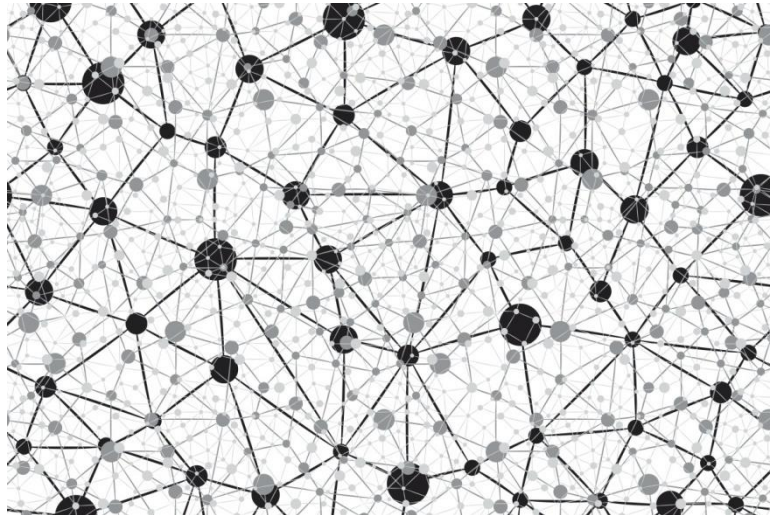




ΕΘΝΙΚΟ ΜΕΤΣΟΒΙΟ ΠΟΛΥΤΕΧΝΕΙΟ

ΣΧΟΛΗ ΕΦΑΡΜΟΣΜΕΝΩΝ ΜΑΘΗΜΑΤΙΚΩΝ
ΚΑΙ ΦΥΣΙΚΩΝ ΕΠΙΣΤΗΜΩΝ

*Εισαγωγικά θέματα αριθμητικής ανάλυσης,
για νευρωνικά δίκτυα*



Διπλωματική εργασία

Επιμέλεια: Μπεκιάρης Ιωάννης

Αριθμός μητρώου: 09113559

Επιβλέπων Καθηγητής: Χρυσάφινος Κωνσταντίνος

Τριμελής Εξεταστική επιτροπή

Καθηγητής Ε.Μ.Π
.....

κ. Χρυσάφινος Κωνσταντίνος

Καθηγητής Ε.Μ.Π
.....

κ. Γεωργούλης Εμμανουήλ

Καθηγητής Ε.Μ.Π.
.....

κ. Κολέτσος Ιωάννης

Περιεχόμενα

1.	Εισαγωγή στα νευρωνικά δίκτυα.....	7
1.1	Ιστορική αναδρομή.....	7
1.2	Τι είναι νευρωνικό δίκτυο	9
1.3	Perceptrons και πολύ-επίπεδα perceptrons.....	14
1.3.1	Νευρωνικά δίκτυα ενός επιπέδου με προς τα εμπρός τροφοδότηση σήματος (perceptrons).....	14
1.3.2.	Πολύ-επίπεδα perceptron	16
2.	Πυκνότητα και προσέγγιση νευρωνικών δικτύων.....	19
2.1	Άμεση προσέγγιση της πυκνότητας	23
2.2	k-Σιγμοειδής συνάρτηση	28
2.3	Δυϊκός χώρος και συνέλιξη μεθόδων της προσέγγισης.....	29
3.	Αριθμητική ανάλυση σε αλγορίθμους μάθησης.....	45
3.1	Αλγόριθμος Delta rule	46
3.2	Η μέθοδος epoch	52
3.3	Γενίκευση σε μη- γραμμικά συστήματα.....	53
4	Αριθμητικές εφαρμογές των νευρωνικών δικτύων	56
4.1	Δίκτυα Hopfield και βελτιστοποίηση γράφων	57
4.2	Το πρόβλημα του πλανόδιου πωλητή (TSP)	59
	Πίνακας εικόνων	64
	Βιβλιογραφία.....	65

Πρόλογος

Σκοπός της παρούσας εργασίας αποτελεί η μελέτη εισαγωγικών θεμάτων της αριθμητικής ανάλυσης, σε έννοιες οι οποίες βασίζονται στα νευρωνικά δίκτυα, σύμφωνα με την επιστημονική δημοσίευση του S.W.Ellacott με τίτλο “Aspects of the numerical analysis of neural networks”. Πιο συγκεκριμένα, ο αναγνώστης μέσα από τα πρώτα κεφάλαια της εργασίας, εισάγεται ομαλά στην έννοια των νευρωνικών δικτύων και στις μαθηματικές προσεγγίσεις που θα ακολουθήσουν στα επόμενα κεφάλαια. Η παρακάτω εργασία χωρίζεται σε 4 κεφάλαια, τα οποία επικεντρώνονται και εμβαθύνουν στις έννοιες των νευρωνικών δικτύων, την πυκνότητα και τη θεωρία προσέγγισης, τους αλγόριθμους μάθησης καθώς και σε ορισμένες εφαρμογές της αριθμητικής ανάλυσης στα συγκεκριμένα δίκτυα, ακολουθώντας τη συλλογική πορεία του συγγραφέα.

Μολονότι οι έρευνες των τεχνητών νευρωνικών δικτύων ξεκίνησαν με σκοπό την κατασκευή ενός μοντέλου που προσομοιώνει το βιολογικό δίκτυο νευρώνων, στο παρόν οι μελέτες αυτές δεν περιορίζονται στα βιολογικά δίκτυα καθώς βρίσκουν εφαρμογή σε ποικίλες περιοχές, όπως για παράδειγμα σε μεθόδους επεξεργασίας δεδομένων, κατηγοριοποίησης δεδομένων, αναγνώριση μοτίβων και άλλα.

Ευχαριστίες

Θα ήθελα να ευχαριστήσω θερμά και κύριο Χρυσafίνο Κωνσταντίνο για την πολύτιμη βοήθεια του στην υλοποίηση της διπλωματικής του εργασίας. Η καθοδήγηση που μου παρείχε με βοήθησε να αναπτύξω τις γνώσεις μου και να κατανοήσω σύνθετες μαθηματικές έννοιες βασισμένοι στην αριθμητική ανάλυση στα νευρωνικά δίκτυα.

Επίσης, θα ήθελα να ευχαριστήσω τα μέλη της εξεταστικής επιτροπής, τον κύριο Γεωργούλη Εμμανουήλ και τον κύριο Κολέτσο Ιωάννη, οι οποίοι διάβασαν και ασχολήθηκαν με την διπλωματική μου εργασία.

1. Εισαγωγή στα νευρωνικά δίκτυα

1.1 Ιστορική αναδρομή

Η περιοχή των νευρωνικών δικτύων, αποτελεί ένα μαθηματικό εργαλείο δίχως ιδιαίτερα μεγάλη προϋστορία της οποίας η αρχή χρονολογείται στα μέσα του 20^{ου} αιώνα με μία κύρια ώθηση την δεκαετία του 1980.

Το πρώτο μοντέλο νευρωνικού δικτύου, παρουσιάζεται από τους McCulloch και Pitts το 1943, οι οποίοι παρουσίασαν, για πρώτη φορά, τα στοιχεία ενός νευρωνικού δικτύου. Πιο συγκεκριμένα, υποστήριξαν πως ένα νευρωνικό δίκτυο αποτελείται από μια συλλογή ενός μεγάλου αριθμού νευρώνων και έδειξαν ότι οι νευρώνες μπορούν να λειτουργήσουν άμεσα με τις διασυνδέσεις τους. Αργότερα, οι ίδιοι προχώρησαν σε ένα πιο εξελιγμένο πρότυπο και πολύπλοκο δίκτυο με άμεση εφαρμογή στην αναγνώριση σχημάτων. Σε αυτό το σημείο είναι ιδιαίτερα σημαντικό να τονιστεί, πως η θεμελιώδης ιδέα των νευρωνικών δικτύων προέρχεται από την μελέτη λειτουργίας του εγκεφάλου των θηλαστικών καθώς και στην προσπάθεια κατανόησης της ακολουθίας των βημάτων που ακολουθούν τα βιολογικά συστήματα, με σκοπό την εξαγωγή/ πρόβλεψη συμπεριφοράς/αποτελεσμάτων. Επιπροσθέτως, αξίζει να αναφερθεί ότι οι εργασίες των McCulloch–Pitts πιθανότατα να μην είχαν λάβει την αρμόζουσα προσοχή εάν δεν τις χρησιμοποιούσε ο J. Von Neumann ως παράδειγμα για τις υπολογιστικές μηχανές. Αργότερα και συγκεκριμένα το 1957, παρουσιάστηκε για πρώτη φορά από τον Rosenblatt, το μοντέλο του αισθητήρα (perceptron) το οποίο πρόκειται για ένα μοντέλο με δύο μόνο επίπεδα της εισόδου και εξόδου. Ειδικότερα, το μοντέλο αυτό στην αρχή ενθουσίασε την επιστημονική κοινότητα καθώς φαίνεται πως παρουσίαζε μεγάλη ακρίβεια στην επίλυση προβλημάτων. Η ιδέα των νευρωνικών δικτύων, χαρακτηρίστηκε ως η αποτελεσματικότερη τεχνική για την επίλυση προβλημάτων τα οποία έως τότε παρέμεναν άλυτα, γρήγορα όμως αναγνώρισαν ότι τα μοντέλα αυτά έχουν αρκετούς περιορισμούς.

Πολλοί επιστήμονες αποφάσισαν να ασχοληθούν αποκλειστικά με την ιδέα των νευρωνικών δικτύων, παρουσιάζοντας πολλές αναφορές, εργασίες και πρακτικές εφαρμογές, ονόματα όπως οι Widrow και Hoff (1959) των οποίων τα μοντέλα τους χρησιμοποιήθηκαν ως φίλτρα προκειμένου να εξαλειφθεί η ηχώ στις τηλεφωνικές γραμμές. Ο J. Hopfield (1982) με την εργασία του απέδειξε πως τα νευρωνικά δίκτυα μπορούν να χρησιμοποιηθούν και ως αποθηκευτικός χώρος με την δυνατότητα ανάκτησης πληροφορίας. Επιπροσθέτως, εργασίες όπως των McClelland και Rumelhart (1986) παραθέτουν μια διαφορετική εκδοχή των perceptron, ενσωματώνοντας στο υπάρχον μοντέλο και άλλα επίπεδα νευρώνων (εκτός από την είσοδο και την έξοδο) και προτείνοντας την μέθοδο οπισθοδιάδοσης (back-propagation) η οποία παραμένει έως σήμερα η

αποτελεσματικότερη τεχνική εκπαίδευσης δικτύων. Τέλος, επισημαίνεται η εργασία του S.W. Ellacot (1994) που αποτελεί μια απόπειρα σύνδεσης των νευρωνικών δικτύων με τον αυστηρό μαθηματικό φορμαλισμό και ιδιαίτερα με την ερευνητική περιοχή της αριθμητικής ανάλυσης. (βλ. **[1, Κεφάλαιο2, σελ. 1-6]**)

1.2 Τι είναι νευρωνικό δίκτυο

Η πρόοδος στα τεχνητά νευρωνικά δίκτυα, που συχνά αναφέρονται απλώς και ως νευρωνικά δίκτυα, βασίζεται στην θεμελιώδη αναγνώριση, ότι ο εγκέφαλος επεξεργάζεται πληροφορίες με ένα τελείως διαφορετικό τρόπο από τον συμβατικό, με τον οποίο επεξεργάζονται πληροφορίες οι σύγχρονοι ψηφιακοί υπολογιστές.

Ο εγκέφαλος είναι ένα εξαιρετικά περίπλοκο, μη γραμμικό και παράλληλο σύστημα διαδικασιών. Διαθέτει την ικανότητα της οργάνωσης των νευρώνων καθώς και την πραγματοποίηση συγκεκριμένων υπολογισμών, όπως την αναγνώριση προσώπων, αντίληψη και τον έλεγχο μηχανών πολύ πιο γρήγορα από τους σύγχρονους υπολογιστές. Για παράδειγμα, η ανθρώπινη όραση, η οποία πρόκειται για μια διεργασία επεξεργασίας δεδομένων, μια συνάρτηση δηλαδή του οπτικού μας συστήματος με σκοπό την παρουσίαση και την παροχή πληροφοριών του περιβάλλοντος γύρω του υποκειμένου, ώστε να πετύχει την αλληλεπίδραση του με αυτό. Το παραπάνω λειτουργικό σύστημα για να πετύχει την αναγνώριση ενός γνώριμου προσώπου σε ένα ασυνήθιστο τοπίο χρειάζεται περίπου 100 – 200 ms (0,1 – 0,2 s), πολύ πιο γρήγορα από τους συμβατικούς υπολογιστές. (Churchland, 1986)

Αναγνωρίζοντας τον εξαιρετικό τρόπο λειτουργίας του εγκεφάλου, θα λέγαμε ότι ένα νευρωνικό δίκτυο είναι μια μηχανή σχεδιασμένη έτσι ώστε να μοντελοποιεί διεργασίες με τον τρόπο του εγκεφάλου. Για να επιτευχθεί καλή απόδοση, το νευρωνικό δίκτυο χρησιμοποιεί μια ενδοσύνδεση από απλά υπολογιστικά κελιά που καλούνται νευρώνες ή μονάδες επεξεργασίας.

Σύμφωνα με τα παραπάνω, κατά τους (Alexander και Morton (1990)) θα μπορούσαμε να ορίσουμε το νευρωνικό δίκτυο ως:

Ένα νευρωνικό δίκτυο είναι ένας τεράστιος, παράλληλα κατανεμημένος επεξεργαστής που έχει μια φυσική κλίση στην αποθήκευση βιωματικών γνώσεων, κάνοντάς τις διαθέσιμες, για χρήση αργότερα. Το νευρωνικό δίκτυο ομοιάζει με τον ανθρώπινο εγκέφαλο σε δύο τρόπους:

1. Η γνώση του νευρωνικού δικτύου επιτυγχάνεται με μια διαδικασία μάθησης.
2. Οι συνδέσεις των ενδονευρώνων, γνωστές και ως βάρη συνάψεων χρησιμοποιούνται με σκοπό την αποθήκευση της πληροφορίας.

Η ενέργεια που εκτελείται για τη διαδικασία μάθησης, καλείται αλγόριθμος μάθησης. Πρόκειται δηλαδή για μια συνάρτηση που τροποποιεί τα βάρη των

συνάψεων του δικτύου με μεθοδικό τρόπο, ώστε να επιτευχθεί το επιθυμητό αποτέλεσμα.

Γιατί όμως να επιλέξουμε τα νευρωνικά δίκτυα ώστε να προσεγγίσουμε το πρόβλημα μας;

Η απάντηση στο παραπάνω ερώτημα, έπεται στο φαινόμενο ότι ένα νευρωνικό δίκτυο εκμεταλλεύεται την ισχύ ενός ηλεκτρονικού συστήματος για τη χρήση της τεράστιας παράλληλης κατανεμημένης δομής του, καθώς και την ικανότητα της μάθησης και κατ' επέκταση της γενίκευσης των αποτελεσμάτων. Αυτές οι δύο ιδιότητες του δικτύου το καθιστούν ικανό για την επίλυση πολύπλοκων προβλημάτων. Παρόλα αυτά, στην πράξη τα παραπάνω πραγματοποιούνται σε συνεργασία με τον μελετητή και του δικτύου που έχει οριστεί.

Πιο συγκεκριμένα, ένα νευρωνικό δίκτυο κατέχει τις παρακάτω ιδιότητες και ικανότητες:

1) Μη γραμμικότητα (Nonlinearity):

Ο νευρώνας πρόκειται για ένα μη γραμμικό μηχανισμό και κατ' επέκταση, το νευρωνικό δίκτυο που αποτελείται από ένα σύνολο ενδοσυνδεδεμένων νευρώνων μορφοποιείται ως ένα μη γραμμικό σύστημα. Η μη γραμμικότητα του δικτύου πρόκειται για μια σημαντική ιδιότητα αφού πολλά από τα φυσικά προβλήματα που πρέπει να αντιμετωπιστούν είναι εκ φύσεως μη γραμμικά, όπως για παράδειγμα η γενίκευση ενός σήματος (ομιλητικό σήμα, κ.α.).

2) Σχεδιασμός των εισαγόμενων – εξαγόμενων, δεδομένων:

Ένα συνηθισμένο μοντέλο μάθησης, ονομάζεται supervised learning και περιλαμβάνει τροποποιήσεις των βαρών του δικτύου, εφαρμόζοντας μια συλλογή από καταγεγραμμένα δείγματα και παραδείγματα. Κάθε παράδειγμα εμπεριέχει ένα μοναδικό εισαγόμενο σήμα και την κατάλληλη επιθυμητή έξοδο. Το δίκτυο επιλέγει ένα τυχαίο παράδειγμα από το δείγμα και τα βάρη των συνάψεων μετασχηματίζουν κατάλληλα το δίκτυο με σκοπό την ελαχιστοποίηση της απόκλισης ανάμεσα στην επιθυμητή έξοδο και την πραγματική αντίδραση του δικτύου που έχει παραχθεί από το εισαγόμενο δεδομένο σε συμφωνία με ένα κατάλληλο στατιστικό, ή άλλο κριτήριο. Η διαδικασία εκμάθησης επαναλαμβάνεται για πολλά ορίσματα του δείγματος, ούτως ώστε να φτάσει σε μια σταθερή κατάσταση, όπου δε θα σημειώνονται δραματικές αλλαγές στα βάρη των συνάψεων. Έτσι το δίκτυο μας μαθαίνει από ορίσματα με βάση το σχεδιασμό των εισαγόμενων – εξαγόμενων δεδομένων.

3) Προσαρμοστικότητα:

Τα νευρωνικά δίκτυα έχουν την ικανότητα του να προσαρμόζουν τα βάρη των συνάψεων σε αλλαγές του περιβάλλοντος. Ειδικότερα, ένα δίκτυο που είναι σχεδιασμένο να λειτουργεί σε ένα συγκεκριμένο περιβάλλον, μπορεί εύκολα να επανεκπαιδευτεί έτσι ώστε να ενεργεί σε ένα άλλο με μικρές αλλαγές συγκριτικά με τις αρχικές συνθήκες.

4) *Αποδεικτική απόκριση:*

Στο πλαίσιο της ταξινόμησης των μοτίβων, ένα νευρωνικό δίκτυο μπορεί να σχεδιαστεί έτσι ώστε να παρέχει πληροφορίες όχι μόνο για την επιλογή του κατάλληλου μοτίβου αλλά επίσης και για την βεβαιότητα της επιλογής. Αυτή η πληροφορία μπορεί να βοηθήσει στην κατασκευή μοτίβων με σκοπό τη βελτίωση της απόδοσης και της ταξινόμησης από το δίκτυο.

5) *Συνάφεια του συστήματος:*

Η πληροφορία, παρουσιάζεται από την κατασκευή και την κατάσταση ενεργοποίησης του νευρωνικού δικτύου, οπότε κάθε νευρώνας ενδεχομένως να είναι επηρεασμένος από τη συνολική δραστηριότητα όλου του δικτύου. Συνεπώς, η συνάφεια αντιμετωπίζεται με φυσικότητα από ένα νευρωνικό δίκτυο.

6) *Ανοχή λάθους:*

Ένα νευρωνικό δίκτυο, έχει την δυνατότητα να είναι εγγενώς ανεκτικό σε σφάλματα, με την έννοια ότι η απόδοση του μπορεί να υποβιβαστεί κάτω από ορισμένες συνθήκες λειτουργίας. Δηλαδή, εάν ένας νευρώνας ή οι ενώσεις του έχουν υποστεί ζημιά, μπορεί να καταργηθεί ή να ανακαλεστεί ένα αποθηκευμένο μοτίβο.

7) *Ομοιομορφία της ανάλυσης και σχεδιασμός:*

Το νευρωνικό δίκτυο εξ' ορισμού, είναι ένα σύστημα επεξεργασίας δεδομένων με σκοπό την καθολικότητα. Πιο συγκεκριμένα, θα μπορούσαμε να πούμε ότι οι νευρώνες είναι κύριο συστατικό για όλα τα νευρωνικά δίκτυα. Η ταύτιση αυτή, επιτρέπει στα δίκτυα να μοιράζονται θεωρίες και αλγορίθμους μάθησης σε διαφορετικούς κλάδους της επιστήμης. Τέλος, δίκτυα που θέλουμε να κατασκευάσουμε μπορούν να επιτευχθούν από την επεξεργασία ήδη υπαρχόντων. **(βλ. [11, Κεφάλαιο1, σελ. 1-5])**

Παρακάτω παρουσιάζεται ένα απλό παράδειγμα όπου το δίκτυο είναι με τέτοιο τρόπο σχεδιασμένο, ώστε να υπολογίζει την “exclusive-or” (XOR) συνάρτηση. Ουσιαστικά, πρόκειται για ένα κλασικό πρόβλημα στα τεχνητά νευρωνικά δίκτυα με σκοπό την εξαγωγή της τιμής 0 αν τα εισαγόμενα δεδομένα είναι ίδια και την τιμή 1, αλλιώς.

Παράδειγμα1: Υπολογισμός “XOR”

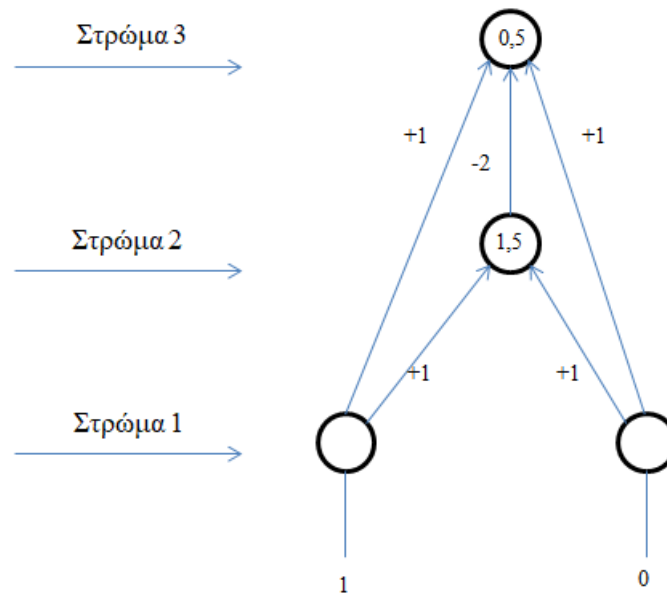
Έστω το σετ διανυσμάτων X που περιλαμβάνει τους συνδυασμούς μεταξύ των τιμών 0 και 1, καθώς και τα βάρη συνάψεων w_{ij} , όπου ο δείκτης i εκφράζει τη θέση του εισαγόμενου δεδομένου και j τη θέση του στρώματος που καταλήγει, όπως φαίνεται παρακάτω .

$$X = \{(0,0), (0,1), (1,0), (1,1)\}.$$

Το δίκτυο πρέπει να εξάγει την τιμή “0”, εάν τα εισαγόμενα δεδομένα είναι ίδια και “1” αλλιώς.

Έστω εισαγόμενο διάνυσμα το (1,0) και τα βάρη συνάψεων $w_{11} = 1$, $w_{21} = 1$, $w_{12} = 1$, $w_{22} = 1$, $w_{23} = -2$.

Το διάνυσμα με το οποίο ασχολούμαστε θα πολλαπλασιαστεί με τα αντίστοιχα βάρη των συνάψεων και θα συγκριθεί με την οριακή τιμή 1,5 (threshold value). Τα στοιχεία αυτά φαίνονται στο παρακάτω σχήμα.



Σχήμα 1

Πιο συγκεκριμένα, η ποσότητα $(1 \times 1) + (0 \times 1) = 1 < 1,5$ και άρα το δίκτυο θα εξάγει τη τιμή 0. Διαφορετικά αν είχαμε ως input το διάνυσμα τιμών (1,1) η πράξη $(1 \times 1) + (1 \times 1) = 2 > 1,5$ τότε το output που θα παίρναμε θα ήταν η τιμή 1.

Εν συνεχεία και αφού έχουμε προχωρήσει στο επόμενο κρυφό στρώμα, τα input και η ποσότητα του πρώτου κρυφού στρώματος πολλαπλασιαζόμενα εκ νέου με τα νέα βάρη πρέπει να συγκριθούν με την τελική οριακή τιμή 0,5 (όπως φαίνεται στο Σχήμα1). Πιο συγκεκριμένα, $(1 \times 1) + (0 \times 1) + (0 \times (-2)) = 1 > 0,5$ και έτσι το output του δικτύου είναι η τιμή 1, οδηγώντας μας στο επιθυμητό αποτέλεσμα. Να σημειωθεί ότι η εξαγόμενη τιμή της σύναψης του εισαγόμενου δεδομένου σε μορφή συνάρτησης, ονομάζεται συνάρτηση ενεργοποίησης της σύναψης.

Στον παραπάνω γράφο, τα σύμβολα που αναγνωρίζονται ως κύκλοι, καλούνται νευρώνες και οι ευθείες γραμμές που ενώνουν τη στήλη των εισαγόμενων δεδομένων καθώς και τα κρυφά στρώματα λέγονται ενώσεις ή συνάψεις. Ουσιαστικά, ο γράφος χρησιμοποιεί τα inputs πολλαπλασιάζοντας τα με τα ανάλογα βάρη των ενώσεων, καταλήγοντας τα στα κρυφά στρώματα. Έπειτα, συγκρίνοντας τις τιμές με τις threshold values καταλήγουμε στο στρώμα εξόδου με το προβλεπόμενο output.

Ορισμός 1.2.1: Έστω συνάρτηση $f: \mathbb{R} \rightarrow \mathbb{R}$ και x_1, x_2 που ανήκουν στο πεδίο ορισμού. Τότε, η συνάρτηση μας καλείται γνησίως αύξουσα αν $\forall x_1, x_2$ με $x_1 \geq x_2$ ισχύει ότι $f(x_1) \geq f(x_2)$.

Ορισμός 1.2.2: Έστω συνάρτηση $f: \mathbb{R} \rightarrow \mathbb{R}$. Η συνάρτηση θα καλείται φραγμένη αν υπάρχουν αριθμοί $A, B \in \mathbb{R}$ τέτοιοι ώστε $A \leq f(x) \leq B$.

Ορισμός 1.2.3: Οι σιγμοειδείς πρόκειται για συναρτήσεις όπου τα όρια τους στο $\pm \infty$ είναι 1 και 0 αντιστοίχως. Η δημοφιλέστερη επιλογή είναι αυτή της λογιστικής συνάρτησης με τον τύπο:

$$\sigma(x) = \frac{1}{1 + e^{-x}}.$$

Πέρα όμως από τη λογιστική συνάρτηση, παραδείγματα σιγμοειδών συναρτήσεων είναι η υπερβολική εφαπτομένη, η τόξο εφαπτομένη κ.α.

$$f(x) = \tanh x = \frac{e^x - e^{-x}}{e^x + e^{-x}}, f(x) = \arctan x.$$

Αξίζει να επισημανθεί σε αυτό το σημείο, ότι το παραπάνω δίκτυο χρησιμοποιείται ως ένα μοτίβο ταξινόμησης. Ειδικότερα, τα εισαγόμενα δεδομένα χωρίζονται σε κλάσεις (στην περίπτωση μας είχαμε τον διαχωρισμό σε δύο κλάσεις, σύμφωνα με το XOR όπου τα inputs ήταν τα 0,1) και έτσι μια απλή binary (0,1) έξοδος θα ήταν επαρκής. Βέβαια σε πολλές περιπτώσεις, οι προγραμματιστές μπορεί να απαιτήσουν περισσότερες από δύο κλάσεις, το οποίο μπορεί να επιτευχθεί επιτρέποντας το στρώμα εξόδου να πάρει τιμές πέρα από binary ή ακόμα να υπάρχει η επιλογή να χρησιμοποιηθούν περισσότερα από ένα output.

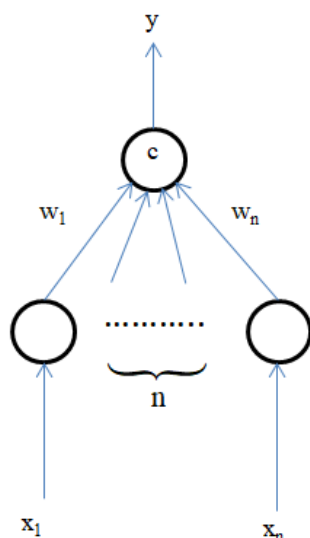
1.3 Perceptrons και πολύ-επίπεδα perceptrons

Η συγκεκριμένη ενότητα θα επικεντρωθεί στο απλούστερο και ένα από τα πιο αναγνωρίσιμα συστήματα μάθησης το perceptron (αισθητήρα) και μια επέκταση αυτού, τα πολύ-επίπεδα perceptrons (multi-layer perceptrons, MLP), (ή και ημι-γραμμικό δίκτυο με πρόσθια τροφοδότηση) που αναπτύχθηκαν από το Rosenblatt το 1957.

1.3.1 Νευρωνικά δίκτυα ενός επιπέδου με προς τα εμπρός τροφοδότηση σήματος (perceptrons)

Το perceptron πρόκειται για την απλούστερη μορφή ενός νευρωνικού δικτύου και χρησιμοποιείται κυρίως για την κατηγοριοποίηση μοτίβων που είναι γραμμικά διαχωρίσιμα. Συγκεκριμένα, αποτελείται από έναν απλό νευρώνα με κατάλληλα βάρη συνάψεων και μια οριακή τιμή. Ο αλγόριθμος που χρησιμοποιήθηκε για τη ρύθμιση των εισαγόμενων δεδομένων του δικτύου, πρωτοεμφανίστηκε σε μια διαδικασία μάθησης κατά τον Rosenblatt (1958-1962). Ο Rosenblatt, απέδειξε ότι αν το μοτίβο (διανύσματα) που χρησιμοποιούνται κατά την εκπαίδευση του αλγορίθμου είναι σχεδιασμένα σε δύο γραμμικά διαχωρίσιμες κλάσεις, τότε ο αλγόριθμος perceptron συγκλίνει σε πεπερασμένο αριθμό βημάτων σχηματίζοντας ένα γραμμικό διαχωρίσιμο σετ δεδομένων. Η απόδειξη της σύγκλισης του αλγορίθμου είναι γνωστή και ως θεώρημα σύγκλισης του perceptron (perceptron convergence theorem).

Στο παρακάτω σχήμα παρουσιάζεται ένα απλό perceptron με μία έξοδο (Σχήμα 2).



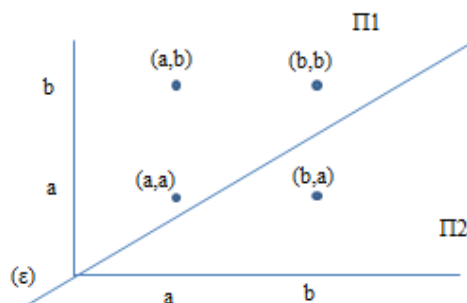
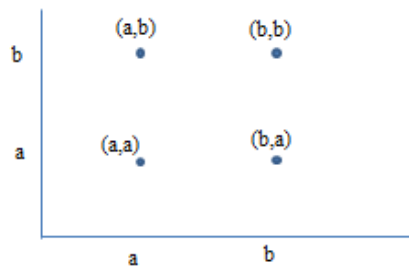
Σχήμα 2

Στο (Σχήμα 2) οι μονάδες περιγράφονται ακριβώς όπως το (σχήμα 1) με παρόμοια χαρακτηριστικά (τα inputs έχουν μια συνάρτηση ενεργοποίησης και το output περνάει από μία οριακή τιμή).

Έτσι, συμπεραίνει κανείς σύμφωνα με τα παραπάνω: Έστω x το διάνυσμα εισόδου, w το διάνυσμα των βαρών ($x, w \in \mathbb{R}^n$) και c η οριακή τιμή. Τότε αν το εσωτερικό γινόμενο $w^T \cdot x \leq c$, η έξοδος y που θα πάρουμε θα είναι 1 καθώς και αν η πράξη $w^T \cdot x \geq c$ τότε το output y θα είναι 0.

Παρατήρηση 1.3.1: Συμπερασματικά θα μπορούσε να παρατηρηθεί πως για ένα δεδομένο διάνυσμα βαρών και μια οριακή τιμή που δεν αλλάζει το δίκτυο διαιρεί με ένα υπερπίπεδο τον \mathbb{R}^n σε δύο χώρους. Συνεπώς, ο αλγόριθμος αυτός έχει την ιδιότητα της δυαδικής (binary) γραμμικής ταξινόμησης του χώρου.

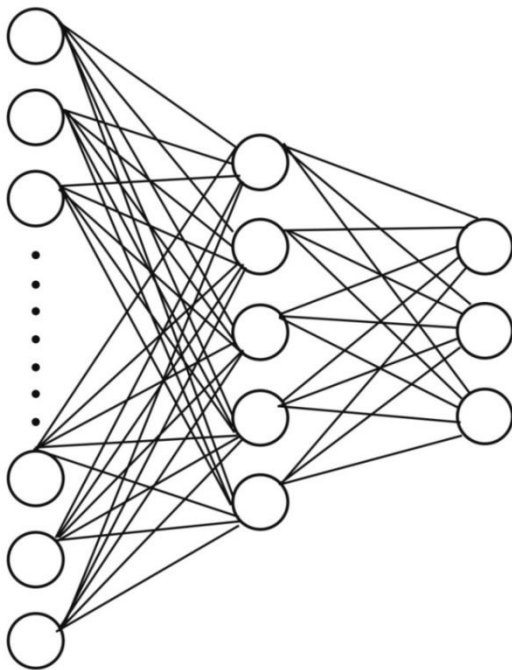
Παρατήρηση 1.3.2 (βλ. [5,σελ. 149 , σελ. 148]): Για την ειδική περίπτωση του \mathbb{R}^2 ($n = 2$) και την XOR συνάρτηση διαπιστώνουμε ότι τα inputs για να παράγουν τις τιμές 0,1 πρέπει να βρίσκονται σε διαγώνιες αντίθετες γωνίες ενός τετραγώνου. Αυτό θα σήμαινε ότι ένα perceptron δεν μπορεί να χωρίσει γραμμικά το χώρο μεταξύ τους και πρέπει να αντιμετωπιστεί με μη- γραμμικά δίκτυα. Παρακάτω έπεται γραφική απεικόνιση ενός σετ διανυσμάτων $X = \{(a, a), (a, b), (b, a), (b, b)\}$ με $a, b \in \mathbb{R}$



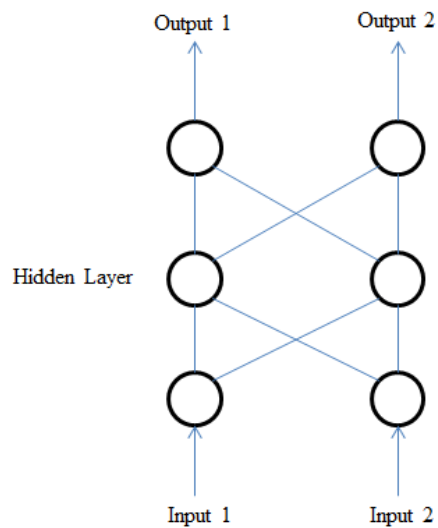
Σύμφωνα με το παραπάνω σχήμα, καταλαβαίνουμε πως δεν έχει την δυνατότητα “γραμμικού διαχωρισμού” του χώρου.

1.3.2. Πολύ-επίπεδα perceptron

Μια γενίκευση των perceptrons, είναι η κλάση νευρωνικών δικτύων που ονομάζεται πολύ- επίπεδα perceptron (multilayer perceptron) (MLP). Τα συγκεκριμένα δίκτυα αποτελούνται από ένα σετ κόμβων που απαρτίζουν τα εισαγόμενα δεδομένα, ένα ή περισσότερα κρυφά στρώματα από υπολογιστικούς κόμβους και το επίπεδο εξόδου. Το εισαγόμενο σήμα διαδίδεται μέσω του δικτύου με εμπρόσθια τροφοδότηση από στρώμα σε στρώμα. Τα MLP έχουν εφαρμοστεί επιτυχώς ώστε να επιλυθούν διάφορα δύσκολα προβλήματα, εκπαιδεύοντας το δίκτυο με έναν αλγόριθμο μάθησης που καλείται οπισθοδρόμηση σφάλματος (Ο αλγόριθμος αυτός βασίζεται στον κανόνα μάθησης που βελτιώνει τη σύγκλιση του. Παρακάτω υπάρχει η γραφική απεικόνιση δύο multilayer perceptron (Σχήμα 3, Σχήμα 4) εκ των οποίων το πρώτο πρόκειται για ένα δίκτυο με n εισαγόμενα δεδομένα, ένα κρυφό στρώμα και τρεις εξόδους, ενώ το δεύτερο για ένα δίκτυο με δύο inputs ένα κρυφό στρώμα και δύο outputs.



Σχήμα 3



Σχήμα 4

Ορισμός 1.3.1: Λογικές συναρτήσεις (Boolean) ονομάζονται εκείνες οι συναρτήσεις βάση των οποίων μπορεί κανείς να συμπεράνει αν μια πρόταση είναι αληθείς ή ψευδής. Η λογική συνάρτηση AND επιστρέφει TRUE όταν όλα τα επιχειρήματα της πρότασης ισχύουν και FALSE διαφορετικά. Η λογική συνάρτηση OR επιστρέφει TRUE όταν τουλάχιστον ένα από τα επιχειρήματα της πρότασης ισχύει και FALSE διαφορετικά.

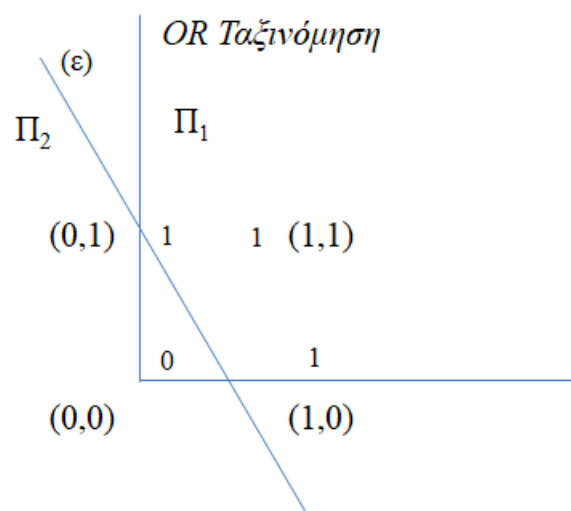
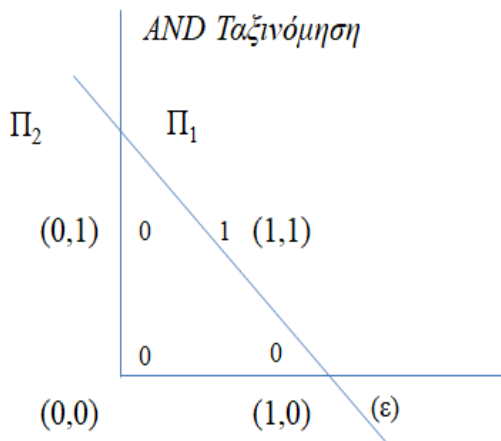
Παρατήρηση 1.3.1: Οι λογικές συναρτήσεις AND και OR έχουν την ιδιότητα του γραμμικού διαχωρισμού του χώρου (linear separator), σε αντίθεση με τη συνάρτηση XOR που όπως δείξαμε πριν δεν έχει την ιδιότητα αυτή.

Πιο συγκεκριμένα, η μαθηματική επέκταση των δύο λογικών συναρτήσεων δίνεται από τα παρακάτω ταμπλό.

X	Y	X AND Y
0	0	0
0	1	0
1	0	0
1	1	1

X	Y	X OR Y
0	0	0
0	1	1
1	0	1
1	1	1

Όπου, τα μηδενικά (0) εκφράζουν τα FALSE επιχειρήματα και οι μονάδες (1) τα TRUE επιχειρήματα των προτάσεων. Έτσι η αναπαράσταση των παραπάνω στο καρτεσιανό σύστημα συντεταγμένων είναι:

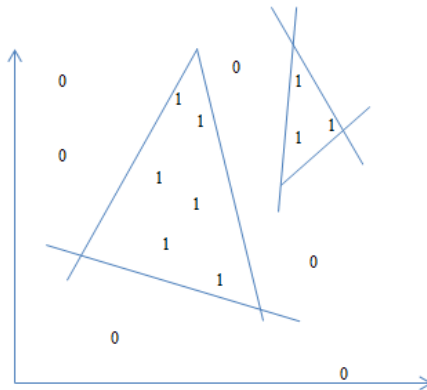


Που πράγματι και στις δύο περιπτώσεις επιτυγχάνεται η γραμμική διαχώριση του χώρου, σε δύο επίπεδα. Το επίπεδο Π_1 περιέχει τις μονάδες (1), ενώ το επίπεδο Π_2 περιέχει τα μηδενικά (0).

Ορισμός 1.3.2: Διαμέριση ενός μη κενού συνόλου A , είναι ένα σύνολο $U = \{U_i\}$ υποσύνολο του A , τα οποία είναι ανά δύο ξένα μεταξύ τους και ισχύει ότι $\forall i \neq j$ $U_i \cap U_j = \emptyset$ και η $\bigcup_i U_i = A$.

Παρατήρηση: Για να κατηγοριοποιηθεί ένα σύνολο πεπερασμένων στοιχείων, χρειαζόμαστε μη κενά πεπερασμένα υποσύνολα του χώρου.

Με βάση τα παραπάνω, παρατηρείται πως για να κατηγοριοποιηθεί κάποιο σετ σημείων σε δύο σετ, μπορούμε να το καταφέρουμε με μια αλληλουχία τριών perceptron. Πιο συγκεκριμένα, το πρώτο perceptron χωρίζει το χώρο σε δύο επίπεδα, η λογική συνάρτηση AND δημιουργεί πολύγωνα και τέλος η λογική συνάρτηση OR αναθέτει αυτά τα πολύγωνα σε κλάσεις. Επομένως, μπορούμε να κατασκευάσουμε οποιαδήποτε διαμέριση του \mathbb{R}^n σε πολύγωνα και να τα κατατάσσουμε τις περιοχές αυτές σε κλάσεις. Τα παραπάνω θα παρουσιαστούν στο (Σχήμα 5) στον χώρο του \mathbb{R}^2 .



Σχήμα 5

Βλέποντας ότι ένα πολύ-επίπεδο perceptron μπορεί να διαχωρίσει οποιοδήποτε πεπερασμένο σετ σημείων στον \mathbb{R}^n μέσω των συναρτήσεων AND και OR καταλαβαίνουμε ότι εφαρμόζεται μόνο για διακριτές λογικές συναρτήσεις.

Παρόλα αυτά, εμείς θα θέλαμε τα δίκτυα μας να εφαρμόζονται και σε συνεχή προβλήματα. Πιο συγκεκριμένα, ως συνήθως τα εισαγόμενα δεδομένα είναι διανύσματα στον \mathbb{R}^n και το output y του δικτύου είναι διάνυσμα του \mathbb{R}^m συνήθως $m \ll n$ (όπου σε πολλές περιπτώσεις $m=1$). Έτσι το δίκτυο μας υπολογίζει μια συνάρτηση $g: \mathbb{R}^n \rightarrow \mathbb{R}^m$, που τη θεωρούμε προσέγγιση μιας συνάρτησης

$f : \mathbb{R}^n \rightarrow \mathbb{R}^m$. Να σημειωθεί ότι η συνάρτηση g που κατασκευάζεται από το δίκτυο, πρόκειται για μια συνεχή συνάρτηση με σκοπό την προσέγγιση οποιασδήποτε συνάρτησης f . Έστω τώρα A, B δύο πεπερασμένα σετ στον \mathbb{R}^n , όπου το σύνολο A περιέχει “1” και το σύνολο B “0” όπως φαίνεται στο (Σχήμα5) και εμείς θέλουμε να παράγουμε το output 1 για τα σημεία στο σύνολο A και 0 για το σύνολο B . Κατασκευάζοντας έτσι ένα πεπερασμένο σετ πολυγώνων ($A \subset Q := \bigcup_j P_j$) και $B \cap P_j = \emptyset$, επιλέγοντας $m=1$ και θέτοντας για χαρακτηριστική συνάρτηση f ως $f(x)=1$ εάν το $x \in Q$ και 0 αλλιώς, το δίκτυο μας μπορεί πετύχει την ταξινόμηση του συνόλου.

Συνοπώς, αναγνωρίζετε ότι τα νευρωνικά δίκτυα μπορούν να εξεταστούν και με μεθόδους της θεωρίας προσέγγισης, γεγονός με το οποίο θα ασχοληθούμε αναλυτικά στο κεφάλαιο 2.

Παρατήρηση 1.3.2: Για τα παραπάνω χρησιμοποιήθηκε μια μη-συνεχής χαρακτηριστική συνάρτηση, παρόλα αυτά αφού τα σύνολα A, B είναι πεπερασμένα μπορεί να αντιμετωπιστεί αυτό με διαδικασίες εξομάλυνσης (smooth process).

2. Πυκνότητα και προσέγγιση νευρωνικών δικτύων

Σε αυτό το κεφάλαιο θα μελετήσουμε το αποτέλεσμα που εξάγει ένα νευρωνικό δίκτυο από τη σκοπιά της θεωρίας προσέγγισης. Πιο συγκεκριμένα θα θέλαμε να ξέρουμε εάν το σύνολο των πιθανών συναρτήσεων g που θα έχουμε ως output, είναι πυκνό σε έναν κατάλληλο συναρτησιακό χώρο όπου θα περιέχει την συνάρτηση f που θέλουμε να προσεγγίσουμε. Χωρίς βλάβη της γενικότητας και επειδή σχεδόν όλα τα αποτελέσματα εξετάζουν την υπόθεση ενός απλού output ($m=1$), θα κάνουμε και εμείς αυτή την απλοποίηση.

Ορισμός 2.0.1 (βλ. [2, Ορισμός 4.1, σελ. 45]): Μετρικός χώρος είναι ένα ζεύγος (X, ρ) όπου X είναι ένα μη κενό σύνολο $\rho : X \times X \rightarrow \mathbb{R}$ μια απεικόνιση που ικανοποιεί τις ιδιότητες :

- $\rho(x, y) \geq 0$ για κάθε $x, y \in X$ και $\rho(x, y) = 0$ αν και μόνο αν $x = y$.
- $\rho(x, y) = \rho(y, x)$ για κάθε $x, y \in X$ (συμμετρική ιδιότητα).
- $\rho(x, z) \leq \rho(x, y) + \rho(y, z)$ για κάθε $x, y, z \in X$ (τριγωνική ιδιότητα).

Η απεικόνιση ρ , ονομάζεται μετρική, τα στοιχεία του συνόλου X ονομάζονται σημεία και ο αριθμός $\rho(x, y)$ ονομάζεται απόσταση του x από το y .

Παραδείγματα μετρικών χώρων είναι:

- Η συνήθης μετρική, στο σύνολο \mathbb{R} των πραγματικών αριθμών που ορίζεται ως $\rho(x, y) = |x - y|$ για κάθε $x, y \in \mathbb{R}$.
- Η ευκλείδεια μετρική ρ_2 στο σύνολο \mathbb{R}^k των διατεταγμένων κ-άδων πραγματικών αριθμών όπου ορίζεται ως εξής: Για δύο στοιχεία $\vec{x} = (x_1, x_2, \dots, x_k)$, $\vec{y} = (y_1, y_2, \dots, y_k)$ του \mathbb{R}^k είναι
$$\rho_2(x, y) = \left(\sum_{i=1}^k (x_i - y_i)^2 \right)^{1/2}.$$

Ορισμός 2.0.2 (βλ. [2, Ορισμός 5.1, σελ. 56]): Έστω (X, ρ) μετρικός χώρος και $x_0 \in X$. Τότε για κάθε $\varepsilon > 0$ ορίζουμε την ανοιχτή σφαίρα κέντρου x_0 και ακτίνας ε , ως το σύνολο $S(x_0, \varepsilon) = \{x \in X : \rho(x_0, x) < \varepsilon\}$.

Ορισμός 2.0.3 (βλ. [2, Ορισμός 9.2, σελ. 103]): Έστω (X, ρ) μετρικός χώρος και $A \subset X$. Μια οικογένεια $\{G_i\}_{i \in I}$ υποσυνόλων του X λέγεται κάλυμμα του A αν $A \subset \bigcup_{i \in I} G_i$. Αν επιπλέον για κάθε $i \in I$ το G_i είναι ανοιχτό, το $\{G_i\}_{i \in I}$ λέγεται ανοιχτό κάλυμμα, ενώ στην περίπτωση που το I είναι πεπερασμένο το $\{G_i\}_{i \in I}$ πεπερασμένο κάλυμμα. Αν $J \subset I$ και $A \subset \bigcup_{i \in J} G_i$ το $\{G_i\}_{i \in J}$ λέγεται υποκάλυμμα του $\{G_i\}_{i \in I}$ (για το A).

Ορισμός 2.0.4 (βλ. [2, Ορισμός 9.2, σελ. 103]): Έστω (X, ρ) μετρικός χώρος και $K \subset X$. Το K θα καλείται συμπαγές αν κάθε ανοιχτό κάλυμμα του K έχει πεπερασμένο υποκάλυμμα, δηλαδή αν για κάθε οικογένεια $\{G_i\}_{i \in I}$ ανοιχτών υποσυνόλων του X με $K \subset \bigcup_{i \in I} G_i$ υπάρχουν $n \in \mathbb{N}$ και $i_1, i_2, \dots, i_n \in I$ ώστε $K \subset \bigcup_{k=1}^n G_{i_k}$. Ειδικότερα αν $K = X$ τότε ο X θα καλείται συμπαγής μετρικός χώρος.

Ορισμός 2.0.5 (βλ. [2, Ορισμός 6.1, σελ. 65]): Έστω (X, ρ) μετρικός χώρος και $D \subset X$. Ένα σημείο $x \in X$ καλείται οριακό του συνόλου D αν $\forall \varepsilon > 0$ ισχύει $S(x, \varepsilon) \cap D \neq \emptyset$.

Ορισμός 2.0.6 (βλ. [2, Ορισμός 6.2 σελ. 65]): Έστω (X, ρ) μετρικός χώρος και $D \subset X$. Η κλειστότητα του D ορίζεται ως $\bar{D} = \{x \in X : x \text{ είναι οριακό σημείο του } D\}$.

Ορισμός 2.0.7 (βλ. [2, Ορισμός 7.9, σελ. 85]): Έστω (X, ρ) μετρικός χώρος και $D \subset X$. Το D λέγεται πυκνό υποσύνολο του X αν η κλειστότητα του D ισούται με το χώρο X ($\overline{D} = X$).

Εναλλακτικός ορισμός του πυκνού συνόλου (βλ. [7, Ορισμός 2.6, σελ. 361]): Ένα υποσύνολο S ενός μετρικού χώρου (X, ρ) , είναι πυκνό σε ένα υποσύνολο A , εάν $\forall \varepsilon > 0$ και για κάθε $a \in A$, υπάρχει ένα $s \in S$ τέτοιο ώστε $\rho(s, a) < \varepsilon$. Με άλλα λόγια, ένα στοιχείο του συνόλου S μπορεί να προσεγγίσει ένα στοιχείο του συνόλου A , σε οποιοδήποτε επιθυμητό βαθμό.

Ορισμός 2.0.8 (βλ. [7, Ορισμός 2.7, σελ. 361]): Ένα υποσύνολο B καλείται ομοιόμορφα πυκνό στον συμπαγή χώρο $C(K)$, εάν για κάθε συμπαγές υποσύνολο του $B \subset \mathbb{R}^n$, το S είναι ρ_β -πυκνό στον $C(K)$, όπου η μετρική $\rho_\beta(f, g) \equiv \sup |f(x) - g(y)|$ για όλα τα $x \in B$ και $f, g \in C(K)$.

Μία ακολουθία συναρτήσεων $\{f_n\}$ συγκλίνει ομοιόμορφα, εάν για όλους τους συμπαγείς χώρους $B \subset \mathbb{R}^n$ η μετρική $\rho_\beta(f_n, f) \rightarrow 0$ όταν το $n \rightarrow \infty$.

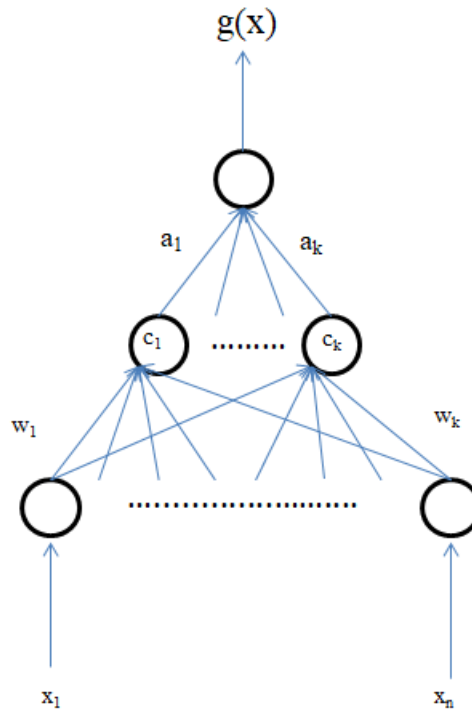
Ενόψει της δυσκολίας που θα προκύψει χρησιμοποιώντας μη-διαφορίσιμες και ασυνεχείς συναρτήσεις, συνηθίζεται να δουλεύουμε με λείες συναρτήσεις ενεργοποίησης ($\sigma(x)$). Συναρτήσεις ενεργοποίησης όπως ήδη έχουμε εξηγήσει στο πρώτο κεφάλαιο είναι σιγμοειδείς συναρτήσεις που δέχονται ως όρισμα το άθροισμα των inputs πολλαπλασιαζόμενα με τα αντίστοιχα βάρη, με κατεύθυνση την έξοδο.

Το αποτέλεσμα το οποίο θα συζητηθεί αργότερα αναφέρεται σε ένα multilayer perceptron με ένα κρυφό στρώμα k μονάδων (Σχήμα 6). Ειδικότερα ορίζουμε ως w_j το διάνυσμα βάρους που σχετίζεται με τα inputs του “ j ” κρυφού νευρώνα. Έτσι για ένα δοσμένο διάνυσμα x , υποθέτοντας ότι κάθε κρυφός νευρώνας έχει την ίδια συνάρτηση ενεργοποίησης και επιτρέψουμε κάποια πραγματική μετατόπιση της οριακής τιμής στη θέση “ j ” (c_j) το όρισμα της συνάρτησης ενεργοποίησης θα είναι $w_j^T \cdot x$, συνεπώς η τιμή που θα έχουμε σαν έξοδο από τον “ j ” νευρώνα θα είναι η $\sigma(w_j^T \cdot x + c_j)$.

Συναρτήσεις αυτής της μορφής καλούνται ridge functions, συναρτήσεις δηλαδή που είναι σταθερές στο υπερεπίπεδο $w_j^T \cdot x = \text{constant}$. Είναι φανερό ότι για την ειδική περίπτωση των δύο διαστάσεων το περίγραμμα της συνάρτησης σχηματίζει ευθείες κορυφογραμμές (ridges).

Συνοψίζοντας τα παραπάνω η output συνάρτηση g του δικτύου (Σχήμα 6) παίρνει την παρακάτω μορφή.

$$g(x) = \sum_j^k a_j \sigma(w_j^T \cdot x + c_j)$$



Σχήμα 6

Ορισμός 2.0.9 (βλ. [3, σελ. 25]): Ορίζουμε L_p χώρο, έναν συναρτησιακό χώρο μετρήσιμων συναρτήσεων, όπου η " p " δύναμη της απόλυτης τιμής της συνάρτησης είναι Lebesgue ολοκληρώσιμη. Πιο συγκεκριμένα, για $1 \leq p < \infty$ και έναν μετρήσιμο χώρο S η p -νόρμα

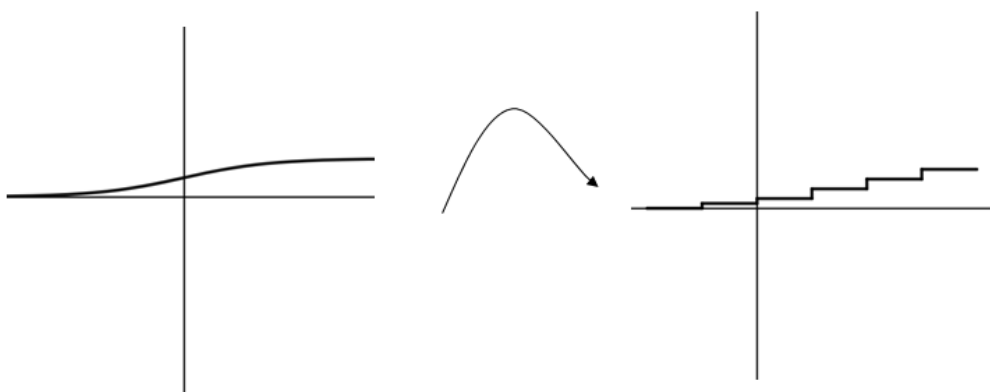
$$\|f_p\| \equiv \left(\int_S |f_p| d\mu \right)^{1/p} < \infty.$$

Η συσχέτιση των χώρων L_p με τη μετρική χώρου ορίζεται ως $\rho_p = \|f - g\|_p$.

2.1 Άμεση προσέγγιση της πυκνότητας

Σκοπός της υποενότητας αυτής είναι η προσέγγιση της απόδειξης της πυκνότητας στη μονοδιάστατη περίπτωση. Μια άμεση προσέγγιση για αυτήν έχει σχεδιαστεί από τους Chen και Liu (1991), όπου είναι γνωστή ως quasi interpolant και πρακτικά σημαίνει ότι η απόδειξη τους είναι βασισμένη σε γραμμικό συνδυασμό συναρτησιακών τιμών. Η βασική ιδέα είναι η προσέγγιση της συνάρτησης f από τμηματικές σταθερές, όπου οι συναρτήσεις ενεργοποίησης αντιμετωπίζονται ως ομαλές σταθερές.

Έπεται γραφική αναπαράσταση:



Ορισμός 2.1.1 (βλ. [5, σελ. 154]): Έστω K συμπαγής χώρος στον \mathbb{R} με $f \in C(K)$ και $\delta > 0$. Ορίζεται ως μέτρο συνέχειας η ποσότητα $\omega(f, \delta) = \sup |f(x) - f(y)|$, για τα όλα $x, y \in K$ με $|x - y| < \delta$. Στο παραπάνω μέτρο αφού η συνάρτηση f είναι συνεχής, η ποσότητα ω είναι πεπερασμένη και τείνει μονοτονικά στο 0 όπως τείνει το δ στο 0.

Ορισμός 2.1.2: Έστω K ανοιχτό υποσύνολο του \mathbb{R} . Μια συνάρτηση f καλείται Lipschitz αν για κάθε $x, y \in K$ και $L > 0$ ισχύει ότι $|f(x) - f(y)| < L|x - y|$.

Έτσι συνδυάζοντας τους ορισμούς των συνεχών και των Lipschitz συναρτήσεων άμεσα προκύπτει ότι $\omega(f, \delta) < L\delta$. Ομοίως η συνέπεια αυτή εφαρμόζεται και για k -Lipschitz και διαφορίσιμες συναρτήσεις.

Ορισμός 2.1.3: Έστω χώρος X και $A \subset X$. Ορίζουμε δείκτρια την συνάρτηση της μορφής:

$$\mathbb{1}_A(x) = \begin{cases} 1, & x \in A \\ 0, & x \notin A. \end{cases}$$

Ορισμός 2.1.4: Έστω K ένα υποσύνολο του \mathbb{R} . Ορίζουμε ως κλιμακωτή συνάρτηση (step function, Heaviside Function) μια συνάρτηση που μπορεί να γραφτεί ως πεπερασμένος γραμμικός συνδυασμός δείκτριων συναρτήσεων στο υποσύνολο αυτό.

Το μέτρο της συνέχειας μας δίνει μια εκτίμηση για το πόσο καλά μπορεί να προσεγγιστεί η συνάρτηση μας από τμηματικές σταθερές. Έτσι χωρίς βλάβη της γενικότητας, επιλέγουμε ως K ένα υποσύνολο του $[0,1]$ και $n \in \mathbb{N}$.

Θεωρούμε κλιμακωτή συνάρτηση $h_n(x) = f\left(\frac{r}{n}\right)$, $x \in \left[\frac{r}{n}, \frac{r+1}{n}\right]$, $r = 0, \dots, n-1$.

Προφανώς από τα παραπάνω προκύπτει ότι η άπειρη νόρμα $\|f - h_n\|_\infty \leq \omega\left(f, \frac{1}{n}\right)$.

Θα χρησιμοποιήσουμε την τμηματικά σταθερή συνάρτηση:

$h_n(x) := f(0) + \sum_{r=1}^{\mu} \left\{ f\left(\frac{r}{n}\right) - f\left(\frac{r-1}{n}\right) \right\}$, όπου μ είναι ο μεγαλύτερος ακέραιος που δεν υπερβαίνει το nx .

Στη συνέχεια θεωρούμε τη συνεχή, σιγμοειδή συνάρτηση $\sigma(x) = \frac{1}{1+e^{-x}}$ και την παραμετροποίηση του x με την τιμή ax όπου $a > 1$, ώστε να δημιουργήσουμε μια σημειακή σύγκλιση σε μια οριακή συνάρτηση με την τιμή $\sigma(0)$ στο 0. Έτσι όταν η τιμή $a \rightarrow \infty$ η ποσότητα $\sigma(ax) \rightarrow 0$ ή 1, αναλόγως με το αν το $x < 0$ ή $x > 0$. Οπότε η ποσότητα $\sigma(ax) - \sigma(a(x-1))$ θα προσεγγίσει την δείκτρια συνάρτηση με την τιμή 1 στο υποσύνολο $(0,1)$, και 0 αλλού. Για τις περιπτώσεις ασυνέχειας στα σημεία 0,1 η σύγκλιση δεν μπορεί να είναι ομοιόμορφη.

Ορισμός 2.1.5 (βλ. [2, Ορισμός 8.20, σελ. 98]): Μια συνάρτηση $f : A \rightarrow \mathbb{R}$ είναι ομοιόμορφα συνεχής, αν $\forall \varepsilon > 0, \exists \delta > 0$ τέτοια ώστε $\forall x_0, x \in A$, $|x - x_0| < \delta$ τότε $|f(x) - f(x_0)| < \varepsilon$.

Παρόλα αυτά, παρατηρήθηκε από τους Chen et.al. ότι για την αντιμετώπιση του προβλήματος της ασυνέχειας, συνδυάζοντας την κατασκευή της $h_n(x)$ με τη σιγμοειδή, μπορούσαμε να πετύχουμε μια ομοιόμορφη προσέγγιση αφού οι ασυνέχειες μηδενίζονται όταν το $n \rightarrow \infty$.

Πιο συγκεκριμένα θεωρούμε μια σιγμοειδή συνάρτηση συνεχή στον \mathbb{R} και ορίζουμε A_n τέτοιο ώστε να πετύχουμε τη διακριτοποίηση του χώρου της f .

Διαμερίζουμε το χώρο σε μέγεθος $A_n \leq x \leq A_{n+1}$ με σκοπό η συνάρτηση f να παίρνει την τιμή $f(A_n)$ στο εκάστοτε διάστημα $(A_n, A_{n+1}]$.

Μια απλοϊκή τέτοια διαμέριση του χώρου θα μπορούσε να είναι η $\frac{r}{n} \leq x \leq \frac{r-1}{n}$.

Στην περίπτωση μας ορίζουμε τα διαστήματα A_n τέτοια ώστε $1 - \frac{1}{n} \leq \sigma(x) \leq 1 + \frac{1}{n}$,

$$\forall x \geq A_n \text{ και } |\sigma(x)| \leq \frac{1}{n}, \quad \forall x \leq -A_n.$$

Οπότε ορίζουμε την ψευδοσυνάρτηση παρεμβολής $g_n(x)$ που είναι τμηματικά σταθερή, παίρνοντας μια τιμή παρεμβολής της f στο παραπάνω διάστημα,

$$g_n(x) = f(0) + \sum_{r=1}^n \left\{ f\left(\frac{r}{n}\right) - f\left(\frac{r-1}{n}\right) \right\} \sigma(A_n(nx-r)) \quad \text{για } x \in [0,1].$$

Θεώρημα 2.1.1 (βλ. [5, σελ. 155]): Υπάρχει σταθερά c τέτοια ώστε για $f \in C[0,1]$ ισχύει ότι:

$$\|f - g_n\|_{\infty} \leq c \omega\left(f, \frac{1}{n}\right).$$

Απόδειξη. Ξεκινώντας από το αριστερό μέλος της ανισότητας έχουμε

$$\|f - g_n\|_{\infty} = \|f - h_n + h_n - g_n\|_{\infty}, \quad \text{από τριγωνική ανισότητα έχουμε}$$

$$\|f - h_n + h_n - g_n\|_{\infty} \leq \|f - h_n\|_{\infty} + \|h_n - g_n\|_{\infty}$$

Όπως έχουμε ήδη προαναφέρει η ποσότητα $\|f - h_n\|_{\infty}$ είναι φραγμένη από την ποσότητα $\omega\left(f, \frac{1}{n}\right)$. Συνεπώς πρέπει να αποδειχτεί ότι φράζεται η ποσότητα $\|h_n - g_n\|_{\infty}$.

Για κάθε $x \in [0,1]$ και για μ όπως έχει οριστεί προηγουμένως έχουμε:

$$\begin{aligned} h_n(x) - g_n(x) &= f(0) + \sum_{r=1}^{\mu} \left\{ f\left(\frac{r}{n}\right) - f\left(\frac{r-1}{n}\right) \right\} \\ &\quad - f(0) - \sum_{r=1}^{\mu} \left\{ f\left(\frac{r}{n}\right) - f\left(\frac{r-1}{n}\right) \right\} \sigma(A_n(nx-r)) \\ &\quad + \sum_{r=\mu+1}^n \left\{ f\left(\frac{r}{n}\right) - f\left(\frac{r-1}{n}\right) \right\} \sigma(A_n(nx-r)) \\ &= \sum_{r=1}^{\mu} \left\{ f\left(\frac{r}{n}\right) - f\left(\frac{r-1}{n}\right) \right\} \{1 - \sigma(A_n(nx-r))\} + \sum_{r=\mu+1}^n \left\{ f\left(\frac{r}{n}\right) - f\left(\frac{r-1}{n}\right) \right\} \sigma(A_n(nx-r)). \end{aligned}$$

Όταν $r \leq \mu - 1$ και αφού το μ είναι επιλεγμένο κατάλληλα έτσι ώστε

$$\mu \leq n \cdot x \Rightarrow r \leq n \cdot x - 1 \Rightarrow n \cdot x - r \geq 1 \Rightarrow A_n(nx - r) \geq A_n,$$

που σημαίνει $|1 - \sigma(A_n(nx - r))| \leq \frac{1}{n}$ από τον ορισμό του A_n .

Όταν $r \geq \mu + 2$ έχουμε $|\sigma(A_n(nx - r))| \leq \frac{1}{n}$.

Τελικά, αντικαθιστώντας τις δύο παραπάνω σχέσεις:

$$|h_n(x) - g_n(x)| \leq \omega(f, \frac{1}{n}) + \left| \left\{ f\left(\frac{\mu}{n}\right) - f\left(\frac{\mu-1}{n}\right) \right\} + \left\{ f\left(\frac{\mu+1}{n}\right) \right\} \sigma(A_n(nx - \mu - 1)) \right|.$$

όπου το δεύτερο μέλος της ανισότητας φράζεται από την ποσότητα $2(1+S)\omega(f, \frac{1}{n})$

Συνεπώς η αρχική μας ποσότητα :

$$\|f - g_n\|_\infty \leq \omega(f, \frac{1}{n}) + 2(1+S)\omega(f, \frac{1}{n}) = c\omega(f, \frac{1}{n}). \quad \blacksquare$$

Για να περάσουμε από την μονοδιάστατη στην πολυδιάστατη περίπτωση θα χρησιμοποιήσουμε την μέθοδο του τανυστικού γινομένου προτύπων στη δισδιάστατη περίπτωση.

Ορισμός 2.1.6 (βλ. [2, σελ. 123]): Συναρτήσεις βάσης ονομάζουμε τα στοιχεία που αποτελούν μια βάση ενός συναρτησιακού χώρου.

Έστω δύο σύνολα με συναρτήσεις βάσης, $\{\Phi_1, \Phi_2, \dots, \Phi_\mu\}, \{\Psi_1, \Psi_2, \dots, \Psi_\nu\}$, όπου $\Phi_i, \Psi_j : \mathbb{R} \rightarrow \mathbb{R}$.

Το τανυστικό γινόμενο των βάσεων (διανυσματικών γινόμενο) είναι το σύνολο $\mu \times \nu$ των συναρτήσεων όπως ορίζεται παρακάτω :

$$\zeta_{i,j}(x, y) = \Phi_i(x)\Psi_j(y).$$

Σε αυτό το σημείο να σημειωθεί ότι τα σύνολα $\{\Phi_1, \Phi_2, \dots, \Phi_\mu\}, \{\Psi_1, \Psi_2, \dots, \Psi_\nu\}$ μπορούν να περιέχουν συναρτήσεις από διαφορετικές κλάσεις, συνήθως όμως στην πράξη τα δύο σύνολα έχουν τις ίδιες κλάσεις συναρτήσεων (εκθετικές, τριγωνομετρικές, κλπ).

Στην περίπτωση μας θα χρησιμοποιήσουμε την κατασκευή αυτή σε ridge functions όπου για λόγους απλότητας θεωρούμε ότι εφαρμόζεται η ίδια σιγμοειδή συνάρτηση και για το x και για το y .

Έτσι σύμφωνα με τα παραπάνω μια τυπική μονοδιάστατη ridge function για τις δύο διαστάσεις θα έπαιρνε τη μορφή :

$$\sigma(a_i x + c_i) \text{ και } \sigma(b_j y + d_j).$$

Το τανυστικό γινόμενο βάσης συναρτήσεων θα πάρει τη μορφή:

$$\sigma(a_i x + c_i) \sigma(b_j y + d_j).$$

Γενικά το παραπάνω γινόμενο δεν μας παρέχει μια δισδιάστατη ridge function.

Παρόλα αυτά για την ειδική περίπτωση της εκθετικής συνάρτησης ως ridge function ($\sigma(x) = e^x$) πετυχαίνεται το επιθυμητό αποτέλεσμα αφού :

$$\begin{aligned} \sigma(a_i x + c_i) \sigma(b_j y + d_j) &= e^{(a_i x + c_i)} e^{(b_j y + d_j)} \\ &= e^{(a_i x + c_i + b_j y + d_j)} \\ &= \sigma(a_i x + c_i + b_j y + d_j). \end{aligned}$$

Παρατήρηση 2.1.1: Η παραπάνω τοποθέτηση έχει οδηγήσει πολλούς ερευνητές στην παραγωγή n-διάστατων προσεγγίσεων των ridge functions. Η βασική ιδέα είναι να χρησιμοποιήσουμε την ειδική περίπτωση της εκθετικής συνάρτησης αποδεικνύοντας την πυκνότητα τους και έπειτα να χρησιμοποιήσουμε την μονοδιάστατη περίπτωση όπως το θεώρημα (2.1.1) ώστε να προσεγγίσουμε την εκθετική συνάρτηση από γραμμικό συνδυασμό της επιθυμητής σ .

Ορισμός 2.1.7: Ένα σύνολο συναρτήσεων ονομάζεται θεμελιώδης σε ένα χώρο, εάν ο γραμμικός συνδυασμός τους είναι πυκνός σε αυτό τον χώρο.

Θεώρημα 2.1.2 (βλ. [2, Θεώρημα 11.8, σελ. 126]): Έστω K ένα συμπαγή σύνολο στον \mathbb{R}^n . Τότε το σύνολο E συναρτήσεων της μορφής $\mu(x) = \exp(a^T \cdot x)$, $a \in \mathbb{R}^n$, είναι θεμελιώδης στο $C(K)$.

Όπου $C(K)$ ο συναρτησιακός χώρος του συνόλου K με συνεχείς συναρτήσεις, εμπλουτισμένες με ομοιόμορφα συνεχείς νόρμες.

Για την απόδειξη του θεωρήματος θα χρησιμοποιήσουμε το θεώρημα προσέγγισης των Stone- Weierstrass.

Θεώρημα 2.1.3 (βλ. [9, Stone- Weierstrass theorem, σελ. 364]):

Κάθε συνεχής συνάρτηση ορισμένη σε ένα κλειστό υποσύνολο $[\alpha, \beta]$, μπορεί να προσεγγιστεί ομοιόμορφα από μια πολυωνυμική συνάρτηση.
(Θεώρημα Stone- Weierstrass)

Μια διαφορετική εκδοχή του παραπάνω θεωρήματος είναι η παρακάτω :

Έστω $C(K)$ χώρος συναρτήσεων, εάν ο χώρος K είναι άλγεβρα και έχει την ιδιότητα των διαχωρισμένων σημείων, τότε το $C(K)$ είναι πυκνό.

Ορισμός 2.1.8 (βλ. [9, σελ. 364]) : Ένα σύνολο, ενός χώρου λέμε ότι έχει την ιδιότητα των διαχωρισμένων σημείων, εάν ανά δύο διαφορετικά σημεία x, y του χώρου υπάρχει συνάρτηση p τέτοια ώστε

$$p(x) \neq p(y)$$

Συνδυάζοντας όλα τα παραπάνω οδηγούμαστε στο θεώρημα (2.1.4) όπου θα μας δώσει το επιθυμητό αποτέλεσμα για την δισδιάστατη περίπτωση.

Πιο συγκεκριμένα :

Θεώρημα 2.1.4 (βλ. [5, σελ. 157]) : Έστω K συμπαγές σύνολο στον \mathbb{R}^n . Τότε το

σύνολο F συναρτήσεων της μορφής $g(x) = \sum_j^k a_j \sigma(w_j^T \cdot x + c_j)$ όπου σ μια

συνεχής σιγμοειδής συνάρτηση είναι πυκνό στον $C(K)$.

Απόδειξη : Έστω συνάρτηση $f \in C(K)$. Τότε από Θεώρημα (2.1.2), για κάθε $\varepsilon > 0$ υπάρχει ένας πεπερασμένος αριθμός m από διανύσματα a_i , τέτοια ώστε :

$$\left\| f - \sum_{i=1}^m e^{a_i^T \cdot x} \right\|_{\infty} < \frac{\varepsilon}{2}.$$

Αφού υπάρχουν m βαθμωτά μεγέθη της μορφής $a_i^T \cdot x \in \mathbb{R}$ μπορούμε να βρούμε ένα διάστημα που να τα περιέχει όλα.

Άρα υπάρχει αριθμός Γ τέτοιος ώστε : $e^{a_i^T \cdot x} = e^{\Gamma \cdot y}$, όπου $y = \frac{a_i^T \cdot x}{\Gamma} \in [0,1]$.

Οπότε από το θεώρημα (2.1.1) καταλαβαίνουμε ότι η συνάρτηση $e^{\Gamma \cdot y}$ προσεγγίζεται από γραμμικό συνδυασμό συναρτήσεων της μορφής $\sigma(w_j^T \cdot x + c_j)$, με ομοιόμορφο σφάλμα μικρότερο του $\frac{\varepsilon}{2m}$, δίνοντας μας το επιθυμητό αποτέλεσμα. ■

2.2 k-Σιγμοειδής συνάρτηση

Το 1992 ο Mhaskar με τον Miccheli παρουσίασαν την ιδέα της k- σιγμοειδούς συνάρτησης. Μια k- σιγμοειδής συνάρτηση ορίζεται όπως παρακάτω :

$$\begin{cases} \sigma(x)/x^k \rightarrow 1, & x \rightarrow \infty \\ \sigma(x)/x^k \rightarrow 0, & x \rightarrow -\infty. \end{cases}$$

Για την περίπτωση που $k = 0$ πρόκειται για την κλασική περίπτωση της σιγμοειδούς συνάρτησης.

Η ιδέα τους ήταν να αντικαταστήσουν την τμηματικά σταθερή συνάρτηση, από μια συνάρτηση τύπου spline (δηλαδή μια συνάρτηση η οποία ορίζεται από τμηματικά πολυώνυμα), μεγαλύτερου βαθμού, ώστε να πετύχουν προσεγγίσεις καλύτερου βαθμού για ομαλές συναρτήσεις. Προφανώς μια σιγμοειδή αυτής της μορφής δεν πρόκειται για τη συμβατική της μορφή που χρησιμοποιείται στα νευρωνικά δίκτυα, συνεπώς πρέπει παρουσιαστούν νέοι τρόποι προσέγγισης. Αρχικά, ασχολήθηκαν με την πολυπαραγοντική περίπτωση προσεγγίζοντας την συνάρτηση f με τανυστικό γινόμενο συναρτήσεων τύπου spline.

Έπειτα ασχολήθηκαν με το πρόβλημα προσέγγισης της f από ένα νευρωνικό δίκτυο, με σταθερό αριθμό νευρώνων διευθετημένα σε παραπάνω από ένα στρώμα. Ακόμα ο Mhaskar απέδειξε ότι οι πολυωνυμικοί όροι σε ένα spline μπορούν να αποσυντεθούν ως σύνθεση γραμμικών συναρτήσεων που προσεγγίζονται από σιγμοειδής συναρτήσεις.

Από τα παραπάνω καταλαβαίνουμε ότι αρκετά ενδιαφέροντα αποτελέσματα εκπέμπουν οι k -σιγμοειδής συναρτήσεις (για την περίπτωση $k > 0$), παρόλα αυτά οι k -σιγμοειδής συναρτήσεις δεν είναι αυτές που χρησιμοποιούνται στα νευρωνικά δίκτυα.

2.3 Δυϊκός χώρος και συνέλιξη μεθόδων της προσέγγισης

Στην ενότητα αυτή θα ασχοληθούμε με την προσέγγιση της πυκνότητας του χώρου με μεθόδους από τον ολοκληρωτικό λογισμό, βασισμένοι σε δυϊκούς χώρους και συνέλιξεις των μεθόδων.

Ορισμός 2.3.1 (βλ. [3, σελ. 33]): Ως γραμμικό συναρτησοειδές ορίζουμε τις γραμμικές απεικονίσεις $f : X \rightarrow \mathbb{R}$, όπου X είναι διανυσματικός χώρος.

Ορισμός 2.3.2 (βλ. [5, σελ. 160]): Έστω διανυσματικός χώρος με νόρμα $(X, \|\cdot\|)$ στον \mathbb{R} τότε ο δυϊκός του X που ορίζεται ως X^* , είναι ο χώρος που εμπεριέχει όλα τα φραγμένα γραμμικά συναρτησοειδή του X εφοδιασμένος με τη φυσική νόρμα $\|l\| = \sup_{x \in X, \|x\|=1} |l(x)|$

Ορισμός 2.3.3 (βλ. [2, σελ. 91]): Μια ακολουθία $\{x_n\}$ στοιχείων μετρικού χώρου (X, ρ) λέγεται Cauchy αν για κάθε $\varepsilon > 0$ υπάρχει $n_0 \in \mathbb{N}$ ώστε για κάθε $n, m \in \mathbb{N}$ με $n, m \geq n_0$ ισχύει $\rho(x_n, x_m) < \varepsilon$.

Ορισμός 2.3.4 (βλ. [2, σελ. 91]): Ένας μετρικός χώρος (X, ρ) λέγεται πλήρης αν κάθε Cauchy ακολουθία συγκλίνει.

Παρατήρηση 2.3.1: Οι χώροι Banach πρόκειται για διανυσματικούς χώρους πάνω στον \mathbb{R} εφοδιασμένοι με νόρμα οι οποίοι είναι πλήρης.

Στο παρακάτω θεώρημα θα συσχετίσουμε τη πυκνότητα ενός υπόχωρου $V \subset X$ στον X με τον δυϊκό του.

Πόρισμα 2.3.1 (βλ. [3, σελ. 57]): Έστω χώρος X με νόρμα και Y υπόχωρος του X . Αν $f: Y \rightarrow \mathbb{R}$ είναι φραγμένο γραμμικό συναρτησοειδές και $M \geq 0$ ώστε $f(y) \leq M \|y\|$ για όλα τα $y \in Y$, τότε υπάρχει $\tilde{f}: X \rightarrow \mathbb{R}$ γραμμική επέκταση της f , ώστε $\tilde{f}(x) \leq M \|x\|$ για όλα τα $x \in X$. (Βασική εφαρμογή του θεωρήματος Hahn-Banach)

Θεώρημα 2.3.1 (βλ. [5, σελ.160]): Έστω χώρος X και V υπόχωρος του. Τότε ο V είναι πυκνός στον X αν και μόνο αν, το γραμμικό συναρτησοειδές του $l \in X^*$ για το οποίο ισχύει $l(v) = 0, \forall v \in V$ είναι η τετριμμένη συνάρτηση $l(x) \equiv 0$.

Απόδειξη: \Rightarrow

Έστω ο χώρος V πυκνός στον X και l ένα γραμμικό συναρτησοειδές με $l(v) = 0, \forall v \in V$.

Έστω $x \in X$ αφού ο V είναι πυκνός στον X για κάθε $\varepsilon > 0$ υπάρχει $v \in V$ τέτοιο ώστε $\|x - v\| < \varepsilon$.

$$\begin{aligned} \text{Οπότε } |l(x)| &= |l(x) - l(v)| \\ &= |l(x - v)| \text{ (αφού το } l \text{ πρόκειται για ένα γραμμικό συναρτησοειδές)} \\ &\leq \|l\| \|x - v\| \text{ (από τον ορισμό του συνεχούς γραμμικού συναρτησοειδούς)} \\ &< \|l\| \varepsilon. \text{ (λόγω πυκνότητας)} \end{aligned}$$

Άρα καταλήξαμε ότι $|l(x)| < \|l\| \varepsilon$, για κάθε $\varepsilon > 0$, συνεπώς πρόκειται για τετριμμένη συνάρτηση $l(x) \equiv 0$.

\Leftarrow

Έστω ότι ο χώρος V δεν είναι πυκνός στον X .

Τότε υπάρχουν $x \in X$ και $\delta > 0$ τέτοια ώστε $\|x - v\| > \delta$ για κάθε $v \in V$.

Έστω χώρος W ορισμένος ως ένα $\text{span}\{x, V\}$.

Επειδή ο αριθμός a είναι μοναδικά ορισμένος εάν $a_1x + v_1 = a_2x + v_2$, έχουμε $(a_1 - a_2)x = v_2 - v_1$, οπότε αναγκαστικά θα έπρεπε να έχουμε $a_1 = a_2$ αφού $x \notin V$.

Ορίζουμε τότε, ένα γραμμικό συναρτησοειδές του χώρου W ως εξής:

$l(ax + v) = a$. Τότε έχουμε ότι $l(x) = 1$ και $l(v) = 0$, για κάθε $v \in V$.

Παρατηρούμε ότι:

- Για $a \neq 0$

$$\|ax + v\| = |a| \|x + a^{-1}v\|$$

$$\text{Αφού } \|x + a^{-1}v\| \geq \delta$$

$$(\text{γιατί } \|x - (-a^{-1}v)\| > \delta, \text{ επειδή } -a^{-1}v \in V) \Rightarrow$$

$$\geq |l(w)|\delta \Rightarrow$$

$$|l(w)| \leq \frac{\|w\|}{\delta}$$

- Για $a = 0$

$$|l(w)| = 0 \text{ και άρα η παραπάνω ανισότητα ισχύει.}$$

Η παραπάνω ανισότητα μας δείχνει ότι το l πρόκειται για ένα μη-τετριμμένο φραγμένο συναρτησοειδές του W .

Οπότε από τη βασική εφαρμογή του θεωρήματος Hahn- Banach το l μπορεί να επεκταθεί ως ένα φραγμένο γραμμικό συναρτησοειδές σε όλο το χώρο X .

Άτοπο. ■

Σύμφωνα με τις παραπάνω προσεγγίσεις για την πυκνότητα ενός υπόχωρου V ενός χώρου X , αρκούσε να δείξουμε ότι κάθε γραμμικό συναρτησοειδές του V ήταν το τετριμμένο. Δεδομένου της δυσκολίας αντιμετώπισης του παραπάνω ζητήματος για οποιονδήποτε χώρο X , θα δουλέψουμε για τους συνήθεις διανυσματικούς χώρους

που χρησιμοποιούνται και στην πράξη για τους οποίους μπορούμε να βρούμε μια διακριτή παρουσίαση του δυϊκού του χώρου X^* και της νόρμας του.

Πιο συγκεκριμένα μπορεί να δειχθεί (Kreyszig, 1978, σελ.227) ότι για την ειδική περίπτωση του χώρου $X = C[a, b]$, κάθε γραμμικό συναρτησοειδές μπορεί να γραφτεί ως :

$$l(f) = \int_a^b f(x)d(w)$$

Όπου w είναι μια συνάρτηση φραγμένης απόκλισης.

Συνεπώς για να αποδείξουμε την πυκνότητα και στην ειδική περίπτωση των μονοδιάστατων σιγμοειδών συναρτήσεων αρκεί να δείξουμε ότι αν το παραπάνω ολοκλήρωμα εκμηδενίζεται όποτε η f πρόκειται για σιγμοειδή συνάρτηση τότε αναγκαστικά η συνάρτηση w είναι σταθερή.

Θεώρημα 2.3.2 (Θεώρημα κυριαρχημένης σύγκλισης):

Έστω $f_n : X \rightarrow \mathbb{R}$ ακολουθία μετρήσιμων συναρτήσεων. Υποθέτουμε ότι $f_n \rightarrow f$ σχεδόν παντού και ότι υπάρχει $g : X \rightarrow [0, +\infty)$ ολοκληρώσιμη, ώστε για κάθε $n \in \mathbb{N}$, $|f_n| \leq g$ σχεδόν παντού. Τότε οι f_n και η f είναι ολοκληρώσιμες καθώς και $\lim_{n \rightarrow \infty} \int f_n d\lambda = \int f d\lambda$.

Ξεκινώντας, θεωρούμε την ποσότητα $\int_a^b \sigma(kx+l)d\omega(x) = 0$ για όλα τα $k, l \in \mathbb{Z}$ και υιοθετούμε την ιδέα των κλιμακωτών συναρτήσεων κάνοντας το k αρκετά μεγάλο.

Στην συνέχεια για κάθε $p, q \in \mathbb{Z}$ με $\frac{p}{q} \in [a, b]$, ορίζουμε τη συνάρτηση :

$$\begin{cases} 0, & a < p/q \\ \sigma(l), & x = p/q \\ 1, & p/q < x < b. \end{cases}$$

Μελετώντας την έκφραση $\sigma(nq(t - \frac{p}{q}))$, όταν το $n \rightarrow \infty$, η σχέση συγκλίνει σημειακά στο $r, r \in [a, b]$.

Έτσι από το θεώρημα κυριαρχημένης σύγκλισης του Lebesgue έχουμε ότι :

$$0 = \int_a^b r(x)dw(x) = \int_{P/q^+}^b dw(x) + \sigma(l)(w(P/q^+) - w(P/q^-)).$$

Ο τελευταίος όρος χαρακτηρίζει το ‘άλμα’ της συνάρτησης w στο σημείο P/q .

Στη συνέχεια παρατηρούμε ότι ο όρος του ολοκληρώματος δεν εξαρτάται από το l , έτσι στέλνοντας το $l \rightarrow -\infty$, εξ’ ορισμού της σιγμοειδούς συνάρτησης έχουμε ότι $\sigma(l) \rightarrow 0$.

Επομένως έχουμε ότι :

$$\int_t^s dw(x) = 0, \text{ όπου } s, t \text{ πρόκειται για ρητούς αριθμούς.}$$

Συνεπώς με βάση το θεώρημα κυριαρχημένης σύγκλισης του Lebesgue χρησιμοποιώντας ακολουθίες ρητών αριθμών για να συγκλίνει μονοτονικά στα απαιτούμενα τελικά σημεία καταλήγουμε ότι η σχέση ισχύει για όλα τα $t, s \in [a, b]$.

Οπότε το ολοκλήρωμα ισούται με $w(s) - w(t)$ δείχνοντας δηλαδή ότι πρόκειται για σταθερά. **(βλ. [5, σελ.162])** ■

Μολονότι η παραπάνω προσέγγιση καταλήγει στο συμπέρασμα του Θεωρήματος (2.1.1), μας περιορίζει καθώς δεν παράγει την εκτίμηση του όρου $\omega(f, 1/n)$, συνεπώς δεν είναι ξεκάθαρο πως τέτοιες εκτιμήσεις μπορούν να εξασφαλιστούν από την παραπάνω διαδικασία.

Στη συνέχεια θα εργαστούμε με την ιδέα της συνέλιξης (convolution). Πιο συγκεκριμένα η βασική ιδέα της συνέλιξης είναι η κατασκευή ενός πυρήνα βασισμένη στις σιγμοειδής συναρτήσεις με σκοπό την προσέγγιση της ιδιότητας αναπαραγωγής της γενικευμένης συνάρτησης Dirac, όπου τελικά αφού πλησιαστούν με τον τύπο τετραγωνισμού (πρόκειται για μια φόρμουλα προσέγγισης με σκοπό τον υπολογισμό ορισμένων ολοκληρωμάτων) θα παραχθεί η αναγκαία προσέγγιση.

Παρατήρηση 2.3.2: Η γενικευμένη συνάρτηση Dirac (συνάρτηση δέλτα) ορίζεται με βάση τις ιδιότητες:

$$i. \quad \delta(x) = \begin{cases} +\infty, & x = 0 \\ 0, & x \neq 0. \end{cases}$$

$$ii. \quad \int_{-\infty}^{+\infty} \delta(x)dx = 1.$$

Ξεκινάμε με τη βασική ιδέα αναπαραγωγής συναρτήσεων, όπου θα αναλυθεί στο παρακάτω θεώρημα. Όπως παρατηρήθηκε από τους Xu κ.α. (Stein και Weiss, 1971, σελ. 10) η συνάρτηση f που θα χρησιμοποιηθεί πρέπει να είναι ομοιόμορφα συνεχής στον \mathbb{R}^n . Επίσης καθιστάτε πιο φυσικό στη θεωρία προσέγγισης να συζητάμε για ένα συμπαγή υποσύνολο K έτσι πρέπει να επεκτείνουμε την συνάρτηση f από το K στον \mathbb{R}^n με έναν κατάλληλο τρόπο. Τέλος θα συμβολίσουμε το διαφορικό $dV(x)$ για να δηλώσουμε τα στοιχεία $dx_1 dx_2 \cdots dx_n$, καθώς και το συμβολισμό $\|\cdot\|$ για την κοινή ευκλείδεια νόρμα στον \mathbb{R}^n .

Επεκτείνοντας τον συντελεστή συνέχειας στον \mathbb{R}^n ορίζουμε παρακάτω ως:

$$\omega(f, \delta) = \sup_{\substack{x, y \in K \\ \|x-y\| < \delta}} |f(x) - f(y)| .$$

Στον παραπάνω ορισμό το σύνολο $K \subset \mathbb{R}^n$ δεν χρειάζεται να είναι συμπαγές, όμως αναγκαία η f πρέπει να είναι ομοιόμορφα συνεχής K .

Παρακάτω παρουσιάζονται τα αναγκαία θεωρήματα της συνέλιξης καθώς και μερικά αποτελέσματα αυτών που θα μας φανούν χρήσιμα σε πολλές περιπτώσεις.

Θεώρημα 2.3.3 (βλ. [5, σελ. 163]): Έστω συνάρτηση f φραγμένη και ομοιόμορφα συνεχής στον \mathbb{R}^n και συνάρτηση $g \in L_1(\mathbb{R}^n)$ με την ιδιότητα:

$$\int_{\mathbb{R}^n} g(x) dV(x) = 1.$$

Ορίζουμε $g_m(x) = m^n g(mx)$. Τότε:

- i. $f * g_m(x)$ συγκλίνει ομοιόμορφα στην f όταν $m \rightarrow \infty$
- ii. Για κάθε $R > 0$

$$\|f * g_m - f\|_\infty \leq \omega(f, \frac{2R}{m}) \|g\|_1 + 2\|f\|_\infty \int_{\|s\| > R} |g(s)| dV(s), \text{ όπου η } \|\cdot\|_\infty \text{ είναι πάνω στον } \mathbb{R}^n.$$

- iii. Ως μια εναλλακτική συνθήκη του ερωτήματος (ii), υποθέτουμε ότι η f είναι Lipschitz με σταθερά Λ και $M = \int_{\mathbb{R}^n} \|x\| |g(x)| dV(x) < \infty$

τότε:

$$\|f * g_m - f\|_\infty \leq \frac{M\Lambda}{m} .$$

Απόδειξη:

i. Ξεκινώντας από τη σχέση:

$$\int_{\mathbb{R}^n} g_m(x) dV(x) = \int_{\mathbb{R}^n} m^n g(mx) dV(x) = \int_{\mathbb{R}^n} g(mx) dV(mx) \Rightarrow$$

(Αφού το διαφορικό $dV(x)$ εκφράζει τον όγκο των στοιχείων $dx_1 dx_2 \cdots dx_n$ στον \mathbb{R}^n , έχουμε $dV(mx) = m \cdot m \cdots m \cdot dV(x) \equiv m^n dV(x)$)

Θέτοντας $y = mx$, έχουμε:

$$\Rightarrow \int_{\mathbb{R}^n} g(y) dV(y) = 1.$$

Στη συνέχεια παρατηρούμε ότι λόγω της συνέλιξης των συναρτήσεων f, g_m και του ολοκληρωτικού μετασχηματισμού ότι $((f * g)(x) \stackrel{\Delta}{=} \int_{\mathbb{R}^D} f(x-t)g(t)dt)$ ότι:

$$(f * g_m)(x) - f(x) = \int_{\mathbb{R}^n} (f(x-t) - f(x))g_m(t) dV(t).$$

Άρα:

$$\begin{aligned} |(f * g_m)(x) - f(x)| &\leq \int_{\mathbb{R}^n} |f(x-t) - f(x)| |g_m(t)| dV(t) = \\ &= \int_{\mathbb{R}^n} |f(x-t) - f(x)| |g(mt)| m^n dV(t) \\ &= \int_{\mathbb{R}^n} \left| f\left(x - \frac{s}{m}\right) - f(x) \right| |g(s)| dV(s) \quad (2.1) \\ &\leq \int_{\mathbb{R}^n} \omega\left(f, \frac{\|s\|}{m}\right) |g(s)| dV(s) \end{aligned}$$

Παρατηρούμε σε αυτό το σημείο ότι η ποσότητα $\omega\left(f, \frac{\|s\|}{m}\right)$ είναι μια ολοκληρώσιμη συνάρτηση, συναρτήσεως της f και συγκλίνει μονοτονικά στο 0 όταν το $m \rightarrow \infty$, συνεπώς από το θεώρημα μονοτονικής σύγκλισης το ολοκλήρωμα τείνει στο 0. Οπότε η συνάρτηση $(f * g_m)(x)$ συγκλίνει στην f όταν το $m \rightarrow \infty$.

ii. Για να οδηγηθούμε στο αποτέλεσμα του δεύτερου ερωτήματος σύμφωνα με τη σχέση (2.1) χωρίζουμε το ολοκλήρωμα στις περιοχές $\|s\| \leq R$ και $\|s\| > R$, όπου φράζουμε κάθε όρο.

iii. Αφού η f είναι Lipschitz με σταθερά Λ , θα ισχύει ότι:

$$|f(x) - f(y)| < \Lambda |x - y|, \text{ άρα από τη σχέση (2.1)}$$

\Rightarrow

$$\int_{\mathbb{R}^n} \left| f\left(x - \frac{s}{m}\right) - f(x) \right| |g(s)| dV(s) < \int_{\mathbb{R}^n} \Lambda \left| x - \frac{s}{m} - x \right| |g(s)| dV(s)$$

$$= \frac{\Lambda}{m} \int_{\mathbb{R}^n} \|s\| |g(s)| dV(s) = \frac{\Lambda}{m} M. \quad \blacksquare$$

Παρατήρηση 2.3.3: Με βάση το δεύτερο σκέλος του θεωρήματος, παρατηρούμε ότι το φράγμα δεν τείνει στο 0 όταν το $m \rightarrow \infty$. Για να επιτευχθεί κάτι τέτοιο θα έπρεπε η συνάρτηση g να έχει ένα συμπαγές στήριγμα και να επιλέγαμε τον όρο R κατάλληλα ώστε ο δεύτερος όρος να μηδενιστεί.

Εφοδιασμένοι πλέον με τα βασικά εργαλεία της ομοιόμορφης σύγκλισης των συνελιξέων των Xu και άλλων, αρχικά κατασκευάζουμε έναν πυρήνα συνέλιξης (με τον όρο πυρήνα συνέλιξης εννοούμε την επεξεργασία μιας συνάρτησης/πίνακα μέσω μιας άλλης που λέγεται πυρήνας) g της μορφής:

$$g(x) = \frac{1}{a_{n-1}} \int_{S^{n-1}} \varphi(x^T \cdot u) dS^{n-1}(u). \quad (2.2)$$

όπου S^{n-1} η μοναδιαία μπάλα στον \mathbb{R}^n (Σφαίρα πολλαπλότητας $n-1$) (Το σύνολο $\{u \in \mathbb{R}^n \mid \|u\| = 1\}$ και a_{n-1} το εμβαδόν της επιφάνειας που λαμβάνει από την τιμή του ολοκληρώματος με $\phi \equiv 1$).

Πάνω σε έναν κατάλληλο συμπαγή σύνολο K του \mathbb{R}^n , με K το ν -διάστατο υπόχωρο $K = [-a, a]^n$ για κάποια πραγματικά a .

Στη συνέχεια θεωρούμε ότι η $f \in C(K)$ και επεκτείνουμε την f στον \mathbb{R}^n με τέτοιο τρόπο ώστε να ισχύει $f(x) = 0$ όταν $x \in 2K$, όπου για κάθε $t > 0$ έχουμε:

$$tK = \{x \in \mathbb{R}^n \mid \frac{x}{t} \in K\}.$$

Αφού η f είναι συνεχής στο φραγμένο σύνολο $2K$ συμπεραίνουμε ότι είναι ομοιόμορφα συνεχής και φραγμένη στον \mathbb{R}^n .

Με βάση τα παραπάνω έχουμε το θεμελιώδη θεώρημα.

Θεώρημα 2.3.4 (βλ. [5, σελ. 164]): Έστω $K = [-a, a]^n$ για κάποια $a \in \mathbb{R}$ και έστω $\phi \in C(\mathbb{R})$ ομοιόμορφα συνεχής. Θεωρούμε επίσης τη συνάρτηση g ορισμένη όπως τη σχέση (2.2). Τότε:

Αν, η $g \in L_1(\mathbb{R}^n)$ και $\int_{\mathbb{R}^n} g(x)dV(x) \neq 0$ τότε το σύνολο των συναρτήσεων της μορφής $\phi(x^T \cdot a + c)$, $a \in \mathbb{R}^n$, $c \in \mathbb{R}$ είναι θεμελιώδης στον $C(K)$.

Απόδειξη: Έστω $f \in C(K)$, επεκτείνουμε την f ως μια φραγμένη, ομοιόμορφα συνεχή συνάρτηση του \mathbb{R}^n και θεωρούμε ότι το ολοκλήρωμα της g είναι 1.

Έστω $\varepsilon > 0$ και με βάση το θεώρημα (2.3.3), επιλέγουμε ένα m τέτοιο ώστε $\|f * g - f\|_\infty \leq \varepsilon/3$, όπου περιορίσαμε την $\|\cdot\|_\infty$ στο σύνολο K .

Στη συνέχεια με βάση τον τύπο τετραγωνισμού προσεγγίζουμε τη συνέλιξη, όπως παρακάτω:

$$(f * g_m)(x) = \int_{2K} g_m(x-y)f(y)dV(y) = \int_{2mK} g(mx-z)f\left(\frac{z}{m}\right)dV(z). \quad (2.3)$$

Έτσι αφού η προς ολοκλήρωση συνάρτηση είναι συνεχής και το σύνολο $2mK$ είναι συμπαγές, μπορούμε να την προσεγγίσουμε χρησιμοποιώντας αριθμητική ολοκλήρωση για οποιαδήποτε τιμή του x

Για να βρούμε ένα ομοιόμορφο φράγμα θα εφαρμόσουμε ένα απλό άθροισμα Riemann για την προσέγγιση.

Οπότε για κάθε $\delta > 0$ θεωρούμε P μια διαμέριση του συνόλου $2mK$ από πεπερασμένα, ξένα μεταξύ τους σύνολα Borel, με διάμετρο το πολύ δ . **(βλ. [5, σελ. 165])**

Για κάθε $A \in P$ επιλέγουμε στοιχείο $z_A \in A$ και ορίζουμε ως: $b_A = \int_A f\left(\frac{z}{m}\right)dV(z)$.

Έτσι:

$$\begin{aligned} & \left| \int_{2mK} g(mx-z)f\left(\frac{z}{m}\right)dV(z) - \sum_{A \in P} b_A g(mx-A) \right| \\ & \leq \sum_{A \in P} \int_A |g(mx-z) - g(mx-z_A)| \left| f\left(\frac{z}{m}\right) \right| dV(z) \\ & \leq \omega(g, \delta) \sum_{A \in P} \int_A \left| f\left(\frac{z}{m}\right) \right| dV(z) \end{aligned}$$

$$= \omega(g, \delta) m^n \int_{2K} |f(y)| dV(y).$$

Επομένως μπορούμε να επιλέξουμε δ και P ώστε το φράγμα να είναι μικρότερο του $\frac{\varepsilon}{3}$.

Στη συνέχεια εφαρμόζουμε όμοια επιχειρήματα στη σχέση (2.2).

Για κάθε $\theta > 0$ θεωρούμε Q μια διαμέριση του S^{n-1} ως μια πεπερασμένη, ξένη συλλογή από σύνολα Borel μέγιστης διαμέτρου θ .

Για κάθε $B \in Q$, θέτουμε:

$$c_B = \frac{1}{a_{n-1} B} \int dS^{n-1}(u) \text{ και επιλέγοντας } u_b \in B, \text{ καταλήγουμε στη σχέση:}$$

$$\left| g(mx - z_A) - \sum_{B \in Q} c_B \varphi((mx - z_A)^T u_B) \right| \leq \omega(\varphi, \|mx - z_A\| \theta)$$

Όμως $z_A \in 2mK$ οπότε $\frac{z_A}{m} \in 2K$, άρα για κάθε $x \in K$:

$$\|mx - z_A\| = m \left\| x - \frac{z_A}{m} \right\| \leq 3R, \text{ όπου } R \text{ η διάμετρος του } K.$$

Συνεπώς μπορούμε να επιλέξουμε θ, B κατάλληλα ώστε το δεξί μέρος της παραπάνω ανισότητας να είναι μικρότερο του $\frac{\varepsilon}{3}$, δηλαδή:

$$\left| \sum_{A \in P} b_A (g(mx - z_A) - \sum_{B \in Q} c_B \varphi((mx - z_A)^T u_B)) \right| \leq \frac{\varepsilon}{3}.$$

Τελικά συνδυάζοντας τις τρεις προσεγγίσεις για κάθε $x \in K$ έχουμε:

$$\left| f(x) - \sum_{A \in P} \sum_{B \in Q} b_A c_B \varphi(mx^T u_B - z_A^T u_B) \right| \leq \varepsilon.$$

Επομένως για να αποδείξουμε ότι οι συναρτήσεις $\phi(x^T \cdot a + c)$ είναι θεμελιώδεις πρέπει να δείξουμε ότι η συνάρτηση g ορισμένη από τη σχέση (2.2), ανήκει στο χώρο $L_1(\mathbb{R}^n)$ και δεύτερον πρόκειται για μη-μηδενικό ολοκλήρωμα.

Για την περάτωση του πρώτου ελέγχου, πρέπει να φράξουμε την $g(x)$, από μια κατάλληλη δύναμη του $r = \|x\|$, όπου το r είναι μεγάλο, λαμβάνοντας υπ'όψιν το πόσο γρήγορα τείνει η g στο 0.

Το επόμενο λήμμα μας δείχνει ποια δύναμη του r απαιτείται.

Λήμμα 2.3.1 (βλ. [5, σελ. 166]): Για $r = \|x\|$, $x \in \mathbb{R}^n$ και $q, R \in \mathbb{R}$, με $R > 0$ έχουμε:

$$\int_{\|x\| \geq R} r^{-q} dV(x) < \infty, \text{ αν και μόνο αν } q > n.$$

Απόδειξη: Συμβολίζοντας με S_r^{n-1} την σφαίρα ακτίνας r , του \mathbb{R}^{n-1} για $\rho > R$ έχουμε:

$$\int_{\rho \leq \|x\| \leq R} r^{-q} dV(x) = \int_{\rho}^R r^{-q} \int_{S_r^{n-1}} dS_r^{n-1} dr.$$

Αφού η σφαίρα S_r^{n-1} πρόκειται για $(n-1)$ -διάστατη πολλαπλότητα έχουμε ότι:

$$\int_{S_r^{n-1}} dS_r^{n-1} = a_{n-1} r^{n-1},$$

όπου a_{n-1} ο όγκος σφαίρας ακτίνας 1.

Συνεπώς:

$$\int_{\|x\| \geq R} r^{-q} dV(x) = a_{n-1} \int_R^{\rho} r^{n-1} r^{-q} dr = a_{n-1} \int_R^{\rho} r^{n-1-q} dr = a_{n-1} \left[\frac{r^{n-q}}{n-q} \right]_R^{\rho}, q \neq n.$$

Οπότε το όριο για $\rho \rightarrow \infty$ υπάρχει όταν $q > n$. Για την περίπτωση που $q = n$ θα

$$\text{είχαμε } a_{n-1} \int_R^{\rho} r^{-1} dr = a_{n-1} \ln\left(\frac{\rho}{R}\right) \xrightarrow{\rho \rightarrow \infty} \infty. \quad \blacksquare$$

Για να δείξουμε ότι η $g \in L_1(\mathbb{R}^n)$ αρκεί να δείξουμε ότι $g(x) = o(\|x\|^{-n})$ όταν $\|x\| \rightarrow \infty$, συγκλίνει δηλαδή πολυωνυμικά γρήγορα στο 0.

Στη συνέχεια ορίζουμε έναν πυρήνα επεξεργασίας της g ακτινικής μορφής.

Λήμμα 2.3.2 (βλ. [5, σελ. 167]): Έστω g ορισμένη ως :

$$g(x) = \frac{1}{a_{n-1}} \int_{S^{n-1}} \varphi(x^T \cdot u) dS^{n-1}(u).$$

Τότε: $g(x) = g_0(\phi, r)$ όπου $r = \|x\|$ και

$$g_0(\phi, r) = \frac{a_{n-2}}{a_{n-1}} \int_{-1}^1 \varphi(rs)(1-s^2)^{\frac{n-3}{2}} ds \quad (2.4)$$

$$= \frac{a_{n-2}}{a_{n-1}} \int_{-r}^r r^{2-n} \varphi(t)(r^2-t^2)^{\frac{n-3}{2}} dt, \quad r \neq 0. \quad (2.5)$$

Απόδειξη: Για την απόδειξη του λήμματος λειτουργούμε με $\|x\| \neq 0$, αφού οι σχέσεις (2.2), (2.4) που θα χρησιμοποιηθούν εξαρτώνται από το x .

Θεωρούμε ένα σύστημα συντεταγμένων με πόλο την κατεύθυνση του x .

Έστω w ένα μοναδιαίο διάνυσμα στην κατεύθυνση x , όπου $w = \frac{x}{r}$. Τότε κάθε

σημείο $u \in S^{n-1}$ μπορεί να εκφραστεί ως $u = w \cos \theta + v$, όπου $\cos \theta = \frac{u^T \cdot x}{r}$ και v

ένα μοναδιαίο διάνυσμα κάθετο στο x με $\|v\| = \sin \theta$.

Συνεπώς μπορούμε να καλύψουμε όλη τη σφαίρα S^{n-1} όταν το θ κυμαίνεται από 0 έως π και το v τρέχει σε όλες τις κατευθύνσεις ορθογώνια του x .

Άρα:

$$\begin{aligned} g(x) &= \frac{1}{a_{n-1}} \int_{S^{n-1}} \varphi(x^T \cdot u) dS^{n-1}(u) \\ &= \frac{1}{a_{n-1}} \int_0^\pi \varphi(r \cos \theta) \int_{S_{\sin \theta}^{n-2}} dS_{\sin \theta}^{n-2}(v) d\theta. \end{aligned}$$

Ο όρος $S_{\sin \theta}^{n-2}$ εκφράζει τη $(n-2)$ -διάστατη σφαίρα με ακτίνα $\sin \theta$, έτσι:

$$\int_{S_{\sin \theta}^{n-2}} dS_{\sin \theta}^{n-2}(v) d\theta = a_{n-2} \sin^{n-2}(\theta).$$

Οπότε μπορούμε να ορίσουμε τη $g_0(\phi, r)$, αφού η g πρόκειται για μια ακτινική συνάρτηση, ως εξής:

$$\begin{aligned}
g_o(\varphi, r) &= \frac{1}{a_{n-1}} \int_0^\pi \varphi(r \cos \theta) a_{n-2} \sin^{n-2}(\theta) d\theta \\
&= \frac{a_{n-2}}{a_{n-1}} \int_0^\pi \varphi(r \cos \theta) \sin^{n-3}(\theta) \sin(\theta) d\theta.
\end{aligned}$$

Θέτοντας $s = \cos \theta$, $ds = -\sin(\theta)d\theta$ και $\sin^{n-3}(\theta) = (1-s^2)^{\frac{n-3}{2}}$

Τέλος με την αντικατάσταση $t = rs$ η g_o θα πάρει τη μορφή:

$$\begin{aligned}
g_o(\varphi, r) &= \frac{a_{n-2}}{a_{n-1}} \int_r^{-r} \varphi(t) \left(1 - \frac{t^2}{r^2}\right)^{\frac{n-3}{2}} \frac{1}{r} (-1) dt \\
&= \int_{-r}^r \varphi(t) r^{2-n} (r^2 - t^2)^{\frac{n-3}{2}} dt, r \neq 0.
\end{aligned}$$

■

Στη συνέχεια πρέπει να επιλέξουμε κατάλληλα στη συνάρτηση ϕ ώστε η g να ανήκει στο χώρο $L_1(\mathbb{R}^n)$.

Η ουσιώδης προϋπόθεση είναι ότι η συνάρτηση ϕ πρέπει να τείνει γρήγορα στο 0, για $\pm\infty$.

Αρχικά, δεδομένου ότι η σιγμοειδής συναρτήσεις τείνουν στο 1 στο $+\infty$ ορίζουμε τη συνάρτηση $\varphi(t)$ όπως παρακάτω:

$$\phi(t) := \psi(t) = \sigma(1+t) + \sigma(1-t) - 1.$$

Λήμμα 2.3.3: Έστω $\sigma(x) = \frac{1}{1+e^{-x}}$ σιγμοειδής συνάρτηση. Τότε η συνάρτηση $\varphi(t) := \psi(t) = \sigma(1+t) + \sigma(1-t) - 1$ τείνει στο 0 στο $\pm\infty$ και πρόκειται για άρτια συνάρτηση.

Απόδειξη:
$$\psi(t) = \frac{1}{1+e^{-(1+t)}} + \frac{1}{1+e^{-(1-t)}} - 1.$$

Για $t \rightarrow +\infty$ ο πρώτος όρος τείνει στο 1, ενώ ο δεύτερος όρος τείνει στο 0. Συνεπώς,

$$\lim_{t \rightarrow +\infty} \psi(t) = 1+0-1=0. \text{ Ομοίως για } t \rightarrow -\infty \lim_{t \rightarrow -\infty} \psi(t) = 0+1-1=0.$$

Η συνάρτηση είναι άρτια, αφού για κάθε $t \in \mathbb{R}$, ισχύει ότι

$$\psi(t) = \psi(-t).$$

Επιπροσθέτως περιμένουμε ότι η συνάρτηση $\psi(t)$ τείνει γρήγορα στο 0, αφού αναμένουμε ότι η συνάρτηση σ προσεγγίζει την κλιμακωτή συνάρτηση.

Πιο συγκεκριμένα, θεωρούμε ότι η συνάρτηση σ είναι συνεχής και ισχύει:

$$|\psi(t)| \leq K \cdot |t|^{-q}, \quad q > n - 2 \text{ για κάποια } K \in \mathbb{R} \quad (2.6)$$

Παρατηρούμε ότι για την $\sigma(x) = \frac{1}{1+e^{-x}}$ η παραπάνω συνθήκη ισχύει για κάθε $q > 0$ και η $\psi(t)$ τείνει στο 0 εκθετικά.

Το επόμενο βήμα είναι να επεκτείνουμε τον πυρήνα της σχέσης (2.5). Συμβολίζουμε $\lambda = \frac{n-3}{2}$ και για $n > 2$ έχουμε:

Σε πρώτη φάση θα εργαστούμε για n περιττό, άρα το λ πρόκειται για μη-αρνητικό αριθμό.

Ο όρος του πυρήνα $(r^2 - t^2)^\lambda$ είναι ένα πολυώνυμο μεταξύ του r και του t άρα το ολοκλήρωμα της σχέσης (2.5) παίρνει τη μορφή:

$$g_o(\phi, r) = r^{2-n} \sum_{j=0}^{\lambda} \beta_j(r) r^{2\lambda-2j}. \quad (2.7)$$

$$\text{όπου, } \beta_j(r) = \frac{a_{n-2}}{a_{n-1}} {}^\lambda C_j \int_{-r}^r \phi(t) t^{2j} dt. \quad (2.8)$$

και ${}^\lambda C_j$ ο συνηθισμένος διωνυμικός συντελεστής: με ${}^\lambda C_j = \binom{\lambda}{j} \frac{\lambda!}{j!(\lambda-j)!}$

Από τον ορισμό της ϕ και την (2.6) έχουμε ότι όλα τα β_j συγκλίνουν όταν το $r \rightarrow \infty$.

Λήμμα 2.3.4 (βλ. [5, σελ. 168]): Έστω n περιττός και $\psi(t)$ ικανοποιεί τη συνθήκη (2.6) καθώς και g ορισμένη όπως η σχέση (2.2) με $\psi = \phi$. Τότε μια αναγκαία και ικανή συνθήκη της g για να ανήκει $L_1(\mathbb{R}^n)$ είναι:

$$\int_0^\infty \psi(t) t^{2j} dt = 0, \quad j = 0, \dots, \frac{n-3}{2}. \quad (2.9)$$

Απόδειξη: Για την περίπτωση όπου κάποια j δεν ικανοποιείται η παραπάνω σχέση, βλέπουμε από την (2.7), ότι όσο το $r \rightarrow \infty$ η $g_o(\phi, r)$ θα συμπεριφέρεται σαν την r^{-p} , όπου $p = n - 2 - 2\lambda + 2j \leq n - 2 < n$.

Συνδυάζοντας τις σχέσεις (2.7) και (2.9) έχουμε:

$$g_o(\varphi, r) = -r^{2-n} \sum_{j=0}^{\lambda} \gamma_j(r) r^{2\lambda-2j}, \quad (2.10)$$

$$\gamma_j(r) = \frac{a_{n-2}}{a_{n-1}} {}^{\lambda}C_j \int_{|t|>R} \psi(t) t^{2j} dt.$$

Στην παραπάνω σχέση έχουμε επιπλέον,

$$\begin{aligned} \gamma_j(r) &\leq \frac{a_{n-2}}{a_{n-1}} {}^{\lambda}C_j \int_{|t|>R} K |t|^{-q} t^{2j} dt \\ &= 2K \frac{a_{n-2}}{a_{n-1}} {}^{\lambda}C_j \int_r^{\infty} t^{2j-q} dt \\ &= -2K \frac{a_{n-2}}{a_{n-1}} {}^{\lambda}C_j r^{2j-q+1} \text{ αφού από υπόθεση } q > n-2 > 2j. \end{aligned}$$

Αντικαθιστώντας το παραπάνω φράγμα στη σχέση (2.10) βρίσκουμε ότι το g_o τείνει στο 0 τουλάχιστον τόσο γρήγορα όσο το r στη δύναμη $2-n-2\lambda-2j+2j-q-1 = 2-n+n-3-q+1 = -q$, $q > n$.

Άρα η $g(\|x\|)$ τείνει στο 0 ικανά γρήγορα. **(βλ. [5, σελ. 169])** ■

Τέλος απομένει να αποδείξουμε ότι η g δεν μπορεί να έχει μηδενικό ολοκλήρωμα. Ο Xu κ.α. έδειξαν το παρακάτω αποτέλεσμα.

Λήμμα 2.3.5 (βλ. [5, σελ. 169]): Έστω n περιττός και η $\psi(t)$ ικανοποιεί τις συνθήκες του λήμματος (2.3.3) και την σχέση (2.9). Θεωρούμε ότι η $\psi(t)$ είναι άρτια, τότε:

$$\int_{\mathbb{R}} g(x) dV(x) = -2a_{n-2} r_n \int_0^{\infty} \psi(t) t^{n-1} dt, \text{ όπου } r_n = \int_0^1 r(1-r^2)^{\frac{n-3}{2}} dr > 0.$$

Απόδειξη:

$$\begin{aligned} \int_{\mathbb{R}} g(x) dV(x) &= a_{n-1} \int_0^{\infty} r^{n-1} g_o(r) dr \quad (\text{Δουλεύοντας όπως στο λήμμα 2.3.1}) \\ &= a_{n-2} \int_0^{\infty} \int_{-r}^r \psi(t) (r^2 - t^2)^{\frac{n-3}{2}} dt dr \quad (\text{Από τη σχέση 2.5}) \end{aligned}$$

$$= 2a_{n-2} \int_0^{\infty} r \int_0^r \psi(t) (r^2 - t^2)^{\frac{n-3}{2}} dt dr \quad (\text{Αφού η } \psi(t) \text{ είναι άρτια})$$

$$= -2a_{n-2} \int_0^{\infty} r \int_r^{\infty} \psi(t) (r^2 - t^2)^{\frac{n-3}{2}} dt dr.$$

Η παραπάνω ισότητα προκύπτει από τη σχέση (2.9) του Λήμματος (2.3.3) σύμφωνα με την αντικατάσταση:

$$\int_0^{\infty} \psi(t) t^{2j} dt = \int_0^r \psi(t) t^{2j} dt + \int_r^{\infty} \psi(t) t^{2j} dt = 0$$

$$= -2a_{n-2} \int_0^{\infty} \psi(t) \int_r^{\infty} r (r^2 - t^2)^{\frac{n-3}{2}} dr dt. \quad (\text{Από το θεώρημα Fubini})$$

Παρατηρούμε ότι το εσωτερικό ολοκλήρωμα είναι μια σταθερά πολλαπλασιαζόμενη με μια δύναμη t^{n-1} , έτσι θέτοντας $t=1$ βλέπουμε ότι η τιμή του εσωτερικού ολοκληρώματος ταυτίζεται με τη σταθερά r_n της υπόθεσης, όπου είναι αυστηρά θετική, όπως και το εσωτερικό του ολοκληρώματος εκτός από τις τιμές που παίρνει για τα 0,1. ■

Συμπληρώνοντας, ο περιορισμός της $\psi(t)$ ώστε να πρόκειται για άρτια συνάρτηση μπορεί να αντιμετωπιστεί αντικαθιστώντας από τη σχέση $\psi(t) = \sigma(1+t) + \sigma(1-t) - 1$

Συνδυάζοντας όλα τα παραπάνω καταλήγουμε στο παρακάτω θεώρημα.

Θεώρημα 2.3.5 (βλ. [5, σελ. 170]): Έστω $\sigma \in C(\mathbb{R})$ και $\psi(t) = \sigma(1+t) + \sigma(1-t) - 1$, θεωρούμε n περιττό και το σύνολο $K = [-r, r] \subset \mathbb{R}^n$. Επίσης η $\psi(t)$ ικανοποιεί τη σχέση $|\psi(t)| \leq K |t|^{-q}$, $q > n - 2$ για κάποια $K \in \mathbb{R}$, και

$$\int_0^{\infty} \psi(t) t^{n-1} dt \neq 0.$$

Τότε το σύνολο των συναρτήσεων $\sigma(x^T \cdot a + c)$, $a \in \mathbb{R}^n$, $c \in \mathbb{R}$ είναι θεμελιώδης στο $C(K)$, όπου $\sigma(x) = \frac{1}{1 + e^{-x}}$.

Μια σύντομη περιγραφή θα παρουσιαστεί για την περίπτωση όπου ο αριθμός n είναι άρτιος. Παρατηρούμε ότι το λήμμα (2.3.3) ικανοποιείται για $j = 0, 1, \dots, \frac{n-2}{2}$. Στη συνέχεια τα επιχειρήματα της προσέγγισης είναι ουσιαστικά τα ίδια με μερικές προσαρμογές των μεθόδων.

Αρχικά, βλέπουμε ότι η επέκταση του πυρήνα που ορίσαμε σχέση (2.7), δεν είναι πλέον πεπερασμένη συνεπώς το άνω άκρο του αθροίσματος πρέπει να γίνει άπειρο. Έπειτα παρατηρούμε ότι η αυξανόμενη δύναμη του t που πρέπει να εκμηδενιστεί προκαλεί πρόβλημα. Η πρώτη δύναμη του t που δεν μηδενίζεται είναι για $t = n$, έτσι η πρώτη δύναμη του r που δεν τείνει σε μηδέν είναι η $(2 - n) + (-n - 3 + n) = -(n + 1)$.

3. Αριθμητική ανάλυση σε αλγορίθμους μάθησης

Στο κεφάλαιο αυτό, θα στρέψουμε την προσοχή μας στην αλγοριθμική σκοπιά των νευρωνικών δικτύων, βασισμένοι σε τεχνικές της αριθμητικής ανάλυσης και της γραμμικής άλγεβρας. Ειδικότερα θα ασχοληθούμε με τον delta rule algorithm, ένας από τους πιο κλασικούς αλγορίθμους μάθησης που υπάγεται στην οικογένεια οπισθοδιάδοσης καθώς και την αιτιολόγηση της γραμμικότητας ως μια προσέγγιση σε μη- γραμμικούς αλγορίθμους οπισθοδιάδοσης. Πιο συγκεκριμένα ο παραπάνω αλγόριθμος ικανοποιεί το ουσιώδες χαρακτηριστικό ότι επιτρέπουμε στα ζευγάρια των εισαγόμενων δεδομένων (x_j, y_j) να παρουσιάζονται διαδοχικά στο σύστημα μας (αφού, δεν θα μας δίνονται όλα τα δεδομένα εξ'αρχής). Ο παραπάνω περιορισμός έπεται από δύο λόγους, έναν φιλοσοφικό και έναν πρακτικό.

Ο πρώτος λόγος προέρχεται από τη διαδικασία μάθησης των θηλαστικών την έμπνευση δηλαδή των νευρωνικών δικτύων, τα οποία "μαθαίνουν" από παραδείγματα όπου παρουσιάζονται διαδοχικά, αφού δεν διατίθεται ολόκληρο το σετ των πληροφοριών από την αρχή.

Ο δεύτερος και σημαντικότερος λόγος, είναι ότι ο όγκος των δεδομένων πολλές φορές είναι πολύ μεγάλος, πράγμα που μπορεί να δημιουργήσει προβλήματα στο υλικολογισμικό (hardware) του συστήματος.

Περιγραφικά στα μοντέλα που θα ασχοληθούμε θα μας δίνεται ένα σύνολο σημείων $\{x_1, x_2, \dots, x_t\}$ στον \mathbb{R}^n , όπου για κάθε σημείο x_j , υπάρχει ένα επιθυμητό output $y_j \in \mathbb{R}$. (Στην περίπτωση μας, θεωρώντας ότι το δίκτυο μας παράγει ένα απλό output, μπορεί να θεωρηθεί και ως μονόμετρο). Ουσιαστικά το output y_j πρόκειται για τιμές μιας συνάρτησης f , όπου θέλουμε το δίκτυο μας να "μάθει", αναγνωρίζοντας μερικά γενικά χαρακτηριστικά της τοπολογίας των σημείων, ούτως ώστε όταν έχουμε μια άγνωστη τιμή x να προσεγγίζει την τιμή y . Αρχικά θα χρειαστούμε μια μονάδα μέτρησης σφάλματος για να καταλάβουμε την αποδοτικότητα του αλγορίθμου. Για την μοντελοποίηση των προβλημάτων μάθησης θα χρησιμοποιήσουμε τη μέθοδο των ελαχίστων τετραγώνων, διαδικασία που αποτελεί την πιο συνηθισμένη επιλογή, για την ακέραιη αποτελεσματικότητα και πληρότητα της.

Με βάση όλα τα παραπάνω σκοπός μας, είναι η παρακάτω μοντελοποίηση:

Έστω, η πραγματική τιμή εξόδου που θέλουμε για ένα δεδομένο σει παραμέτρων, είναι η $g(x)$. Δουλεύοντας στο ορισμένο μοντέλο perceptron ενός κρυφού στρώματος με παραμέτρους $a_j \in \mathbb{R}$, $w_j \in \mathbb{R}^n$ και $c_j \in \mathbb{R}$, $j = 1, \dots, k$ η συνάρτηση ορίζεται ως:

$$g(x) = \sum_{j=1}^k a_j \sigma(w_j^T \cdot x + c_j), \quad x \in \mathbb{R}^n$$

Σκοπός μας θα ήταν να ελαχιστοποιήσουμε την ποσότητα $\sum_{j=1}^t (y_j - g(x_j))^2$, $j = 1, \dots, t$ (αφού έχουμε t inputs), όπου στην πράξη σημαίνει θα επιζητήσουμε μια κοντινή εκτίμηση του ελαχίστου και όχι απαραίτητα την πραγματική της τιμή.

3.1 Αλγόριθμος Delta rule

Η ανάλυση του αλγορίθμου Delta rule, θα εξεταστεί στην περίπτωση του multilayer perceptron, όπου αντί για ένα output θα έχουμε ένα διάνυσμα εξαγόμενων δεδομένων και τα βάρη των συνάψεων θα σχηματίζουν έναν πίνακα. Για λόγους οικονομίας, θα συμβολίζουμε πλέον τα εισαγόμενα δεδομένα με σκοπό τη μάθηση ως x και τα output ως y , καθώς και θα παραβλέψουμε την οριακή τιμή (threshold), αφού θα χειριστούμε το πρόβλημα ως ένα πρόβλημα προσέγγισης.

Ο βασικός μας στόχος είναι να βρούμε έναν πίνακα (έστω πίνακα W) τέτοιον ώστε $y = Wx$ για κάθε ζεύγος διανυσμάτων (x, y) . Να βρούμε δηλαδή έναν πίνακα ο οποίος θα μετασχηματίζει κατάλληλα τα inputs με σκοπό την προσέγγισης στα επιθυμητά outputs.

Η βασική ιδέα των αλγορίθμων μάθησης, πηγάζει από την παρακάτω διαδικασία (διαδικασία εκμάθησης αλγορίθμου):

Έστω ένα σύνολο εισαγόμενων δεδομένων $\{x_1, x_2, \dots, x_t\}$ στον \mathbb{R}^n , όπου για κάθε x_i αντιστοιχεί ένα output y_i . Το σύστημα χρησιμοποιεί αυτά τα ζευγάρια με σκοπό την ανανέωση της εκτίμησης του επιθυμητού πίνακα $W \in \mathbb{R}^{n \times n}$, έτσι όταν ολοκληρωθεί η εκμάθηση το σύστημα θα χρησιμοποιήσει καινούργια δεδομένα x ώστε να προβλέψει τα κατάλληλα y .

Πιο συγκεκριμένα, θεωρούμε ότι ο πίνακας W ανανεώνεται μετά από κάθε μοτίβο εκμάθησης, έτσι η αλλαγή του πίνακα ορίζεται ως:

$$(\delta W)_{ij} = \eta (y_j - (Wx)_j) x_i$$

Όπου η είναι ο βαθμός εκμάθησης και $(Wx)_j$ εκφράζει το j -στοιχείο του Wx .

Σημείωση: Ο βαθμός εκμάθησης, πρόκειται για μια παράμετρο που χρησιμοποιείται στην εκμάθηση του δικτύου, έχοντας μια μικρή θετική συνεισφορά (συνήθως μεταξύ του 0 και του 1), που ελέγχει το πόσο γρήγορα το μοντέλο μας θα προσαρμοστεί στο πρόβλημα.

Ουσιαστικά καταλαβαίνουμε ότι, όταν ο όρος σφάλματος της παρένθεσης είναι θετικός, θα πρέπει να προστεθεί μια συνιστώσα του x σε κάθε σειρά του W , ώστε να αυξηθεί η εξαγόμενη τιμή του δικτύου για το συγκεκριμένο μοτίβο. Αντίστοιχα, όταν η τιμή είναι αρνητική θα αφαιρέσουμε μια συνιστώσα του x ώστε να μειώσουμε το output για τα δεδομένα x . Για λόγους απλούστευσης παρατηρούμε ότι δεν υπάρχει σύζευξη ανάμεσα στις γραμμές του W , δηλαδή η καινούργια γραμμή j εξαρτάται αποκλειστικά από την παλιά γραμμή j , συνεπώς μπορούμε να αφαιρέσουμε το δείκτη από τους όρους. Έτσι το output y_j γίνεται y και η j γραμμή του πίνακα W γίνεται w^T . Οπότε χωρίς βλάβη της γενικότητας και επιστρέφοντας στο perceptron με ένα output, έχουμε:

$$\delta w_i = \eta(y - w^T x)x_i .$$

Συνεπώς μετά από μια εναλλαγή από τα διανύσματα βαρών w_k δημιουργούμε την επανάληψη:

$$\begin{aligned} w_{k+1} &= w_k + \delta w_k \\ &= w_k + \eta(y - w_k^T x)x \\ &= (I - \eta x x^T)w_k + \eta y x . \end{aligned} \tag{3.1}$$

Σημείωση: Ο δείκτης k της εξίσωσης αναφέρεται στην k επανάληψη και όχι στο k στοιχείο.

Παρατήρηση 3.1.1: Από τη δεύτερη ισότητα των παραπάνω, παρατηρούμε τη χρήση του delta rule algorithm, την προσθήκη δηλαδή μιας κατάλληλης συνιστώσας στο μοτίβο των x στο στάδιο επανάληψης που βρίσκεται ο διανυσματικός πίνακας βαρών.

Παρακάτω θα παρουσιαστούν κάποια αποτελέσματα των επαναλήψεων.

Έστω $B = I - \eta x x^T \in \mathbb{R}^{n \times n}$.

Λήμμα 3.1.1: Ο πίνακας B έχει μόνο δύο ιδιοτιμές, $1 - \eta \|x\|^2$ που αντιστοιχεί στο ιδιοδιάνυσμα x και το 1 που αντιστοιχεί στον υπόχωρο διανυσμάτων ορθογώνιο του x . (Όπου $\| \cdot \|$ η ευκλείδεια νόρμα ($\|x\|_2^2 = \sum_{i=1}^n |x_i|^2$)).

Λήμμα 3.1.2: Έστω $0 \leq \eta \leq \frac{2}{\|x\|^2}$.

Τότε έχουμε πως: $\|B\|_2 = \rho(B) = 1$, όπου $\rho(B)$ η φασματική ακτίνα του B και η φυσική νόρμα που επάγει η διανυσματική $\|\cdot\|_2$

Λήμμα 3.1.3: Έστω ότι έχουμε t πρότυπα διανύσματα x_1, x_2, \dots, x_t . Θεωρούμε ότι το ανάπτυγμα του χώρου των εισαγόμενων δεδομένων, περιέχει n γραμμικά ανεξάρτητα inputs.

Για κάθε μοντέλο διανυσμάτων x_p , θα έχουμε ένα διαφορετικό πίνακα B_p , με $B_p = I - \eta x_p x_p^T$.

Έστω $\Lambda = B_t B_{t-1} \dots B_1$, τότε αν $0 < \eta < \frac{2}{\|x_p\|^2}$ για κάθε μοντέλο x_p και το x_p αναπτύσσεται, τότε $\|\Lambda\| < 1$.

Απόδειξη: Εξ' ορισμού υπάρχει διάνυσμα v , τέτοιο ώστε $\|\Lambda\| = \|\Lambda v\|$ και $\|v\| = 1$.

Έτσι $\|\Lambda\| = \|B_t B_{t-1} \dots B_1 v\| \leq \|B_t\| \|B_{t-1}\| \dots \|B_1 v\|$. Θεωρούμε τις περιπτώσεις:

- Αν $v^T x_1 \neq 0$, έχουμε $\|B_1 v\| < 1$, αφού η συνιστώσα του v στην κατεύθυνση του x μειώνεται. (σύμφωνα με το Λήμμα 3.1.1) και έτσι έχουμε $\|\Lambda\| = \|B_t B_{t-1} \dots B_1 v\| \leq \|B_t\| \|B_{t-1}\| \dots \|B_2\| = 1$.
- Αν $v^T x_1 = 0$, έχουμε $B_1 v = v$ (Με βάση το Λήμμα 3.1.1)
Έτσι $\|\Lambda\| = \|B_t B_{t-1} \dots B_2 v\|$, όπου μπορούμε να συνεχίσουμε τη διαδικασία αυτή μέχρι να εφαρμοστεί η πρώτη περίπτωση.

Τέλος παρατηρούμε ότι το v δεν μπορεί να είναι ορθογώνιο σε όλα τα x_p , αφού από υπόθεση αναπτύσσονται. ■

Ένας κοινός τρόπος εφαρμογής του delta rule, είναι να εφαρμόσουμε τα διανύσματα x_1, x_2, \dots, x_t με τη σειρά και να ξανά ξεκινήσουμε τη διαδικασία κυκλικά από το x_1 . Η αναπαράσταση ενός πλήρη κύκλου των μοντέλων (διανυσμάτων x_1, x_2, \dots, x_t) ονομάζεται epoch.

Λαμβάνοντας υπ' όψιν αυτή τη στρατηγική, η επανάληψη της σχέσης (3.1), παράγει το αποτέλεσμα:

$$w_{k+t} = \Lambda w_k + \eta h . \quad (3.2.a)$$

όπου $\Lambda = B_t B_{t-1} \cdots B_1$, ενώ

$$h = y_1 (B_t B_{t-1} \cdots B_2) x_1 + \cdots + y_{t-1} B_t x_{t-1} + y_t x_t . \quad (3.2.b)$$

Στα παραπάνω οι τιμές των y_p εκφράζουν τον στόχο y για το p -μοντέλο διανυσμάτων (x_p) και όχι για το p -σημείο του διανύσματος, καθώς και παρατηρούμε ότι τα B_s και τα h_s εξαρτώνται από τα x_s και το η και όχι από τη δεδομένη κατάσταση του πίνακα W .

Αφού η ποσότητα δw στο delta rule αναφέρεται στο σφάλμα των outputs, παίρνουμε ένα σταθερό σημείο της επανάληψης (3.1), μόνο όταν τα σφάλματα είναι μηδέν, πρακτικά όμως κάτι τέτοιο δεν συμβαίνει. Έτσι καταλαβαίνουμε ότι η επανάληψη στην πραγματικότητα δε συγκλίνει με τη συνήθη έννοια. Από την άλλη, έχουμε δείξει στο Λήμμα (3.1.2) ότι τα x_p αναπτύσσουν το χώρο των εισαγόμενων δεδομένων, έτσι για κατάλληλα μικρό η , έπεται $\|\Lambda\| < 1$.

Στη συνέχεια δίνουμε τον ορισμό της συστολικής απεικόνισης

Ορισμός 3.1.1: Έστω (X, d) μετρικός χώρος. Μια απεικόνιση $F : X \rightarrow X$ είναι συστολική απεικόνιση ή συστολή αν υπάρχει σταθερά Λ , με $0 < \Lambda < 1$ έτσι ώστε:

$$d(F(x), F(y)) \leq \Lambda d(x, y), \text{ για κάθε } x, y \in X .$$

Θεώρημα 3.1.1: Έστω (X, d) μετρικός χώρος και μια απεικόνιση $F : X \rightarrow X$ να είναι συστολή. Τότε το θεώρημα συστολικής απεικόνισης μας λέει ότι η F έχει μοναδικό σταθερό σημείο στον X .

Μπορούμε να θεωρήσουμε την απεικόνιση $F(w) = \Lambda w + \eta h$, όπου ικανοποιεί τη σχέση:

$$\|F(w) - F(v)\| = \|\Lambda(w - v)\| \leq \|\Lambda\| \|w - v\| .$$

Πρόκειται δηλαδή για μια συστολή, με παράμετρο συστολής $\|\Lambda\|$.

Επακόλουθο είναι με βάση το θεώρημα συστολικής απεικόνισης ότι η επανάληψη της σχέσης (3.2.a) έχει μοναδικό σταθερό σημείο. Έτσι αν υπάρχει πίνακας W που μηδενίζει όλα τα σφάλματα, θα μπορούσε να θεωρηθεί ως ένα σταθερό σημείο της σχέσης (3.1), οπότε και της σχέσης (3.2.a). Από την άλλη η σχέση (3.1) δεν έχει σταθερά σημεία και το σταθερό σημείο της (3.2.a) εξαρτάται από το η , συνεπώς ο πίνακας ορίζεται ως $w(\eta)$.

Στην οριακή κατάσταση, όσο η επανάληψη (3.1) τρέχει το μοντέλο των inputs, θα παραχθεί ένας κύκλος ορίων από διανύσματα w_k , που θα επιστρέψει στο $w(\eta)$ όταν ολοκληρωθεί ο κύκλος μετά από t πρότυπα.

Στη συνέχεια αφού το $w(\eta)$ πρόκειται για ένα σταθερό σημείο της (3.2.α) για να δείξουμε την εξάρτηση της σχέσης από το η , έχουμε:

$$w(\eta) = \Lambda(\eta)w(\eta) + \eta h(\eta) \quad (3.3)$$

Μια πιο λεπτομερής προσέγγιση για το $w(\eta)$ δίνεται παρακάτω:

Αρχικά θεωρούμε w^* το διάνυσμα βαρών που ελαχιστοποιεί την έκφραση του σφάλματος $\varepsilon^2 = \sum_{p=1}^t (y_p - w^T x_p)^2$.

$$(3.4)$$

(Το w^* είναι μοναδικό αφού το x_p αναπτύσσεται, περιέχει δηλαδή $\{x_1, x_2, \dots, x_t\}$ γραμμικά ανεξάρτητα διανύσματα).

Έστω ο πίνακας X όπου οι στήλες του είναι τα στοιχεία x_1, x_2, \dots, x_t και

$$L = XX^T = \sum_{p=1}^t x_p x_p^T. \quad (3.5)$$

Τότε το ελάχιστο w^* ικανοποιεί τις κανονικές εξισώσεις, δηλαδή:

$$\begin{aligned} Lw^* &= \sum_{p=1}^t y_p x_p \\ &= h(0), \text{ (προκύπτει από τη σχέση 3.2.b),} \end{aligned}$$

όπου w^* , x_p , $h(0)$, διανύσματα.

Παρατηρούμε ότι όλοι οι πίνακες B_s , τείνουν στον ταυτοτικό I όταν $\eta \rightarrow 0$

Από την άλλη σύμφωνα με τη σχέση (3.3), έχουμε

$$H(\eta)w(\eta) = h(\eta), \text{ όπου } H(\eta) = \frac{I - \Lambda}{\eta}.$$

Ακόμα από τη στιγμή που τα μοντέλα των αναπτύσσονται ο πίνακας L^{-1} υπάρχει.

Ορίζουμε τον αριθμό $\kappa(L) := \|L^{-1}\| \|L\|$, όπου από τη στιγμή που ο πίνακας L είναι συμμετρικός και θετικός, εξ' ορισμού ο αριθμός $\kappa(L)$ ισούται με το λόγο της μεγαλύτερης ως προς τη μικρότερη ιδιοτιμή.

Συνδυάζοντας τα παραπάνω ένα αποτέλεσμα ευστάθειας των γραμμικών εξισώσεων δεδομένου $\|L - H(\eta)\| < \frac{1}{\|L^{-1}\|}$ είναι:

$$\frac{\|w(\eta) - w^*\|}{\|w^*\|} \leq \frac{\kappa(L)}{1 - \|L^{-1}\| \|L - H(\eta)\|_2} \left(\frac{\|h(\eta) - h(0)\|}{\|h(0)\|} + \frac{\|L - H(\eta)\|}{\|L\|} \right).$$

Όμως,

$$\begin{aligned} \Lambda(\eta) &= \prod_{p=1}^t (I - \eta x_p x_p^T) \\ &= I - \eta \sum_{p=1}^t x_p x_p^T + o(\eta^2) \\ &= I - \eta L + o(\eta^2). \end{aligned}$$

Έτσι $H(\eta) = L + o(\eta)$ και βάσει την (σχέση 3.2.b), έχουμε:

$$h(\eta) = h(0) + o(\eta).$$

Συγκεντρώνοντας όλα τα παραπάνω έχουμε το παρακάτω θεώρημα:

Θεώρημα 3.1.2 (βλ. [5,σελ.180]): Έστω τα διανύσματα x_p ώστε $\mathbb{R}^n = \text{span}\{x_1, x_2, \dots, x_p\}$ και w^* ο πίνακας βαρών που ελαχιστοποιεί όλα τα σφάλματα των εξόδων με τη μέθοδο των ελαχίστων τετραγώνων. Αν ο αλγόριθμος delta rule εφαρμοστεί με ένα σταθερό η που ικανοποιεί το Λήμμα (3.1.3) τότε τα βάρη θα συγκλίνουν σε ένα όριο κυκλικά.

Έστω $w(\eta)$ οποιαδήποτε επανάληψη του κύκλου, τότε όσο $\eta \rightarrow 0$:

- $\|w(\eta) - w^*\| = o(\eta)$.
- Αν $\varepsilon(\eta)$ είναι η ρίζα της μέσης τετραγωνικής απόκλισης του $w(\eta)$ και ε^* αντιστοιχεί στο σφάλμα του w^* τότε:
 $\varepsilon(\eta) - \varepsilon^* = o(\eta^2)$.

3.2 Η μέθοδος epoch

Δεδομένου ότι έχουμε ένα σταθερό και πεπερασμένο σύνολο προτύπων x_p όπου $p = 1, \dots, t$ μια εναλλακτική στρατηγική είναι η ανανέωση των βαρών όταν όλος ο κύκλος των επαναλήψεων των προτύπων (epoch), έχει παρουσιαστεί και όχι μετά από κάθε διάνυσμα. Η στρατηγική αυτή, παράγει μια πιο ακριβή κατεύθυνση για τα σφάλματα ελαχίστων τετραγώνων και καλείται μέθοδος epoch.

Τα παραπάνω μας οδηγούν στην επανάληψη:

Έστω $X = [x_1 : x_2 : \dots : x_t]$

$$\begin{aligned} w_{k+1} &= w_k - \eta \sum_{p=1}^t (x_p x_p^T) w_k + \eta \sum_{p=1}^t (y_p x_p) \\ &= \Omega w_k - \eta \sum_{p=1}^t (y_p x_p). \end{aligned} \quad (3.6)$$

Όπου $\Omega = I - \eta XX^T = I - \eta L$.

Η παραπάνω σχέση είναι ισοδύναμη της σχέσης (3.2.α) αλλά όχι με την (3.1), αφού πρόκειται για έναν πλήρη κύκλο επαναλήψεων των προτύπων. Παρατηρούμε ότι δεν υπάρχει πρόβλημα στο κυκλικό όριο και πράγματι ένα σταθερό σημείο w^* θα ήταν ένα πραγματικό ελάχιστο, ελαχίστων τετραγώνων.

Τα παραπάνω γίνονται ξεκάθαρα, θέτοντας $w_{k+1} = w_k = w^*$ στην σχέση (3.6) και βλέπουμε ότι παίρνουμε την κανονική εξίσωση (σχέση 3.5) για το πρόβλημα ελαχίστων τετραγώνων.

Ωστόσο υπάρχει ένα πρόβλημα στη μέθοδο.

Αρχικά πρέπει να βρούμε τις ιδιοτιμές του πίνακα Ω . (Ο πίνακας $L = XX^T$ είναι συμμετρικός και θετικά ορισμένος αφού τα x_p αναπτύσσουν το χώρο (τα $\{x_1, x_2, \dots, x_t\}$ περιέχουν δηλαδή n γραμμικά ανεξάρτητα διανύσματα), οπότε θα έχουμε πραγματικές μη-αρνητικές ιδιοτιμές)

Οι ιδιοτιμές του Ω είναι οι $(1 - \eta) \times$ (τις ιδιοτιμές του L), όπου για έναν αυστηρά θετικό ορισμένο πίνακα όλες οι ιδιοτιμές πρέπει να είναι αυστηρά θετικές. Έτσι για η κατάλληλα μικρό, $\rho(\Omega) = \|\Omega\| < 1$.

Συνεπώς για ένα εισαγόμενο πρότυπο διανυσμάτων που αναπτύσσονται στο χώρο και για κατάλληλα μικρό η , η επανάληψη (3.6) συγκλίνει. Μολαταύτα η επιλογή του πόσο μικρό πρέπει να είναι το η που πρέπει να θέσουμε, εξαρτάται από πιο ακριβείς

προσεγγίσεις για το φάσμα του L και της νόρμας του Ω . Με βάση τα παραπάνω θα δούμε γιατί η μέθοδος epoch στην πραγματικότητα δεν δουλεύει πολύ καλά.

Έστω ο πίνακας $L = XX^T$, $X = [x_1 : x_2 : \dots : x_t]$ έχει ιδιοτιμές λ_j , $j = 1, \dots, n$ με $0 < \lambda_n < \lambda_{n-1} < \dots < \lambda_1 = \rho(XX^T) = \|XX^T\| = \|X^T\|^2$.

Παρατηρούμε ότι για μικρές τιμές του n ο πίνακας Ω είναι θετικά ορισμένος.

Οι ιδιοτιμές του Ω είναι:

$$(1 - \eta\lambda_1) \leq (1 - \eta\lambda_2) \leq \dots \leq (1 - \eta\lambda_n) \text{ και } \rho(\Omega) = \max\{|1 - \eta\lambda_1|, |1 - \eta\lambda_n|\}.$$

$$\text{Επειδή } \lambda_1 = \|X^T\|^2 = \max_{\|v\|=1} \|X^T v\| = \max_{\|v\|=1} v^T X X^T v \leq \sum_{p=1}^t \|x_p\|^2. \quad (3.7)$$

Από την άλλη μπορούμε να βρούμε ένα κάτω φράγμα, αντικαθιστώντας ένα συγκεκριμένο v στην έκφραση (3.7).

Για παράδειγμα, έχουμε για κάθε $k = 1, \dots, t$

$$\lambda_1 \geq \frac{1}{\|x_k\|} \cdot \left(\sum_{p=1}^t x_k^T x_p \right)^2 \geq \|x_k\|. \quad (3.8)$$

3.3 Γενίκευση σε μη- γραμμικά συστήματα

Όπως έχουμε ήδη σχολιάσει στο κεφάλαιο 1, η χρήση των νευρωνικών δικτύων περιορίζεται καθώς πολλά πρότυπα διανύσματα δεν είναι γραμμικά διαχωρίσιμα. Συνεπώς πρέπει να γενικεύσουμε τα συστήματα μας σε μη- γραμμικά για την αντιμετώπιση των προβλημάτων αυτών. Παρακάτω θα ορίσουμε τη μορφή του μη- γραμμικού delta rule algorithm, όπου από αυτόν πηγάζει και η ειδική περίπτωση του αλγορίθμου οπισθοδιάδοσης (backpropagation algorithm).

Παρατήρηση 3.3.1: Είναι ξεκάθαρο, ότι για περιπτώσεις μη- γραμμικών αναλύσεων, πρέπει να περιοριστούμε σε τοπικά αποτελέσματα του χώρου, αφού μια γενική συμπεριφορά θα ταίριαζε περισσότερο σε προσεγγίσεις δυναμικών συστημάτων και θεωρία ελέγχου.

Πιο συγκεκριμένα όπως έχουμε δει για τη γραμμική περίπτωση, η διάσταση των inputs και ο αριθμός των βαρών είναι η ίδια. (n όπως έχουμε ήδη ορίσει)

Στη μη- γραμμική περίπτωση θεωρούμε M το συνολικό αριθμό των βαρών και n τη διάσταση των inputs. Δηλαδή, τα εισαγόμενα πρότυπα x του δικτύου ανήκουν στον \mathbb{R}^n και το διάνυσμα των παραμέτρων w στον \mathbb{R}^M . Συνεπώς για ένα perceptron ενός

στρώματος με m εξόδους, ο πίνακας w θα είναι $m \times n$ διαστάσεων και άρα το $M = m \cdot n$. Καταλαβαίνουμε δηλαδή για το πολύ- στρωματικό (multilayer) perceptron, η διάσταση του πίνακα w είναι το καρτεσιανό γινόμενο, των πινάκων με τα βάρη του κάθε στρώματος.

Έτσι για ένα σύστημα με m outputs το δίκτυο θα υπολογίζει μια συνάρτηση $g : \mathbb{R}^M \times \mathbb{R}^n \rightarrow \mathbb{R}^m$.

Ορίζουμε ως $v = g(w, x)$, όπου $v \in \mathbb{R}^m$

Αρχικά εφοδιάζουμε τον \mathbb{R}^n , \mathbb{R}^M και \mathbb{R}^m με την ευκλείδεια νόρμα και θεωρούμε τα εισαγόμενα πρότυπα ως x_p , που αντιστοιχούν στα outputs v_p . Δηλαδή:

$$v_p = g(w, x_p).$$

Θεωρούμε ότι η g είναι Frechét διαφορίσιμη ως προς το w και συμβολίζουμε $D = D(w, x)$ τον $m \times M$ πίνακα των παραγώγων ως προς τη κανονική μας βάση.

Οπότε για ένα σταθερό διάνυσμα x και μια μικρή αλλαγή δw έχουμε από τον ορισμό της παραγώγου:

$$g(w + \delta w, x) = g(w, x) + D(w, x)\delta w + o(\|\delta w\|). \quad (3.9)$$

Από την άλλη, για δοσμένο w , που αντιστοιχεί σε ένα συγκεκριμένο πρότυπο διανυσμάτων x_p έχουμε το επιθυμητό output y_p και άρα το σφάλμα ε_p των προτύπων x_p που δίνεται ως: $\varepsilon_p^2 = (y_p - v_p)^T (y_p - v_p) = q_p^T q_p$. (3.10)

Επομένως το συνολικό σφάλμα δίνεται από το άθροισμα των σφαλμάτων όλων των προτύπων: $\varepsilon^2 = \sum_{p=1}^t \varepsilon_p^2$.

Ο βέλτιστος αλγόριθμος θα επιδίωκε, την ελαχιστοποίηση του σφάλματος ε^2 , παρόλα αυτά ο αλγόριθμος που εξετάζεται δεν κατασκευάζει “κατεύθυνση για το ε^2 ” αλλά μια “κατεύθυνση για το ε_p^2 ”. Έτσι, για μια μικρή αλλαγή στο q_p θα είχαμε:

$$\begin{aligned} \delta \varepsilon_p^2 &= (q_p + \delta q_p)^T (q_p + \delta q_p) - q_p^T q_p \\ &= 2\delta q_p^T q_p + \delta q_p^T \delta q_p. \end{aligned}$$

Αφού το y_p είναι σταθερό και $q_p = y_p - v_p$ έπεται ότι $\delta q_p = -\delta v_p = -D(w, x_p)\delta w + o(\|\delta w\|)$. (από τη σχέση 3.9)

Άρα, συνδυάζοντας τις παραπάνω σχέσεις έχουμε:

$$\begin{aligned}\delta \varepsilon_p^2 &= -2(D(w, x_p) \delta w)^T (y_p - g(w, x_p)) + o(\|\delta w\|) \\ &= -2\delta w^T (D(w, x_p))^T (y_p - g(w, x_p)) + o(\|\delta w\|).\end{aligned}$$

Συνεπώς, αγνοώντας τους όρους τάξης $o(\|\delta w\|)$ και για μια σταθερή μικρή αλλαγή δw η μεγαλύτερη μείωση του ε_p^2 επιτυγχάνεται θέτοντας:

$$\delta w = \eta (D(w, x_p))^T (y_p - g(w, x_p)).$$

Συγκρίνοντας τώρα τα παραπάνω με την απλή περίπτωση του γραμμικού perceptron με ένα output, όπου ο δεύτερος όρος της έκφρασης είναι βαθμωτός με $g(w, x_p) = w^T x_p$ και η παράγωγος είναι το διάνυσμα κλίσης, που αναπτύσσεται ως προς w , καταλήγουμε ότι η μέθοδος πρόκειται για μια γενίκευση του delta rule algorithm και άρα δεδομένου του k -οστού διανύσματος βάρους έχουμε την επαναληπτική μέθοδο:

$$\begin{aligned}w_{k+1} &= w_k + \delta w_k \\ &= w_k + \eta (D(w_k, x_p))^T (y_p - g(w_k, x_p)).\end{aligned}\tag{3.11}$$

Στη συνέχεια στρέφουμε την προσοχή μας σε μια περιοχή ενός τοπικού ελαχίστου w^* της μεθόδου των ελαχίστων τετραγώνων. (Δεν μπορούμε να περιμένουμε ένα ολικό ελάχιστο αφού το σφάλμα έχει πολλά τοπικά ελάχιστα, όπως είναι γνωστό από την περίπτωση οπισθοδιάδοσης)

Έτσι από τις σχέσεις (3.10) και (3.11) έχουμε για μια περιοχή του w^* :

$$\begin{aligned}w_{k+1} &= w_k + \eta (D(w_k, x_p))^T (y_p - g(w^*, x_p) - D(w^*, x_p)(w_k - w^*)) \\ &\quad + o(\|w_k - w^*\|) \\ &= (I - \eta (D(w_k, x_p))^T D(w^*, x_p)) \cdot w_k + \eta (D(w_k, x_p))^T \times \\ &\quad \times (y_p - g(w^*, x_p) + D(w^*, x_p)w^*) + o(\|w_k - w^*\|).\end{aligned}\tag{3.12}$$

Παρατηρούμε ότι ο πίνακας επανάληψης δεν είναι συμμετρικός παρόλα αυτά θα έπρεπε να γίνει για w_k κοντά στο w^* . Πιο συγκεκριμένα θεωρούμε ο $D(w, x)$ είναι Lipschitz συνεχής στο w^* , ομοιόμορφος πάνω στο χώρο των εισαγόμενων προτύπων x .

Οπότε από την σχέση (3.12) έχουμε:

$$w_{k+1} = (I - \eta(D(w^*, x_p))^T D(w^*, x_p)) \cdot w_k + \eta(D(w^*, x_p))^T \times \\ \times (y_p - g(w^*, x_p) + D(w^*, x_p)w^*) + o(\|w_k - w^*\|). \quad \blacksquare$$

Έστω τώρα ότι εφαρμόζουμε τα πρότυπα x_1, \dots, x_t κυκλικά, όπως τη γραμμική περίπτωση. Αν καταφέρουμε και αποδείξουμε ότι το γραμμικό κομμάτι της απεικόνισης $w_k \rightarrow w_{k+t}$ είναι συστολή, λόγω της συνέχειας ότι υπάρχει μια περιοχή του w^* όπου όλη η απεικόνιση είναι συστολή. (Αφού έχουμε πεπερασμένα πρότυπα)

Αρχικά παρατηρούμε ότι $D(w^*, x_p)^T D(w^*, x_p)$ είναι θετικός ή μηδέν. Έτσι για κατάλληλα μικρό η :

$$\|I - \eta D(w^*, x_p)^T D(w^*, x_p)\| \leq 1.$$

Έτσι αναλύοντας το χώρο των διανυσμάτων βαρών στο ανάπτυγμα των ιδιοδιανυσμάτων που αντιστοιχούν στις μηδενικές και μη-μηδενικές ιδιοτιμές αντίστοιχα, παρατηρούμε ότι οι χώροι αυτοί είναι ορθογώνια συμπληρώματα μεταξύ τους αφού ο πίνακας είναι συμμετρικός. Συνεπώς καταφέρνουμε στο δεδομένο χώρο να παράγουμε τη συστολή για

$$\eta < \frac{1}{\rho(D(w^*, x_p)^T D(w^*, x_p))}.$$

Άρα με παρόμοιο τρόπο του λήμματος (3.1.3) καταφέρνουμε τη συστολή των υποχώρων ανάμεσα σε κάθε πρότυπο διανυσμάτων να αναπτύσσει όλο το χώρο των βαρών.

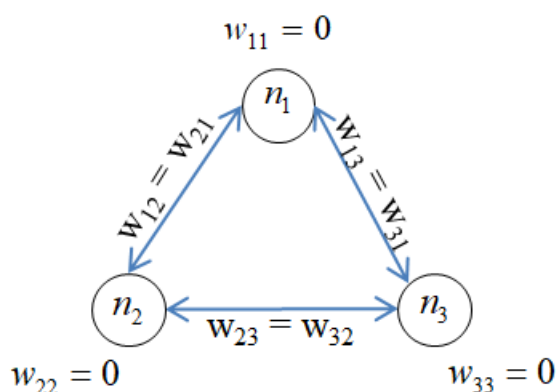
4 Αριθμητικές εφαρμογές των νευρωνικών δικτύων

Οι παραπάνω ιδιότητες προσέγγισης των νευρωνικών δικτύων δίνουν, τη δυνατότητα κατασκευής μοντέλων, με σκοπό την εφαρμογή τους σε πραγματικά προβλήματα. Εφαρμογές όπως, η κατηγοριοποίηση δεδομένων (έχοντας ένα σετ δεδομένων με φωτογραφίες από γάτες και σκύλους, το δίκτυο προβλέπει αν το input είναι γάτα ή σκύλος), εύρεση ανωμαλιών (με βάση τις συναλλαγές ενός ατόμου, το δίκτυο εντοπίζει αν η συναλλαγή είναι απάτη ή όχι), αναγνώριση προτύπων (εκπαιδευόμε τον αλγόριθμο ώστε να αναγνωρίζει εικόνες, ήχους κ.α., π.χ. Face id, Siri, Google assistant κ.α.), μη- γραμμικό έλεγχο (προβλέπει τη συμπεριφορά των δυναμικών συστημάτων που έχουμε ορίσει) κ.α.

Στο κεφάλαιο αυτό, θα ορίσουμε τα δίκτυα Hopfield και μέσω αυτών θα περάσουμε στο πρόβλημα του πλανόδιου πωλητή (Travelling salesman problem- TSP).

4.1 Δίκτυα Hopfield και βελτιστοποίηση γράφων

Το δίκτυο (net) Hopfield, πρόκειται για το πρώτο εφαρμοσμένο μη- γραμμικό δυναμικό νευρωνικό δίκτυο με εφαρμογή στη βελτιστοποίηση γραφημάτων. Σε αντίθεση με το multilayer perceptron, το Hopfield net δεν είναι στρωματικό, πράγμα που σημαίνει ότι κάθε νευρώνας μεταχειρίζονται ισάξια. Κάθε μονάδα είναι συνδεδεμένη με όλες τις άλλες, με αμφίδρομες συνδέσεις που σημαίνει ότι το βάρος της σύνδεσης από μια μονάδα i σε μια μονάδα j είναι $w_{ij} = w_{ji}$, καθώς και δεν επιτρέπονται συνδέσεις από μια μονάδα στον εαυτό της, δηλαδή $w_{ii} = 0$. Καταλαβαίνουμε δηλαδή ότι η τοπολογία των δικτύων Hopfield πρόκειται για ένα διάγραμμα με n κορυφές.



Σχήμα 7

Με βάση τη συνδεσμολογία των δικτύων κάθε κόμβος μπορεί να δεσμευτεί και από μεμονωμένα inputs καθώς και η έξοδος από κάθε μονάδα σχηματίζεται σε μια κατάσταση ενός διανύσματος x όπου το στοιχείο i του x αντιστοιχεί στην έξοδο i της μονάδας. Συνήθως, καθώς και για την ειδική περίπτωση του TSP το x πρόκειται για binary διάνυσμα (τα στοιχεία του είναι 0 ή 1) και αφού οι συνδέσεις του δικτύου μας είναι αναδρομικές θα εξελίσσονται χρονικά. Τέλος στο βασικό Hopfield net κάθε νευρώνας αναβαθμίζεται σε μια μονάδα χρόνου σε τυχαία σειρά ή κυκλικά. Παρακάτω θα ορίσουμε την αναδρομή.

Σχηματίζουμε το συμμετρικό πίνακα βαρών W και θεωρούμε q το διάνυσμα των inputs. Μετά από k χρονικές περιόδους έχουμε την κατάσταση του διανύσματος x_k , έτσι μετά από $k + 1$ χρονικές περιόδους, το input της μονάδος j είναι το j -στοιχείο

του πίνακα $Wx_k + q$. Η ενεργοποίηση της μονάδος συμβαίνει κατά την κορυφή της, θεωρούμε δηλαδή ως p_j την ακραία τιμή της j μονάδος και σχηματίζουμε το διάνυσμα p . Έτσι συμβολίζοντας το j -στοιχείο ενός διανύσματος y ως y_j αναπτύσσουμε τον j νευρώνα, αναβαθμίζοντας το $(x_{k+1})_j$:

$$(x_{k+1})_j = \begin{cases} 1, & (Wx_k + q - p)_j > 0 \\ 0, & (Wx_k + q - p)_j < 0 \\ (x_k)_j, & (Wx_k + q - p)_j = 0 \end{cases} \quad (4.1)$$

Στην περίπτωση που θέλουμε να αναβαθμίσουμε μια απλή μονάδα, εφαρμόζουμε την παραπάνω διαδικασία για το συγκεκριμένο j , αφήνοντας τα υπόλοιπα στοιχεία αμετάβλητα, καθώς και αν θέλουμε να κάνουμε καθολική αναβάθμιση του νευρώνα το εφαρμόζουμε για όλα τα j .

Παρατηρούμε ότι η κατάσταση του διανύσματος x του δικτύου είναι ένα ενεργειακό συναρτησοειδές:

$$E(x) = -\frac{1}{2} x^T W x + q^T x - p^T x. \quad (4.2)$$

Βλέπουμε, ότι παρόλο τη συσχέτιση του δικτύου ως ένα ενεργειακό συναρτησοειδές, ο πίνακας W δεν είναι θετικά ορισμένος. Επίσης, η συνθήκη $w_{ii} = 0$ μας δίνει ότι το $\text{trace}(W) = 0$, έτσι ο W πρέπει να έχει αρνητικές ιδιοτιμές. Συνεπώς το ελάχιστο της E δεν υπάρχει ως ένα μη-στάσιμο (σταθερό) σημείο, αλλά σε κάποια κορυφή του υπερκύβου που ορίζεται από τη συνθήκη $0 \leq (x)_j \leq 1$, $j = 1, \dots, n$, πράγμα που σημαίνει ότι πρέπει να περιοριστούμε σε ένα τοπικό ελάχιστο.

Η διαφορά της ενέργειας μας δίνει:

$$\begin{aligned} \delta E &= E(x_{k+1}) - E(x_k) \\ &= -(Wx_k + q - p)^T (x_{k+1} - x_k) - \frac{1}{2} (x_{k+1} - x_k)^T W (x_{k+1} - x_k). \end{aligned} \quad (4.3)$$

Θεωρούμε το εσωτερικό γινόμενο $(Wx_k + q - p)^T (x_{k+1} - x_k)$, όπου ίσως για μερικά j η ποσότητα $(x_{k+1} - x_k)_j = 0$, συνεπώς δεν θα υπάρχει συμμετοχή του j όρου στο εσωτερικό γινόμενο. Ακόμα για την επανάληψη πρέπει να υπάρχει ένα ή περισσότερα j_s τέτοια ώστε $(x_{k+1} - x_k)_{j_s} \neq 0$. Έτσι έχουμε δύο κλάσεις αποτελεσμάτων.

Αρχικά αν, $(Wx_k + q - p) > 0$ από την (4.1) έχουμε $(x_{k+1})_j = 1$, έτσι αναγκαστικά $(x_k)_j = 0$ και $(x_{k+1} - x_k)_j = 1$ καταλαβαίνουμε δηλαδή ότι ο j όρος σε αυτή την περίπτωση δίνει μια θετική τιμή στο εσωτερικό γινόμενο. Αντίστοιχα αν $(Wx_k + q - p) < 0$ τότε $(x_{k+1})_j = 0$ και $(x_k)_j = 1$ οπότε $(x_{k+1} - x_k)_j = -1$ και άρα ασκείται πάλι μια θετική συμμετοχή. Συνεπώς ο πρώτος όρος της (4.3) παρέχει μια μείωση στην ενέργεια εκτός αν $(x_{k+1} - x_k)_j = 0$. Στην περίπτωση αυτή δεν πραγματοποιούμε αλλαγές στην συνάρτηση της ενέργειας αλλά προχωράμε στον επόμενο νευρώνα. Επίσης αφού ο πίνακας W δεν είναι θετικά ορισμένος δεν εγγυόμαστε ότι ο δεύτερος όρος της (4.3) είναι αρνητικός.

4.2 Το πρόβλημα του πλανόδιου πωλητή (TSP)

Στις εφαρμογές βελτιστοποίησης τα Hopfield nets δεν είναι εκπαιδευμένα, συνεπώς πρέπει να ορίσουμε ένα ενεργειακό συναρτησοειδές που θα αντιστοιχεί στη συνάρτηση που πρέπει να βελτιστοποιήσουμε, καθώς και συναρτήσεις ποινών για κάθε περιορισμό και έναν πίνακα βαρών που θα ορίζει τη συνάρτηση σαν αυτή της σχέσης (4.2). Ορίζοντας τη διαδικασία κατά αυτόν τον τρόπο, ξεκινάμε το δίκτυο περιμένοντας να οδηγηθούμε στο ελάχιστο.

Θα αναπτύξουμε την παραπάνω διαδικασία στο πρόβλημα του πλανόδιου πωλητή (TSP), που πρωτοπαρουσιάστηκε από τους Hopfield και Tank.

Έστω ότι έχουμε m πόλεις, σκοπός του αλγορίθμου είναι να επισκεφτούμε κάθε πόλη ακριβώς μια φορά, επιστρέφοντας στην αρχική μας θέση καλύπτοντας την ελάχιστη απόσταση.

Έστω d_{ij} η απόσταση μεταξύ της i και j πόλης, με $d_{ij} = d_{ji}$ και δημιουργούμε ένα Hopfield net κατανέμοντας m νευρώνες και κάθε πόλη. (Το δίκτυο θα έχει $n = m^2$ νευρώνες).

Ταξινομούμε, τους νευρώνες σε ένα ταμπλό με m γραμμές και m μονάδες, όπου κάθε γραμμή αντιστοιχεί σε μία πόλη. Έτσι για μια επιτυχημένη διαδρομή πρέπει να υπάρχει ακριβώς μια μονάδα (1) σε κάθε γραμμή και αφού δεν γίνεται να βρισκόμαστε σε δύο πόλεις ταυτόχρονα πρέπει να υπάρχει ακριβώς μια μονάδα (1) και σε κάθε στήλη. Οπότε πρέπει να υπάρχουν m μονάδες στο ταμπλό για μια επιτυχημένη διαδρομή.

Για να κατασκευάσουμε το ενεργειακό συναρτησοειδές, ορίζουμε τις συναρτήσεις σφαλμάτων για κάθε ένα από τους περιορισμούς του προβλήματος. Έστω το διάνυσμα κατάστασης του δικτύου x , διανεμημένο από τις γραμμές του πίνακα, έτσι ώστε τα πρώτα m στοιχεία αναφέρονται στην 1^η πόλη, τα επόμενα m στοιχεία στην

2^η πόλη κ.ο.κ. Από τα παραπάνω βλέπουμε ότι το διάνυσμα x , περιλαμβάνει μόνο 1 και 0 και εν' συνεχεία ορίζουμε το διάνυσμα e που έχει μόνο 1.

Αρχικά εισάγουμε τον όρο:

$$(e^T x - m)^2 = x^T (ee^T)x - 2mx + m^2 .$$

Και παρατηρούμε ότι ο όρος μηδενίζεται αν και μόνο αν, έχει ακριβώς m μονάδες. Ο τετραγωνικός όρος m^2 μπορεί να αγνοηθεί αφού πρόκειται για σταθερά και βλέπουμε ότι ο πίνακας ee^T δεν έχει μηδενική διαγώνιο.

Έτσι θεωρούμε τον πίνακα: $A = ee^T - I$.

και άρα η ποσότητα $x^T (ee^T)x = x^T Ax + x^T x$.

Επίσης αφού το διάνυσμα x περιέχει μόνο (0) και (1), βλέπουμε ότι $x^T x = e^T x$, άρα η σχέση μας γράφεται ως:

$$(e^T x - m)^2 = x^T Ax + (1 - 2m)e^T x . \quad (4.4)$$

Οι υπόλοιποι περιορισμοί είναι περίπλοκοι καθώς πρέπει να βρεθεί ο τρόπος με τον οποίο θα γίνει η αντιστοίχιση μέρους τους διανύσματος κατάστασης, με συγκεκριμένες πόλεις.

Έστω πίνακας B που ορίζεται ως εξής:

- $b_{ii} = 0, i = 1, \dots, n$
- $b_{ij} = 0$, αν $(x)_i$ και $(x)_j$ εκπροσωπούν διαφορετικές πόλεις. Μαθηματικά μοντελοποιείται ως:

$\left[\frac{i}{m} \right] \neq \left[\frac{j}{m} \right]$, όπου $[r]$ (ακέραιο μέρος του r) εκφράζει τον μεγαλύτερο ακέραιο, αυστηρά μικρότερο από το r .

- $b_{ij} = 1$, αν $(x)_i$ και $(x)_j$ εκπροσωπούν την ίδια πόλη, δηλαδή $\left[\frac{i}{m} \right] = \left[\frac{j}{m} \right]$, για $i \neq j$.

Από την κατασκευή του B βλέπουμε, ότι είναι συμμετρικός, διαγώνιος και κάθε στοιχείο της διαγωνίου αντιστοιχεί σε μια συγκεκριμένη πόλη.

Η τετραγωνική μορφή είναι: $x^T Bx$. (4.5)

Με την παρακάτω ιδιότητα:

Κάθε block ή υποπίνακας του B , πολλαπλασιάζετε μεταξύ τους αθροίζοντας τις διακεκριμένες θέσεις του x που αντιστοιχούν στην ίδια πόλη και κατά συνέπεια η τιμή που θα πάρουμε θα είναι μεγαλύτερη του 0 αν δύο στοιχεία αντιστοιχούν στην ίδια πόλη. Τέλος με αυτή τη διαδικασία αθροίζουμε για όλες τις πόλεις και παρατηρούμε ότι η σχέση (4.4) θα πάρει την ελάχιστη τιμή (0), αν κάθε γραμμή του πίνακα των μονάδων έχει το πολύ ένα 1 (κάθε πόλη θα επισκέπτεται το πολύ μια φορά).

Στην συνέχεια πρέπει να βεβαιωθούμε ότι δεν θα επισκεφτούμε δύο πόλεις σε ένα βήμα, δηλαδή οι στήλες του πίνακα των μονάδων μας θα έχει ένα (1).

Το παραπάνω περιγράφεται ως: $x^T Cx$. (4.6)

Όπου,

$$C_{ii} = 0, C_{ij} = 1 \text{ αν } i \bmod m = j \bmod m \text{ για } i \neq j$$

και $C_{ij} = 0$ αλλιώς.

Ακόμα χρειαζόμαστε για τον αλγόριθμο έναν όρο που θα μειώνει την απόσταση των διαδρομών. Ουσιαστικά, πρέπει να ελέγξουμε την σειρά των διαδρομών που ακολουθούνται από τον πίνακα μας και να αθροίσουμε τις αποστάσεις, ορίζοντας τη διαδικασία αυτή σε τετραγωνική μορφή. Για κάθε πόλη έχουμε μια μονάδα στον πίνακα που μας λέει πότε επισκεφτήκαμε την πόλη αυτή, καθώς και η μονάδα της προηγούμενης στήλης μας δείχνει από πού ήρθαμε. Από τη στιγμή που μπορεί να υπάρχει το πολύ μια μονάδα στην προηγούμενη στήλη, μπορούμε να πολλαπλασιάσουμε την αντίστοιχη απόσταση με το στοιχείο του x που αντιστοιχεί σε αυτή την είσοδο. Έπειτα πολλαπλασιάζουμε τον όρο που προκύπτει με το στοιχείο του x που αντιστοιχεί στην πόλη που βρισκόμαστε και αθροίζουμε. Για να εξασφαλίσουμε την συμμετρία, πρέπει να κοιτάμε ταυτόχρονα την προηγούμενη και την επόμενη στήλη. Με τον τρόπο αυτό, κάθε σκέλος της διαδρομής θα έχει συμπεριληφθεί δύο φορές, παρόλα αυτά δεν μας επηρεάζει αφού δεν θα υπολογίσουμε την παραπάνω ποσότητα αλλά τη χρησιμοποιούμε για να βρούμε τα βάρη του δικτύου.

Τέλος, πρέπει να συμπεριλάβουμε την επιστροφή στην αρχική μας πόλη.

Η τετραγωνική μορφή θα είναι: $x^T Fx$. (4.7)

Ο πίνακας F έχει m^2 blocks ($m \times m$ blocks) και αντιστοιχούμε τους δείκτες p, q του υποπίνακα F_{pq} ως τη διαδρομή από την πόλη p στην πόλη q .

Τότε, $F_{pp} = 0$.

Για $p \neq q$ κάθε στοιχείο του πίνακα F_{pq} θα ισούται με d_{pq} ή 0, όπου τα μη-αρνητικά στοιχεία υπάρχουν, αν όντως έχει γίνει η διαδρομή από την πόλη p στην πόλη q . Έτσι πρέπει να τοποθετήσουμε τις αποστάσεις d_{pq} στην υπό και υπέρ διαγώνιο του F_{pq} . (Υπό- διαγώνιος είναι η πρώτη διαγώνιος κάτω από την κύρια διαγώνιο του πίνακα και Υπέρ- διαγώνιος είναι η πρώτη διαγώνιος πάνω από την κύρια διαγώνιο)

Επίσης, να προσθέσουμε την πιθανότητα ότι η πόλη p θα πρόκειται για την πρώτη πόλη και η πόλη q να είναι η τελευταία και ανάποδα.

Άρα θεωρούμε $(F_{pq})_{1m} = (F_{pq})_{m1} = d_{pq}$ και αφού τοποθετήσουμε τα στοιχεία στον πίνακα F_{pq} με τις παραπάνω διαδικασίες τα υπόλοιπα στοιχεία είναι όλα μηδενικά.

Τέλος, επιλέγουμε παραμέτρους $\alpha \gg \beta$, $\gamma \gg \delta$ και έτσι συγκρίνοντας τη σχέση (4.2) με τις σχέσεις (4.4-4.7), θέτουμε:

$$W = \alpha A + \beta B + \gamma C + \delta F \text{ και } q = (1 - 2m)e^T.$$

Το δίκτυο θα βρει μια διαδρομή με n μονάδες, προκαλώντας τον όρο της (4.4) να εξαφανιστεί. Έπειτα, θα βρει μια μονάδα σε κάθε γραμμή και στήλη εξαλείφοντας τους όρους των (4.5) και (4.6). Στη συνέχεια θα προσπαθήσει να μειώσει την απόσταση της διαδρομής και τέλος θα ορίσουμε ένα όριο (threshold) στον αλγόριθμο (π.χ. $\frac{1}{2}$).

Πίνακας εικόνων

Εικόνα εξωφύλλου: Συγκρότημα νευρωνικών συνδέσεων από τον: Max krasnov

Σχήμα 1 : Exclusive – OR δίκτυο

Σχήμα 2 : Ένα απλό perceptron

Σχήμα 3 : Πολύ επίπεδο νευρωνικό δίκτυο ⁽¹⁾

Σχήμα 4 : Πολύ επίπεδο νευρωνικό δίκτυο ⁽²⁾

Σχήμα 5 : Διαμέριση του \mathbb{R}^2 από πολύγωνα

Σχήμα 6 : Πολύ επίπεδο νευρωνικό δίκτυο ⁽³⁾

Σχήμα 7 : Ένα απλό Hopfield net, με τρία inputs

Βιβλιογραφία

- [1] Argyrakis P., “Neural Networks”, (2010), Computational Physics Group A.U.Th
- [2] Argyros S., “Real Analysis”, (2011), National Technical University of Athens, Greece
- [3] Argyros S., “Functional Analysis”, (2004), National Technical University of Athens, Greece
- [4] Costarelli D., Spigler R., Sigmoidal Functions Approximation and Applications, Department of Mathematics and Physics of Roma Tre University, Italy
- [5] Ellacott S.W., Aspects of the numerical analysis of neural networks, (1994), University of Brighton, England.
- [6] Eugene Isaacson, Herbert Bishop Keller, Analysis of Numerical Methods, (1994), New York University.
- [7] Hoel Paul G., Sidney C. Port and Charles J. Stone, Introduction to Probability Theory, (1971), University of California, Los Angeles.
- [8] Jost J., The Banach fixed point Theorem. The Concept of Banach Space, (2003), Springer, Berlin
- [9] Kurt Hornik, Multilayer Feedforward Networks are Universal Approximators, (1989), University of California, San Diego.
- [10] Simon Haykin, Neural Network A Comprehensive Foundation, (1994), McMaster University, Canada.
- [11] Will Light, Ridge Functions, Sigmoidal Functions and Neural Networks, (1992), University of Leicester, England.
- [12] Ανάργυρος Φελλούρης, Γραμμική άλγεβρα και αναλυτική γεωμετρία, (2009), EMI.

