & \cdots & \sigma_{\lambda v} & \cdots & \sigma_{\lambda N} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \sigma_{N1} & \sigma_{N2} & \cdots & \sigma_{Nv} & \cdots & \sigma_N^2 \end{bmatrix} \quad (5.2.3)$$

Κάθε διαγώνια τιμή  $\sigma_v^2$  του πίνακα συνδιασποράς εκφράζει τη διασπορά των τιμών  $x_{vk}$  (άξονας  $v$ ), ενώ οι υπόλοιπες τη συνδιασπορά  $\sigma_{vk}$  μεταξύ δύο διαφορετικών μεταβλητών  $x_{\mu k}$  και  $x_{\nu k}$ . Ο πίνακας  $\mathbf{C}_x$  είναι συμμετρικός με θετικές τιμές. Είναι φανερό πως για τη διαδικασία της ταξινόμησης είναι επιθυμητές μεγάλες τιμές των διαγωνίων τιμών  $\sigma_v^2$  διότι μαρτυρούν ένα μεγάλο άπλωμα (spreading) των δεδομένων, ενώ είναι επιθυμητές μηδενικές τιμές για τις συνδιασπορές ώστε οι μεταβλητές να είναι μεταξύ τους ασυσχέτιστες. Για να το επιτύχουμε αυτό αναζητούμε ένα τέτοιο μετασχηματισμό  $\mathbf{W}$  που θα εκφράζεται από ένα πίνακα  $N \times N$  ώστε τα διανύσματα  $\mathbf{y} = \mathbf{W}^T \cdot \mathbf{x}$  να έχουν πίνακα συνδιασποράς  $\mathbf{C}_y$  με μηδενικές ετεροσυσχετίσεις. Θα ισχύει:

$$\mathbf{C}_y = E[\mathbf{y}\mathbf{y}^T] = E[\mathbf{W}^T \mathbf{x} (\mathbf{W}^T \mathbf{x})^T] = E[\mathbf{W}^T \mathbf{x} \mathbf{x}^T \mathbf{W}] = \mathbf{W}^T E[\mathbf{x} \mathbf{x}^T] \mathbf{W} = \mathbf{W}^T \mathbf{C}_x \mathbf{W} \quad (5.2.4)$$

Ο πίνακας  $\mathbf{C}_x$  είναι συμμετρικός θετικά ορισμένος και έχει ιδιοτιμές  $\lambda_1 > \lambda_2 > \dots > \lambda_N > 0$  με αντίστοιχα μοναδιαία ιδιοδιανύσματα  $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_v, \dots, \mathbf{e}_N$ . Αν  $\mathbf{A} =$

$$\begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_N \end{bmatrix} \quad \text{και} \quad \mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_v, \dots, \mathbf{q}_N] \quad \text{ισχύει:} \quad \mathbf{Q}^T \cdot \mathbf{Q} = \mathbf{I} \quad \text{και} \quad \mathbf{Q}^T \cdot \mathbf{C}_x \cdot \mathbf{Q} = \mathbf{A}. \quad (\Delta\epsilon\varsigma$$

Παράρτημα)

Αν ο πίνακας  $\mathbf{W} = \mathbf{Q}$  τότε από την (5.2.4) συνεπάγεται ότι  $\mathbf{C}_y = \mathbf{A}$ .

$$\mathbf{C}_y = \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_N \end{bmatrix} = \mathbf{A}, \quad \lambda_1 > \lambda_2 > \dots > \lambda_v > \dots > \lambda_N \quad (5.2.5)$$

Τα διανύσματα  $\mathbf{y}_k = |y_{1k} \ y_{2k} \ \dots \ y_{vk} \ \dots \ y_{Nk}|^T$  που προκύπτουν από την σχέση  $\mathbf{y} = \mathbf{Q}^T \cdot \mathbf{x}$  θα είναι  $\mathbf{y}_k = |\mathbf{x}_k^T \cdot \mathbf{q}_1, \ \mathbf{x}_k^T \cdot \mathbf{q}_2 \ \dots \ \mathbf{x}_k^T \cdot \mathbf{q}_v \ \dots \ \mathbf{x}_k^T \cdot \mathbf{q}_N|^T$ . Τα μοναδιαία ιδιοδιανύσματα  $\mathbf{q}_v$  του  $\mathbf{C}_x$  είναι κάθετα μεταξύ τους και αποτελούν ένα ορθογώνιο σύστημα συντεταγμένων οι άξονες του οποίου ονομάζονται *Κύριες Συνιστώσες (Principal Components)*. Κάθε μοναδιαίο ιδιοδιάνυσμα  $\mathbf{q}_v$  περιγράφει έναν άξονα επί του οποίου προβάλλεται κάθε διάνυσμα  $\mathbf{x}_k$ , η προβολή του είναι  $y_{vk} = \mathbf{x}_k^T \cdot \mathbf{q}_v$ . Οι τιμές των προβολών επί των κύριων αξόνων είναι αναλλοίωτες (αμετάβλητες) ως προς την περιστροφή των δεδομένων από κάποιο μετασχηματισμό στροφής.

Η τιμή  $\lambda_n = \sigma_n'^2$  είναι η διασπορά των τιμών  $y_{nk}$  κατά τον  $n$  άξονα. Το ιδιοδιάνυσμα  $\mathbf{q}_n$  είναι ένα "αποτελεσματικό" χαρακτηριστικό ενώ η τιμή  $\lambda_n$  αξιολογεί τη σπουδαιότητά του. Εκτός του ότι τα διανύσματα  $\mathbf{y}_k$  περιγράφουν τα πρότυπα με τιμές ασυσχέτιστες, με τον μετασχηματισμό μπορούμε να μειώσουμε και το πλήθος των τιμών που τα περιγράφουν αποδεχόμενοι μικρή απώλεια πληροφορίας (συμπύεση). Προς τούτο χρησιμοποιούμε τις προβολές  $y_n$  που αντιστοιχούν σε ιδιοδιανύσματα με τις μεγαλύτερες αντίστοιχες ιδιοτιμές, αποδεχόμενοι ένα μέσο τετραγωνικό σφάλμα  $\bar{\varepsilon}^2(M) = \sum_{n=M+1}^N \lambda_n$ , όπου  $M$  ο δείκτης του ιδιοδιανύσματος μετά από το οποίο διαγράφουμε τις προβολές. Για παράδειγμα αν  $\lambda_1 > \lambda_2 > \lambda_3 > \lambda_4 > \lambda_5$  οι ιδιοτιμές  $\mathbf{q}_1, \mathbf{q}_2, \mathbf{q}_3, \mathbf{q}_4, \mathbf{q}_5$  τα αντίστοιχα ιδιοδιανύσματα και  $\mathbf{y} = [y_1, y_2, y_3, y_4, y_5]^T$  διάνυσμα των προβολών που προέκυψαν από τον μετασχηματισμό  $\mathbf{y} = \mathbf{Q}^T \cdot \mathbf{x}$ , περιγράφουμε τα δεδομένα με διανύσματα  $\mathbf{y}' = [y_1, y_2, y_3]^T$  αποδεχόμενοι ένα μέσο τετραγωνικό σφάλμα  $\bar{\varepsilon}^2(3) = \sum_{n=3+1}^N \lambda_n = \lambda_4 + \lambda_5$ .

Ο KLT μας προσφέρει ασυσχέτιστα χαρακτηριστικά και δυνατότητα μείωσης του πλήθους τους. Μας προσφέρει επιπρόσθετα ανεξαρτησία ως προς την περιστροφή και την μετατόπιση διότι οι τιμές των προβολών επί των κύριων αξόνων είναι αναλλοίωτες (αμετάβλητες) ως προς την περιστροφή των δεδομένων από μετασχηματισμό στροφής και από μετατόπιση. Η ιδιότητα αυτή είναι εξαιρετικά χρήσιμη για δεδομένα μορφών σε δυαδικές ψηφιακές εικόνες (π.χ. χαρακτήρες, γράμματα).

#### ΠΑΡΑΔΕΙΓΜΑ

$$x_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, x_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, x_3 = \begin{bmatrix} 2 \\ 0 \end{bmatrix}, x_4 = \begin{bmatrix} 4 \\ 3 \end{bmatrix}, x_5 = \begin{bmatrix} 5 \\ 2 \end{bmatrix}$$

$$\mu = \begin{bmatrix} (1+1+2+4+5)/5 \\ (0+1+0+3+2)/5 \end{bmatrix} = \begin{bmatrix} 13/5 \\ 6/5 \end{bmatrix} = \begin{bmatrix} 2.6 \\ 1.2 \end{bmatrix}$$

$$\sigma_1^2 = \frac{(1-2.6)^2 + (1-2.6)^2 + (2-2.6)^2 + (4-2.6)^2 + (5-2.6)^2}{5} = 2.64$$

$$\sigma_2^2 = \frac{(0-1.2)^2 + (1-1.2)^2 + (0-1.2)^2 + (3-1.2)^2 + (2-1.2)^2}{5} = 1.36$$

$$\sigma_{21} = \frac{(0-1.2)(1-2.6) + (1-1.2)(1-2.6) + (0-1.2)(2-2.6) + (3-1.2)(4-2.6) + (2-1.2)(5-2.6)}{5} = 1.48$$

$$\sigma_{12} = \sigma_{21} = 1.48$$

$$C_x = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{bmatrix} = \begin{bmatrix} 2.64 & 1.48 \\ 1.48 & 1.36 \end{bmatrix}$$

$$\det(C_x - \lambda \cdot I) = \det \begin{bmatrix} 2.64 - \lambda & 1.48 \\ 1.48 & 1.36 - \lambda \end{bmatrix} = (2.64 - \lambda) \cdot (1.36 - \lambda) - 1.48^2 = 1.36 \cdot 2.64 - 1.36 \cdot \lambda - 2.64 \cdot \lambda + \lambda^2 - 1.48^2 = \lambda^2 - 4 \cdot \lambda - 1.4$$

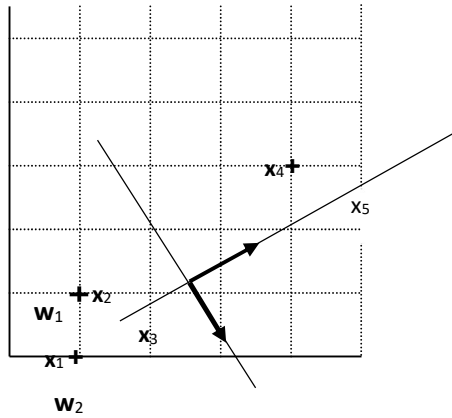
που είναι το χαρακτηριστικό πολυώνυμο.

Για  $\lambda^2 - 4 \cdot \lambda + 1.4 = 0$  προκύπτουν οι ρίζες  $\lambda_1 = 3.61$ ,  $\lambda_2 = 0.39$  που είναι οι ιδιοτιμές του  $C_x$ . Τα αντίστοιχα ιδιοδιανύσματα  $\mathbf{w}_1$ ,  $\mathbf{w}_2$  ικανοποιούν τις σχέσεις  $(C_x - \lambda_1 \cdot I) \cdot \mathbf{q}_1 = 0$  και  $(C_x - \lambda_2 \cdot I) \cdot \mathbf{q}_2 = 0$ .

Για το  $\mathbf{q}_1$  θα έχουμε αναλυτικά

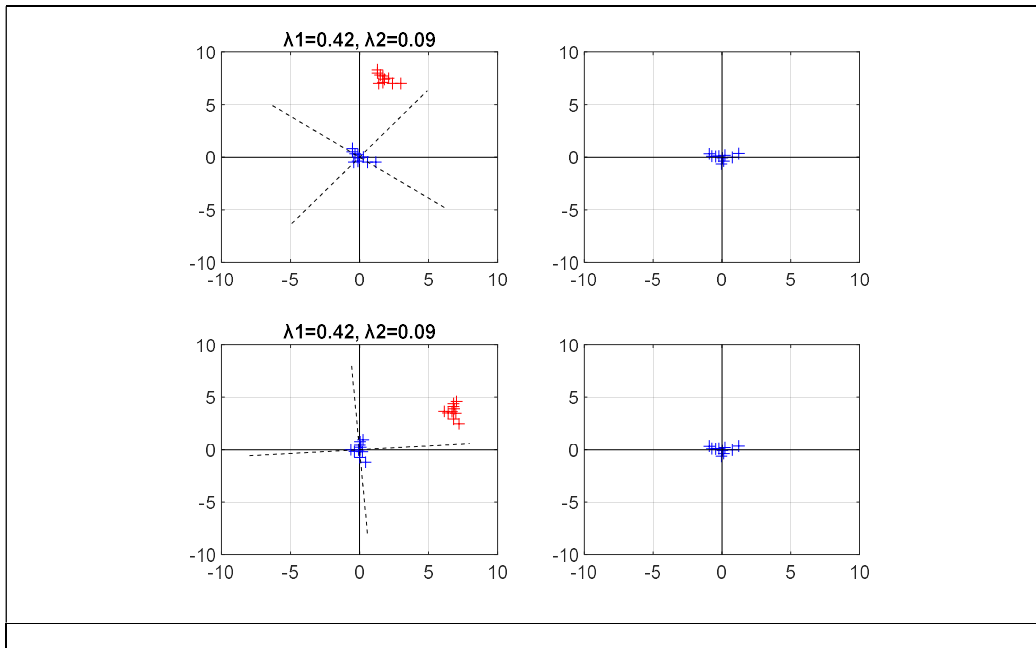
$$\begin{vmatrix} 2.64 - 3.61 & 1.48 \\ 1.48 & 1.36 - 3.61 \end{vmatrix} \cdot \begin{vmatrix} q_{11} \\ q_{21} \end{vmatrix} = 0 \Rightarrow \begin{cases} -0.97 \cdot q_{11} + 1.48 \cdot q_{21} = 0 \\ 1.48 \cdot q_{11} - 2.25 \cdot q_{21} = 0 \end{cases}$$

από το παραπάνω αόριστο σύστημα λύση για  $q_{11}=\alpha$  είναι  $\mathbf{q}_1=[\alpha, \alpha \cdot 0.65]^T$ . Ακόμη για να ισχύει  $\|\mathbf{q}_1\|=1 \Leftrightarrow \sqrt{\alpha^2 + (\alpha \cdot 1.65)^2} = 1 \Leftrightarrow \alpha = \pm 1/\sqrt{1+1.65^2} \approx .84$ , άρα ένα μοναδιαίο ιδιοδιάνυσμα είναι το  $\mathbf{q}_1=[0.84, 0.55]^T$ . Όμοια για την ιδιοτιμή  $\lambda_2$  προκύπτει  $\mathbf{q}_2=[0.55, -0.84]^T$ . Τα διανύσματα είναι κάθετα μεταξύ τους όπως αναμενόταν. Στο Σχ.5.2.1 φαίνονται γεωμετρικά τα αποτελέσματα της μεθόδου.



Σχήμα 5.2.1.

**ΠΑΡΑΔΕΙΓΜΑ ΑΝΕΞΑΡΤΗΣΙΑΣ ΑΠΟ ΣΤΡΟΦΗ**





### Ανάλυση Κύριων Συνιστωσών με Νευρωνικό Δίκτυο.

Ο KLT είναι μία αναλυτική διαδικασία για την εκτίμηση των ιδιοτιμών και των ιδιοδιανυσμάτων από τον πίνακα συνδιασποράς των δεδομένων  $\mathbf{C}_x$ . Ειδικά, εάν το πλήθος των διαστάσεων του χώρου των ανυσμάτων είναι μεγάλο ο υπολογισμός και ο χειρισμός του πίνακα  $\mathbf{C}_x$  είναι ανέφικτος.

Για τον παραπάνω λόγο αντί του KLT χρησιμοποιούμε στο πρώτο μέρος του ταξινομητή ένα γραμμικό νευρωνικό δίκτυο ενός επιπέδου (single-layer linear feed forward neural net), το οποίο επιτελεί την ΑΚΣ (Σχ. 5.2.1). Το δίκτυο αυτό εκπαιδεύεται χωρίς επόπτη (unsupervised) από τον γενικευμένο αλγόριθμο του Hebb που βασίζεται στον κανόνα εκμάθησης του Hebb. Με τον αλγόριθμο αυτό το νευρωνικό δίκτυο συγκλίνει με πιθανότητα ένα στα ιδιοδιανύσματα ενός πίνακα συνδιασποράς  $\mathbf{C}_x$ . Ο υπολογισμός του  $\mathbf{C}_x$  δεν είναι αναγκαίος επειδή τα ιδιοδιανύσματα προκύπτουν κατευθείαν από τα δεδομένα. Η εκπαίδευση του νευρωνικού δικτύου γίνεται ως εξής:

Έστω ότι

- $t=1,2,3,\dots$  μία μεταβλητή για τη μέτρηση της επανάληψης της διαδικασίας,
- $\mathbf{x}(t)$  το διάνυσμα εισόδου στο νευρωνικό δίκτυο τη χρονική στιγμή  $t$ , με συνιστώσες  $x_\nu(t)$ ,  $\nu=0,1,\dots,N$ ,
- $K$  το πλήθος των νευρώνων ( $K \leq N$ ) που ως εκ τούτου είναι και το πλήθος των επιθυμητών κύριων συνιστωσών,
- $k$  ένας δείκτης που αποδίδεται στους νευρώνες  $k=1,\dots,K$ ,
- $w_{k\nu}(t)$  η τιμή του βάρους της σύναψης που συνδέει τον  $k$  νευρώνα με την  $\nu$  είσοδο κατά την επανάληψη  $t$ ,
- $\mathbf{W}_{K \times N}(t)$  ο πίνακας των  $w_{k\nu}(t)$  και
- $\mathbf{y}(t)$  το άνυσμα εξόδου του δικτύου κατά την επανάληψη  $t$ , με συνιστώσες  $y_k(t)$ .

Εκτελούνται τα ακόλουθα βήματα υπολογισμών:

Βήμα 1. Αφαιρείται πρώτα το διάνυσμα  $\boldsymbol{\mu}_x$  της μέσης τιμής από κάθε στοιχείο του συνόλου εκπαίδευσης. Με τον τρόπο αυτό οι τιμές που προκύπτουν έχουν μηδενικό διάνυσμα μέσης τιμής.

Βήμα 2. Αποδίδονται αρχικά ( $t=0$ ) στα βάρη  $w_{k\nu}(0)$  των συνάψεων μικρές τυχαίες τιμές και στην παράμετρο  $\gamma$  του ρυθμού εκπαίδευσης μικρή θετική τιμή (π.χ.  $\gamma=0.007$ ).

Βήμα 3. Υπολογίζεται το διάνυσμα εξόδου  $\mathbf{y}(t)$  με συνιστώσες  $y_k(t)$  για  $\nu=1,\dots,N$ ,  $k=1,\dots,K$  από τη σχέση

$$y_k(t) = \sum_{\nu=1}^N w_{k\nu}(t) \cdot x_\nu(t), \quad \text{ή} \quad \mathbf{y}(t) = \mathbf{W}(t) \mathbf{x}(t) \quad (63)$$

και η ποσότητα  $\Delta w_{k\nu}(t)$  από την σχέση

$$\Delta w_{kv}(t) = \gamma \left( \underbrace{y_k(t) \cdot x_v(t)}_{(\alpha)} - \underbrace{y_k(t) \cdot \sum_{\lambda=1}^k w_{\lambda v}(t) \cdot y_{\lambda}(t)}_{(\beta)} \right) \quad (64)$$

Ο όρος (α) εκφράζει τον απλό κανόνα του Hebb, που λέει ότι εάν δύο νευρώνες που βρίσκονται στα άκρα μιας σύναψης (σύνδεσης), ενεργοποιούνται συγχρόνως, τότε η ισχύς της σύναψης αυξάνεται. Αλλιώς εξασθενεί ή εκφυλίζεται. Ο όρος (β) επιβάλλει ένα όριο στην αύξηση της σύναψης. Η τιμή του  $w_{kv}(t+1)$  προσαρμόζεται σύμφωνα με τη σχέση

$$W_{kv}(t+1) = w_{kv}(t) + \Delta w_{kv} \quad (5.2.7)$$

Βήμα 4. Αύξηση της μεταβλητής  $t$  κατά ένα και επανάληψη από το βήμα 3 έως ότου τα βάρη των συνάψεων  $w_{kv}$  φθάσουν στη σταθερή κατάσταση.

Μετά τη φάση της εκπαίδευσης τα βάρη  $w_{kv}$  του  $k$  νευρώνα συγκλίνουν στην  $n$  συνιστώσα του ιδιοδιανύσματος που αντιστοιχεί στην  $k$  ιδιοτιμή του πίνακα συνδιασποράς  $\mathbf{C}_x$ . Με άλλα λόγια οι γραμμές του πίνακα  $\mathbf{W}$  έχουν προσεγγίσει τα πρώτα  $K$  ιδιοδιανύσματα του  $\mathbf{C}_x$ , ταξινομημένες σε φθίνουσα σειρά.

### 3.7 Ταξινομητές με βάση τον κανόνα πιθανοτήτων του Bayes

Στις προσεγγίσεις που παρουσιάσαμε στα προηγούμενα κεφάλαια δεν λάβαμε υπόψη την πιθανότητα εμφάνισης των προτύπων του συνόλου εκπαίδευσης. Στο κεφάλαιο αυτό θα ασχοληθούμε με ταξινομητές που βασίζονται σ αυτή την πιθανότητα και στην ελαχιστοποίηση του σφάλματος ταξινόμησης που προκύπτει από αυτήν.

Θα ξεκινήσουμε την παρουσίαση με το πρόβλημα των δύο κλάσεων  $C_1, C_2$ . Στόχος μας αρχικά είναι ο υπολογισμός της πιθανότητας το πρότυπο να ανήκει στη κλάση  $C_1$  (ή  $C_2$ ) δεδομένου ότι έχει την τιμή  $\mathbf{x}$ , δηλαδή να βρούμε τις υπό συνθήκη πιθανότητες  $P(C_1|\mathbf{x})$  και  $P(C_2|\mathbf{x})$ . Αν  $P(C_1|\mathbf{x}) > P(C_2|\mathbf{x})$  το πρότυπο ανήκει στη κλάση  $C_1$ , διαφορετικά στη  $C_2$ . Ακολούθως αναζητούμε εκείνες της περιοχές του χώρου των προτύπων για τις οποίες ελαχιστοποιείται το σφάλμα ταξινόμησης με βάση την παραπάνω θεώρηση. Από το κανόνα του Bayes ισχύει ότι

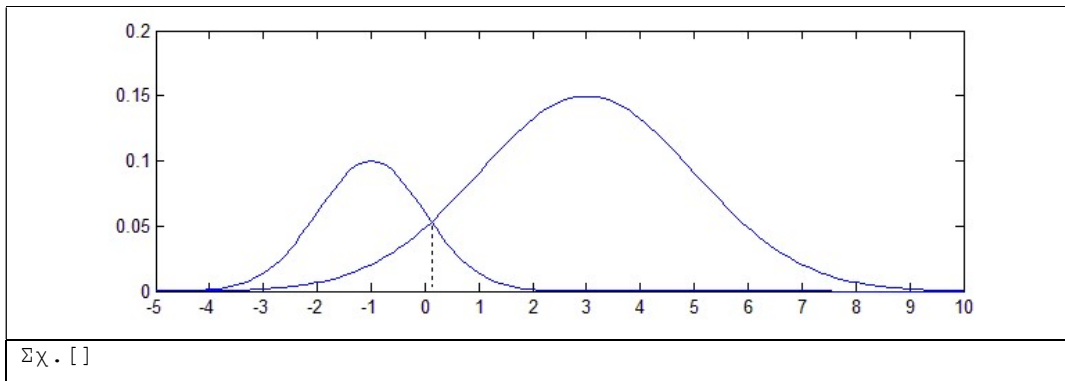
$$P(C_1|\mathbf{x}) = \frac{p(\mathbf{x}|C_1) \cdot P(C_1)}{p(\mathbf{x})} \quad \text{και} \quad P(C_2|\mathbf{x}) = \frac{p(\mathbf{x}|C_2) \cdot P(C_2)}{p(\mathbf{x})} \quad [3.7.1]$$

Όπου  $P(C_1)$  και  $P(C_2)$  οι πιθανότητες το πρότυπο  $\mathbf{x}$  να ανήκει στις κλάσεις  $C_1, C_2$  αντίστοιχα και  $p(\mathbf{x}|C_1), p(\mathbf{x}|C_2)$  συναρτήσεις πυκνότητας πιθανότητας (σ.π.π.) για την γενική περίπτωση που το  $\mathbf{x}$  είναι διανυσματική συνεχής τυχαία μεταβλητή. Για τις σ.π.π. θα χρησιμοποιούμε το μικρό  $p$  και για τις πιθανότητες κεφαλαίο  $P$ .

Άρα το  $\mathbf{x}$  ανήκει

στη  $C_1$  αν  $p(\mathbf{x}|C_1) \cdot P(C_1) > p(\mathbf{x}|C_2) \cdot P(C_2)$ , στη  $C_2$  αν  $p(\mathbf{x}|C_1) \cdot P(C_1) < p(\mathbf{x}|C_2) \cdot P(C_2)$   
 [3.7.2]

Στα προηγούμενα σημειώνεται ότι οι πιθανότητες  $P(C_1)$  και  $P(C_2)$  υπολογίζονται προσεγγιστικά από τις σχέσεις  $P(C_1) = N_1/N$ ,  $P(C_2) = N_2/N$  με  $N_1, N_2$  το πλήθος των προτύπων, οι ποσότητες  $p(\mathbf{x}|C_1)$ ,  $p(\mathbf{x}|C_2)$  θεωρούνται γνωστές ή εκτιμώνται από το σύνολο εκπαίδευσης. Από την σχέση [3.7.2] βρίσκουμε ακολούθως τις περιοχές  $R_1$  και  $R_2$  που διαχωρίζουν τον χώρο των προτύπων. Για παράδειγμα για δύο κλάσεις με  $P(C_1) = 0.25$  και  $P(C_2) = 0.75$ ,  $p(\mathbf{x}|C_1)$ ,  $p(\mathbf{x}|C_2)$  κανονικές κατανομές με μέσους όρους  $\mu_1 = -1$ ,  $\mu_2 = 3$  και διασπορές  $\sigma_1 = 1$ ,  $\sigma_2 = 2$  οι περιοχές  $R_1$  και  $R_2$  θα είναι όπως στο Σχ. \_



Παρακάτω θα δείξουμε ότι με βάση τα παραπάνω το λάθος ταξινόμησης ελαχιστοποιείται. Αν  $R_1$  η περιοχή που περιλαμβάνει τις τιμές του  $\mathbf{x}$  που ταξινομούνται στη κλάση  $C_1$  και  $R_2$  η περιοχή που περιλαμβάνει τις τιμές του  $\mathbf{x}$  που ταξινομούνται στη κλάση  $C_2$  το σφάλμα θα συμβαίνει αν το  $\mathbf{x}$  ανήκει στη  $R_1$  και το πρότυπο στη κλάση  $C_2$  ή αν το  $\mathbf{x}$  ανήκει στη  $R_2$  και το πρότυπο στη κλάση  $C_1$  και η πιθανότητα  $P_e$  να συμβεί θα είναι

$$\begin{aligned}
 P_e &= P(\mathbf{x} \in R_2, C_1) + P(\mathbf{x} \in R_1, C_2) = \\
 &P(C_1) \cdot P(\mathbf{x} \in R_2|C_1) + P(C_2) \cdot P(\mathbf{x} \in R_1|C_2) = \\
 &P(C_1) \cdot \int_{R_2} p(\mathbf{x}|C_1) d\mathbf{x} + P(C_2) \cdot \int_{R_1} p(\mathbf{x}|C_2) d\mathbf{x} = \\
 &\int_{R_2} p(C_1|\mathbf{x}) \cdot p(\mathbf{x}) d\mathbf{x} + \int_{R_1} p(C_2|\mathbf{x}) \cdot p(\mathbf{x}) d\mathbf{x} \quad (3.7.3)
 \end{aligned}$$

Ισχύει ακόμη ότι

$$\int_{-\infty}^{+\infty} p(\mathbf{x}|C_1) d\mathbf{x} = 1 \Rightarrow \int_{R_1} p(\mathbf{x}|C_1) d\mathbf{x} + \int_{R_2} p(\mathbf{x}|C_1) d\mathbf{x} = 1 \Rightarrow$$

$$\frac{1}{P(C_1)} \int_{R_1} p(C_1|\mathbf{x}) \cdot p(\mathbf{x}) d\mathbf{x} + \frac{1}{P(C_1)} \int_{R_2} p(C_1|\mathbf{x}) \cdot p(\mathbf{x}) d\mathbf{x} = 1 \Rightarrow$$

$$\int_{R_1} p(C_1|\mathbf{x}) \cdot p(\mathbf{x}) d\mathbf{x} + \int_{R_2} p(C_1|\mathbf{x}) \cdot p(\mathbf{x}) d\mathbf{x} = P(C_1) \Rightarrow$$

$$\int_{R_2} p(C_1|\mathbf{x}) \cdot p(\mathbf{x}) d\mathbf{x} = P(C_1) - \int_{R_1} p(C_1|\mathbf{x}) \cdot p(\mathbf{x}) d\mathbf{x} \quad (3.7.4)$$

Από τις σχέσεις (3.7.1) και (3.7.2) προκύπτει

$$P_e = P(C_1) - \int_{R_1} (P(C_1|\mathbf{x}) - P(C_2|\mathbf{x})) p(\mathbf{x}) d\mathbf{x}$$

Που σημαίνει ότι το σφάλμα ελαχιστοποιείται όταν η περιοχή  $R_1$  είναι τέτοια ώστε το ολοκλήρωμα να είναι θετικό δηλαδή  $P(C_1|\mathbf{x}) > P(C_2|\mathbf{x})$ . Το αυτό αντίστοιχα ισχύει και για την  $R_2$ .

Στην περίπτωση των πολλών κλάσεων το πρότυπο ανήκει στην κλάση  $C_i$  για την οποία ισχύει  $P(C_i|\mathbf{x}) > P(C_j|\mathbf{x})$  με  $i \neq j$

#### Συναρτήσεις διάκρισης και διαχωριστικές επιφάνειες

Στην περίπτωση που απόφαση λαμβάνεται με κριτήριο την ελαχιστοποίηση του σφάλματος ταξινόμησης για δύο κλάσεις η σχέση μπορούμε σύμφωνα με τα προηγούμενα να ισχυριστούμε ότι το πρότυπο ανήκει στην κλάση  $K_1$  αν  $P(C_1|\mathbf{x}) - P(C_2|\mathbf{x}) > 0$  διαφορετικά ανήκει στην κλάση  $K_2$ . Θα παρουσιάσουμε ακολούθως την δυνατότητα ταξινόμησης του προτύπου μέσω του προσδιορισμού μιας διαχωριστικής επιφάνειας του χώρου που θα ελαχιστοποιεί το σφάλμα ταξινόμησης στην περίπτωση που η συνάρτηση πυκνότητας πιθανότητας (σ.π.π., probability density function, pdf) είναι η κανονική κατανομή (Gaussian).

Για μονοδιάστατη τυχαία μεταβλητή η σ.π.π. της κανονική μεταβλητής δίνεται από τον τύπο

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Όπου  $\mu$  και  $\sigma$  η μέση τιμή και η διασπορά της κατανομής.

Στην γενική περίπτωση που η τυχαία μεταβλητή είναι διάνυσμα  $\mathbf{x}$  ενός  $N$  διάστατου χώρου αντίστοιχη σ.π.π. δίδεται από τη σχέση

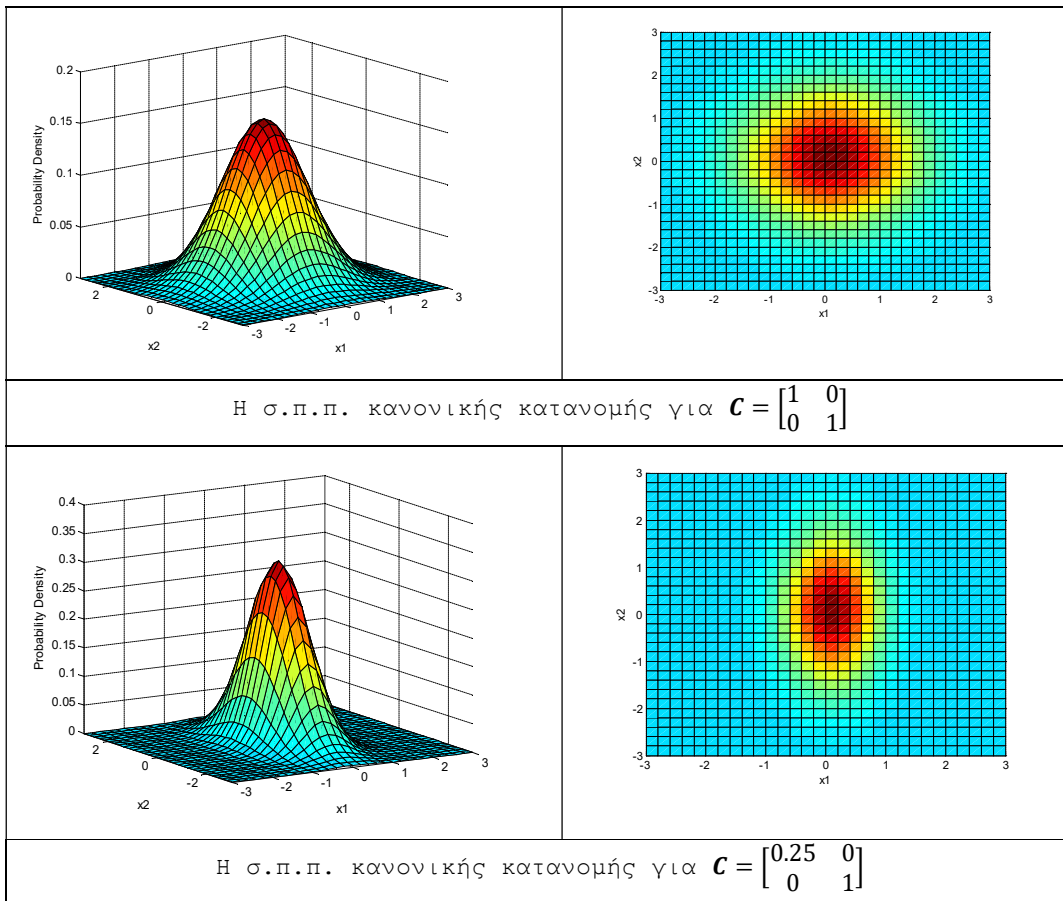
$$p(\mathbf{x}) = \frac{1}{(2\pi)^{N/2} |\mathbf{C}|^{1/2}} e^{-\frac{(\mathbf{x}-\boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x}-\boldsymbol{\mu})}{2}}$$

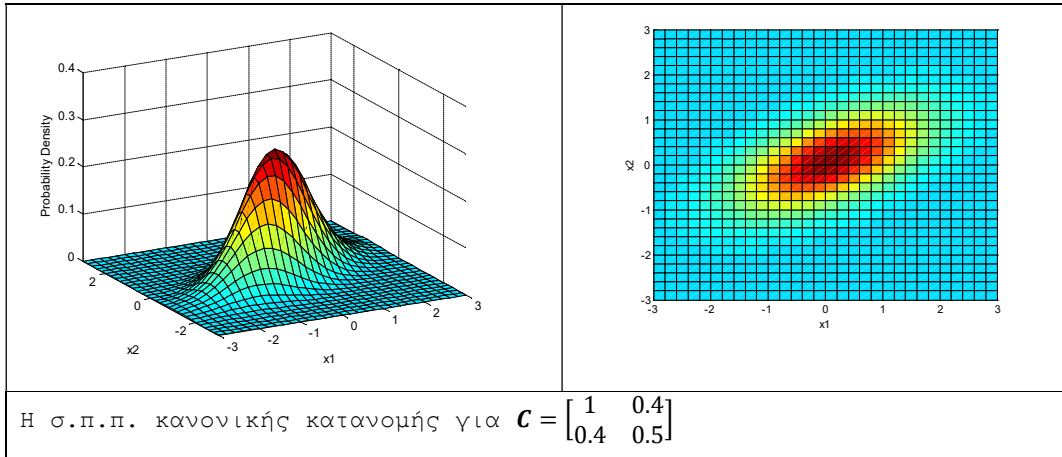
Όπου  $\boldsymbol{\mu}$  μέση τιμή και  $\mathbf{C}$  ο πίνακας συνδιασποράς της σ.π.π. και  $|\mathbf{C}|$  η ορίζουσά του. Προσοχή η χρήση του συμβόλου  $\mathbf{C}$  δεν αναφέρεται πλέον στην κλάση K αλλά στο πίνακα συνδιασποράς.

Για λόγους γεωμετρικής απεικόνισης θα θεωρήσουμε την περίπτωση που η τυχαία μεταβλητή ανήκει στον δισδιάστατο χώρο. Ο πίνακας συνδιασποράς αποτελείται από τις διασπορές και τις ετεροσυσχετίσεις των γραμμών της διανυσματικής μεταβλητής  $\mathbf{x}=[x_1, x_2]^T$  και είναι της μορφής

$$\mathbf{C} = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{bmatrix}$$

Στα ακόλουθα σχήματα βλέπουμε τις σ.π.π κανονικών κατανομών για διαφορετικούς πίνακες συνδιασποράς και μέσους όρους.





Λόγω της εκθετικής μορφής της σ.π.π. της εκθετικής κατανομής μπορούμε να χρησιμοποιήσουμε για απλούστευση των υπολογισμών μία μονότονη λογαριθμική συνάρτηση, δηλαδή η πρόταση

το  $\mathbf{x}$  ανήκει στη  $K_1$  αν  $p(\mathbf{x}|K_1) \cdot P(K_1) > p(\mathbf{x}|K_2) \cdot P(K_2)$ , στη  $K_2$  αν  $p(\mathbf{x}|K_1) \cdot P(K_1) < p(\mathbf{x}|K_2) \cdot P(K_2)$  διατυπώνεται το  $\mathbf{x}$  ανήκει

στη  $K_1$  αν  $g_1(\mathbf{x}) > g_2(\mathbf{x})$ , στη  $K_2$  αν  $g_1(\mathbf{x}) < g_2(\mathbf{x})$

με

$$g_i(\mathbf{x}) = \ln(p(\mathbf{x}|K_i) \cdot P(K_i)) = \ln(p(\mathbf{x}|K_i)) + \ln(P(K_i)) \Rightarrow$$

$$g_i(\mathbf{x}) = -\frac{(\mathbf{x}-\boldsymbol{\mu}_i)^T \mathbf{C}_i^{-1} (\mathbf{x}-\boldsymbol{\mu}_i)}{2} + \ln(P(K_i)) + c_i \Rightarrow$$

$$g_i(\mathbf{x}) = -\frac{1}{2} \mathbf{x}^T \mathbf{C}_i^{-1} \mathbf{x} + \frac{1}{2} \boldsymbol{\mu}_i^T (\mathbf{C}_i^{-1} + (\mathbf{C}_i^{-1})^T) \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_i^T \mathbf{C}_i^{-1} \boldsymbol{\mu}_i + \ln(P(K_i)) + c_i$$

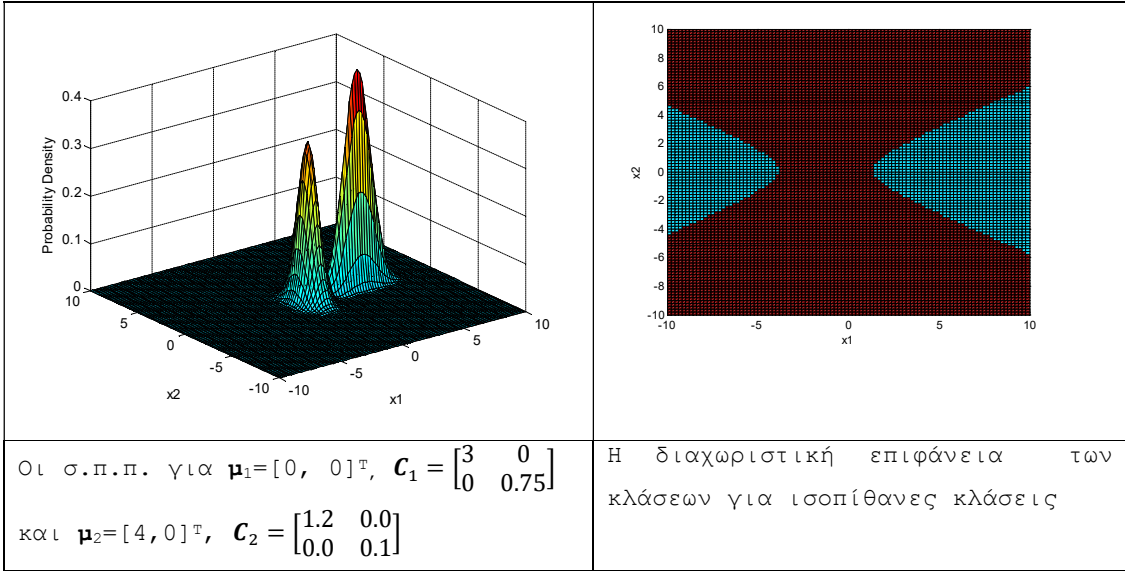
Και  $c_i = -\frac{N}{2} \ln(2\pi) - \frac{1}{2} \ln(|\mathbf{C}_i|)$  σταθερά ποσότητα

Γενικά αυτή είναι μια μη γραμμική τετραγωνική μορφή. Τα σημεία του χώρου για τα οποία

$g_i(\mathbf{x}) = \min\{g_m(\mathbf{x})\}, \forall m \neq i$  καθορίζουν την περιοχή της κλάσης  $K_i$ . Στην περίπτωση των δύο κλάσεων η διαχωριστική τους καμπύλη ικανοποιεί την σχέση

$$d(\mathbf{x}) = g_1(\mathbf{x}) - g_2(\mathbf{x}) = 0$$

Για  $g_1(\mathbf{x})$  και  $g_2(\mathbf{x})$  συναρτήσεις απόφασης δύο κλάσεων Ακολουθούν παραδείγματα για διαφορετικές σ.π.π.



## ΠΑΡΑΡΤΗΜΑ Α: ΒΑΣΙΚΕΣ ΠΡΑΞΕΙΣ ΓΡΑΜΜΙΚΗΣ ΑΛΓΕΒΡΑΣ

Ο ανάστροφος πίνακας  $\mathbf{A}^T$  ενός πίνακα  $\mathbf{A}$  είναι ένας πίνακας που έχει ως γραμμές τις στήλες του  $\mathbf{A}$  με ίδιο δείκτη. Αν ο  $\mathbf{A}_{M \times N}$  έχει  $M$  γραμμές και  $N$  στήλες, τότε ο  $\mathbf{A}^T_{N \times M}$  θα έχει  $N$  γραμμές και  $M$  στήλες. Η  $\nu$  γραμμή του  $\mathbf{A}^T$  είναι η  $\nu$  στήλη του  $\mathbf{A}$  και η  $\mu$  γραμμή του  $\eta \mu$  στήλη του  $\mathbf{A}$ . Για παράδειγμα αν

$$\mathbf{A}_{3 \times 2} = \begin{vmatrix} 1 & 0.5 \\ -1 & 3 \\ 2 & -2 \end{vmatrix} \Leftrightarrow \mathbf{A}^T_{2 \times 3} = \begin{vmatrix} 1 & -1 & 2 \\ 0.5 & 3 & -2 \end{vmatrix}.$$

Πολλαπλασιασμός πινάκων. Εστώ οι πίνακες  $\mathbf{A}_{N \times K}$  και  $\mathbf{B}_{K \times M}$ , το γινόμενο τους είναι ένας πίνακας  $\mathbf{\Gamma}_{N \times M}$  για τον οποίο γράφουμε  $\mathbf{\Gamma} = \mathbf{A} \cdot \mathbf{B}$  και έχει στοιχεία  $\Gamma_{\nu\mu}$  που δίνονται από την σχέση

$$\Gamma_{\nu\mu} = \sum_{\kappa=1}^K A_{\nu\kappa} \cdot B_{\kappa\mu}$$

Για παράδειγμα αν

$$\mathbf{A}_{3 \times 2} = \begin{vmatrix} 1 & 2 \\ 0.5 & 4 \\ 3 & 1 \end{vmatrix} \text{ και } \mathbf{B}_{2 \times 3} = \begin{vmatrix} 3 & -1 & 1 \\ 1 & 0.5 & -0.5 \end{vmatrix}$$

$$\mathbf{A} \cdot \mathbf{B} = \mathbf{\Gamma}_{3 \times 3} = \begin{vmatrix} 1 \cdot 3 + 2 \cdot 1 & 1 \cdot (-1) + 2 \cdot 0.5 & 1 \cdot 1 + 2 \cdot (-0.5) \\ 0.5 \cdot 3 + 4 \cdot 1 & 0.5 \cdot (-1) + 4 \cdot 0.5 & 0.5 \cdot 1 + 4 \cdot (-0.5) \\ 3 \cdot 3 + 1 \cdot 1 & 3 \cdot (-1) + 1 \cdot 0.5 & 3 \cdot 1 + 1 \cdot (-0.5) \end{vmatrix} = \begin{vmatrix} 5 & 0 & 0 \\ 7.5 & 1.5 & -1.5 \\ 10 & -2.5 & 2.5 \end{vmatrix}$$

## ΠΑΡΑΡΤΗΜΑ Γ: ΣΤΑΤΙΣΤΙΚΑ ΧΑΡΑΚΤΗΡΙΣΤΙΚΑ ΤΟΥ ΧΩΡΟΥ ΤΩΝ ΠΡΟΤΥΠΩΝ.

Για ένα πλήθος  $K$  προτύπων που περιγράφονται σε ένα χώρο με  $N$  διαστάσεις η μέση τιμή  $\mu$  είναι ένας πίνακας-στήλης που ορίζεται από τη σχέση

$$\mu = \frac{1}{K} \sum_{\kappa=1}^K x_{\kappa} = [\mu_1, \mu_2, \dots, \mu_N]^T \text{ όπου } \mu_{\nu} = \frac{1}{K} \sum_{\kappa=1}^K x_{\nu\kappa} \quad (2.3.1)$$

Η μέση τιμή  $\mu$  αναφέρεται και ως η μαθηματική προσδοκία  $E(x)$  όπου  $x$  μεταβλητή για τους πίνακες-στήλης των προτύπων. Αν υπάρχουν  $\lambda$  πρότυπα που οι τιμές των πινάκων-στήλης τους είναι ίσες, λαμβάνονται υπόψη  $\lambda$  φορές στο άθροισμα της (9) και ως εκ τούτου δεν χρησιμοποιείται στον τύπο (9) η συνάρτηση συχνότητας εμφάνισης  $f(x)$  κάθε τιμής της μεταβλητής  $x$ . Οι τιμές  $\mu_{\nu}$  όπου  $\nu=1 \dots N$  ορίζουν το μέσο διάνυσμα  $\bar{\mu}$  των  $K$  προτύπων στο  $N$ -διάστατο



χώρο τους. Ο πίνακας συμμεταβλητότητας ή πίνακας συνδιασποράς  $C_{ov}$  ορίζεται από τη σχέση

$$C_{ov} = \frac{1}{K} \sum_{k=1}^K (x_k - \mu) \cdot (x_k - \mu)^T = \frac{1}{K} \sum_{k=1}^K \begin{bmatrix} x_{1k} - \mu_1 \\ x_{2k} - \mu_2 \\ \vdots \\ x_{\nu k} - \mu_\nu \\ \vdots \\ x_{Nk} - \mu_N \end{bmatrix} \cdot [x_{1k} - \mu_1, x_{2k} - \mu_2, \dots, x_{\nu k} - \mu_\nu, \dots, x_{Nk} - \mu_N] =$$

$$\frac{1}{K} \sum_{k=1}^K \begin{bmatrix} (x_{1k} - \mu_1)^2 & (x_{1k} - \mu_1) \cdot (x_{2k} - \mu_2) & \dots & (x_{1k} - \mu_1) \cdot (x_{\nu k} - \mu_\nu) & \dots & (x_{1k} - \mu_1) \cdot (x_{Nk} - \mu_N) \\ (x_{2k} - \mu_2) \cdot (x_{1k} - \mu_1) & (x_{2k} - \mu_2)^2 & \dots & (x_{2k} - \mu_2) \cdot (x_{\nu k} - \mu_\nu) & \dots & (x_{2k} - \mu_2) \cdot (x_{Nk} - \mu_N) \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ (x_{\lambda k} - \mu_\lambda) \cdot (x_{1k} - \mu_1) & (x_{\lambda k} - \mu_\lambda) \cdot (x_{2k} - \mu_2) & \dots & (x_{\lambda k} - \mu_\lambda) \cdot (x_{\nu k} - \mu_\nu) & \dots & (x_{\lambda k} - \mu_\lambda) \cdot (x_{Nk} - \mu_N) \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ (x_{Nk} - \mu_N) \cdot (x_{1k} - \mu_1) & (x_{Nk} - \mu_N) \cdot (x_{2k} - \mu_2) & \dots & (x_{Nk} - \mu_N) \cdot (x_{\nu k} - \mu_\nu) & \dots & (x_{Nk} - \mu_N)^2 \end{bmatrix}$$

$$= \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1\nu} & \dots & \sigma_{1N} \\ \sigma_{21} & \sigma_2^2 & \dots & \sigma_{2\nu} & \dots & \sigma_{2N} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \sigma_{\lambda 1} & \sigma_{\lambda 2} & \dots & \sigma_{\lambda \nu} & \dots & \sigma_{\lambda N} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \sigma_{N1} & \sigma_{N2} & \dots & \sigma_{N\nu} & \dots & \sigma_N^2 \end{bmatrix} = E( (x-\mu) (x-\mu)^T ) \quad (2.3.2)$$

όπου  $\sigma_{\lambda\nu} = \frac{1}{K} \sum_{k=1}^K (x_{\lambda k} - \mu_\lambda) \cdot (x_{\nu k} - \mu_\nu)$

Η ποσότητα  $\sigma_{\lambda\nu}^2$  είναι η μεταβλητότητα ή διασπορά (variance) των μετρήσεων του ν-οστού χαρακτηριστικού των προτύπων. Η ποσότητα  $\sigma_{\lambda\nu} = \sigma_{\nu\lambda}$  λέγεται συμμεταβλητότητα ή συνδιασπορά των τιμών των λ και ν χαρακτηριστικών των προτύπων και χρησιμοποιείται για τον υπολογισμό του συντελεστή συσχέτισης  $R_{\lambda\nu}$  δύο χαρακτηριστικών μ,ν σύμφωνα με τη σχέση

$$R_{\lambda\nu} = \frac{\sigma_{\lambda\nu}}{\sigma_\lambda \cdot \sigma_\nu} \quad (2.3.3)$$

Η ποσότητα  $\sigma_\nu = \sqrt{\sigma_\nu^2}$  λέγεται τυπική απόκλιση. Οι παραπάνω ποσότητες χρησιμοποιούνται για την αξιολόγηση και ανάλυση των μετρούμενων χαρακτηριστικών.

## ΠΑΡΑΡΤΗΜΑ Γ:ΤΑ ΙΔΙΟΔΙΑΝΥΣΜΑΤΑ ΤΟΥ ΠΙΝΑΚΑ ΣΥΝΔΙΑΣΠΟΡΑΣ

Σε ένα συμμετρικό πίνακα  $C_{N \times N}$ ,  $C^T = C$ , με ιδιοτιμές  $\lambda_\mu, \lambda_\nu, \lambda_\mu \neq \lambda_\nu$  και αντίστοιχα μοναδιαία ιδιοδιανύσματα  $w_\mu \neq w_\nu$  ισχύει ότι  $w_\mu^T \cdot w_\nu = 0$ . Δηλαδή τα ιδιοδιανύσματα είναι κάθετα μεταξύ τους.

Απόδειξη: Επειδή  $\lambda_\mu, \lambda_\nu$  ιδιοτιμές και  $\mathbf{w}_\mu, \mathbf{w}_\nu$  ιδιοδιανύσματα ισχύουν οι σχέσεις

$$\left. \begin{aligned} \mathbf{C} \cdot \mathbf{w}_\mu &= \lambda_\mu \cdot \mathbf{w}_\mu \\ \mathbf{C} \cdot \mathbf{w}_\nu &= \lambda_\nu \cdot \mathbf{w}_\nu \end{aligned} \right\} \Rightarrow \begin{aligned} \mathbf{w}_\nu^T \cdot \mathbf{C} \cdot \mathbf{w}_\mu &= \mathbf{w}_\nu^T \cdot \lambda_\mu \cdot \mathbf{w}_\mu = \lambda_\mu \cdot \mathbf{w}_\nu^T \cdot \mathbf{w}_\mu \quad (A) \\ \mathbf{w}_\mu^T \cdot \mathbf{C} \cdot \mathbf{w}_\nu &= \mathbf{w}_\mu^T \cdot \lambda_\nu \cdot \mathbf{w}_\nu = \lambda_\nu \cdot \mathbf{w}_\mu^T \cdot \mathbf{w}_\nu \quad (B) \end{aligned}$$

$$\mathbf{w}_\nu^T \cdot \mathbf{C} \cdot \mathbf{w}_\mu = (\mathbf{w}_\mu^T \cdot \mathbf{C} \cdot \mathbf{w}_\nu)^T \quad (\Gamma)$$

Υπενθυμίζεται ότι  $(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$ .

Από τις σχέσεις (A), (B), (Γ) συνεπάγεται

$$\lambda_\mu \cdot \mathbf{w}_\nu^T \cdot \mathbf{w}_\mu = \lambda_\nu \cdot \mathbf{w}_\mu^T \cdot \mathbf{w}_\nu \Leftrightarrow (\lambda_\mu - \lambda_\nu) \mathbf{w}_\nu^T \cdot \mathbf{w}_\mu = 0$$

επειδή  $\lambda_\mu \neq \lambda_\nu$  έπεται ότι  $\mathbf{w}_\nu^T \cdot \mathbf{w}_\mu = \mathbf{w}_\mu^T \cdot \mathbf{w}_\nu = 0$ .

Από τα προηγούμενα συνεπάγεται ότι για τον πίνακα  $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N]$  ισχύει ότι

$$\mathbf{W}^T \cdot \mathbf{W} = \begin{bmatrix} \mathbf{w}_1^T \\ \mathbf{w}_2^T \\ \vdots \\ \mathbf{w}_N^T \end{bmatrix} [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N] = \begin{bmatrix} \mathbf{w}_1^T \cdot \mathbf{w}_1 & 0 & \dots & 0 \\ 0 & \mathbf{w}_2^T \cdot \mathbf{w}_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \mathbf{w}_N^T \cdot \mathbf{w}_N \end{bmatrix} = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \dots & 0 & \dots & 1 \end{bmatrix} = \mathbf{I}$$

Άρα αντίστροφος του  $\mathbf{W}$  είναι ο  $\mathbf{W}^{-1} = \mathbf{W}^T$ .

Για τα ιδιοδιανύσματα του  $\mathbf{C}$  ισχύει εξ' ορισμού ότι

$$\left. \begin{aligned} \mathbf{C} \cdot \mathbf{w}_1 &= \lambda_1 \cdot \mathbf{w}_1 \\ \mathbf{C} \cdot \mathbf{w}_2 &= \lambda_2 \cdot \mathbf{w}_2 \\ \vdots \\ \mathbf{C} \cdot \mathbf{w}_N &= \lambda_N \cdot \mathbf{w}_N \end{aligned} \right\} \Rightarrow \mathbf{C} \cdot [\mathbf{w}_1, \dots, \mathbf{w}_N] = [\mathbf{w}_1, \dots, \mathbf{w}_N] \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_N \end{bmatrix} = \mathbf{W} \cdot \mathbf{A}$$

όπου  $\mathbf{A} = \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_N \end{bmatrix}$

$$\text{Αφού } \mathbf{C} \cdot \mathbf{W} = \mathbf{W} \cdot \mathbf{A} \Rightarrow \mathbf{W}^{-1} \cdot \mathbf{C} \cdot \mathbf{W} = \mathbf{W}^{-1} \cdot \mathbf{W} \cdot \mathbf{A} = \mathbf{I} \cdot \mathbf{A} = \mathbf{A} \Leftrightarrow \mathbf{W}^T \cdot \mathbf{C} \cdot \mathbf{W} = \mathbf{A}.$$

## ΠΑΡΑΡΤΗΜΑ Δ: ΕΝΤΟΛΕΣ MATLAB & TOOLBOX