

Υπολογιστική Νοημοσύνη
-Δεύτερη Εργασία-

Αλέξανδρος Πετρίδης

Τελευταία ενημέρωση: 8 Μαρτίου 2022

Περιεχόμενα

1	Εφαρμογή σε απλό dataset	3
1.1	Διαδικασία εκπαίδευσης και αξιολόγησης μοντέλων	3
1.2	Ζητούμενα προβλήματος	3
1.3	Σχολιασμός αποτελεσμάτων	7
2	Εφαρμογή σε dataset με υψηλή διαστασιμότητα	7
2.1	Ζητούμενα προβλήματος	8
2.2	Σχολιασμός αποτελεσμάτων	10

Κατάλογος Σχημάτων

1	Τελική μορφή ασαφών συνόλων του μοντέλου TSK_model_1.	3
2	Τελική μορφή ασαφών συνόλων του μοντέλου TSK_model_2.	4
3	Τελική μορφή ασαφών συνόλων του μοντέλου TSK_model_3.	4
4	Τελική μορφή ασαφών συνόλων του μοντέλου TSK_model_4.	5
5	Διαγράμματα μάθησης σφάλματος συναρτήσει του αριθμού των επαναλήψεων του μοντέλου.	5
6	Διαγράμματα μάθησης σφάλματος συναρτήσει του αριθμού των επαναλήψεων του μοντέλου.	6
7	Διαγράμματα αποτύπωσης των σφαλμάτων πρόβλεψης του μοντέλου.	6
8	Διαγράμματα αποτύπωσης των σφαλμάτων πρόβλεψης του μοντέλου.	7
9	Μέσο Σφάλμα ως συνάρτηση τον αριθμό των χαρακτηριστικών και το μέγεθος της ακτίνας	8
10	Προβλέψεις και πραγματικές τιμές μοντέλου.	9
11	Καμπύλη εκμάθησης.	9
12	Τρία ασαφή σύνολα στην αρχική και τελική τους μορφή (περιέχουν 12 κανόνες).	10

Κατάλογος Πινάκων

1	Ταξινόμηση μοντέλων προς εκπαίδευση.	3
2	Δείκτες απόδοσης.	7
3	Δείκτες απόδοσης βέλτιστου μοντέλου	10

1 Εφαρμογή σε απλό dataset

1.1 Διαδικασία εκπαίδευσης και αξιολόγησης μοντέλων

Αφού επιλέχθηκε το repository Airfoil Self-Noise dataset από το UCI repository, πραγματοποιήθηκαν τα παρακάτω:

- **Διαχωρίστηκε σε σύνολα εκπαίδευσης-επικύρωσης-ελέγχου:** όπου δημιουργήθηκαν τρία μη επικαλυπτόμενα υποσύνολα, από τα οποία το πρώτο θα χρησιμοποιηθεί για εκπαίδευση, το δεύτερο για επικύρωση και αποφυγή του φαινομένου υπερεκπαίδευσης και το τελευταίο για τον έλεγχο της απόδοσης του τελικού μοντέλου. Χρησιμοποιήθηκαν: το 60% για το υποσύνολο εκπαίδευσης και από 20% για κάθε ένα από τα εναπομείναντα υποσύνολα.
- **Εκπαιδεύτηκαν TSK μοντέλα με διαφορετικές παραμέτρους:** συγκεκριμένα, εκπαιδεύτηκαν 4 TSK μοντέλα με την υβριδική μέθοδο, σύμφωνα με την οποία οι παράμετροι των συναρτήσεων συμμετοχής βελτιστοποιήθηκαν μέσω της μεθόδου back propagation, ενώ οι παράμετροι της πολυωνυμικής συνάρτησης εξόδου βελτιστοποιήθηκαν μέσω της μεθόδου ελάχιστων τετραγώνων. Οι συναρτήσεις συμμετοχής ήταν bell-shaped και η αρχικοποίηση τους έγινε με τέτοιον τρόπο ώστε τα διαδοχικά ασαφή σύνολα παρουσίασαν σε κάθε είσοδο, βαθμό επικάλυψης περίπου 0.5. Ακόμα η μορφή εξόδου και το πλήθος των συναρτήσεων συμμετοχής για κάθε μεταβλητή εισόδου μεταβάλλονταν όπως φαίνεται στον παρακάτω πίνακα.

Μοντέλο	Πλήθος συναρτήσεων συμμετοχής	Μορφή εξόδου
TSK_model_1	2	Singleton
TSK_model_2	3	Singleton
TSK_model_3	2	Polynomial
TSK_model_4	3	Polynomial

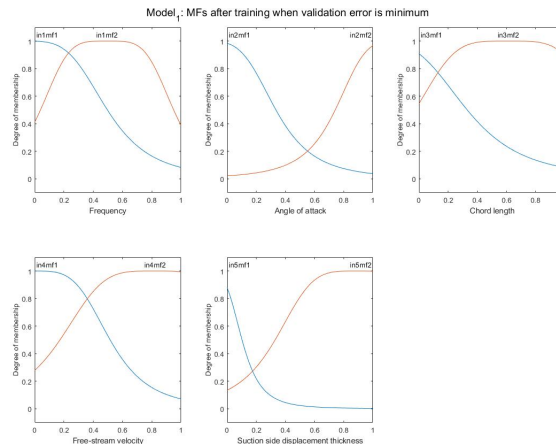
Πίνακας 1: Ταξινόμηση μοντέλων προς εκπαίδευση.

- **Αξιολογήθηκαν τα μοντέλα:** με τους δείκτες αξιολόγησης: μέσου τετραγωνικού σφάλματος (MSE), ρίζα μέσου τετραγωνικού σφάλματος (RMSE), συντελεστή προσδιορισμού R^2 , NMSE και NDEI όπως αυτά περιγράφονται στην εκφώνηση της εργασίας.

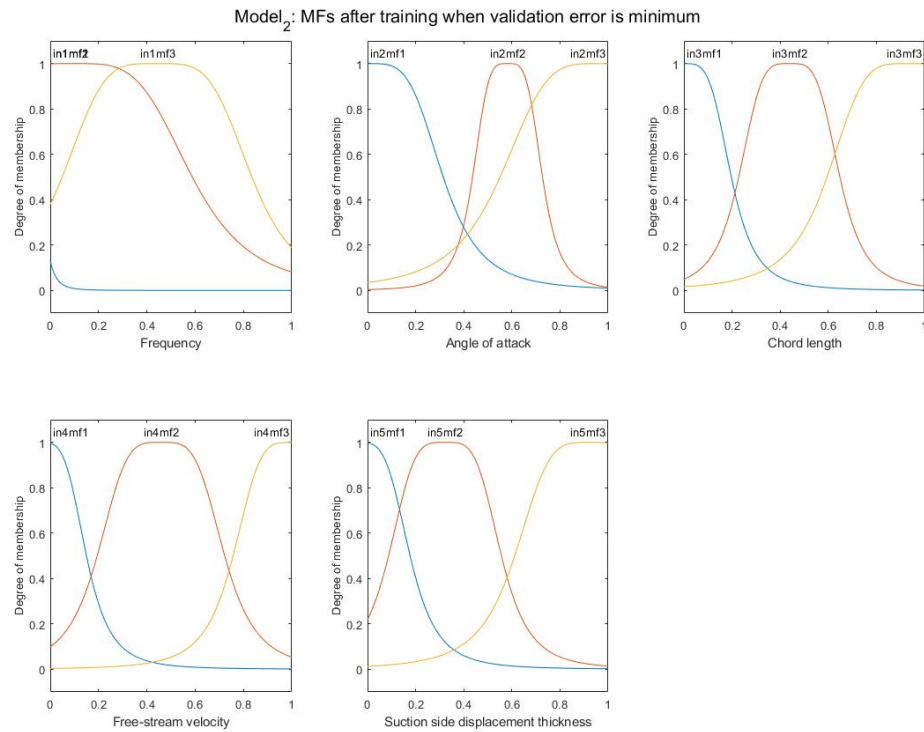
1.2 Ζητούμενα προβλήματος

Για κάθε ένα από τα 4 TSK μοντέλα που περιγράφονται στον παραπάνω πίνακα και τις κατάλληλες αρχικοποιήσεις έγινε εκπαίδευση των μοντέλων με τις παραμέτρους που περιγράφηκαν παραπάνω. Ως τελικό μοντέλο επιλεγόταν πάντα εκείνο το οποίο αντιστοιχούσε στο μικρότερο σφάλμα στο σύνολο επικύρωσης. Για τις τέσσερις περιπτώσεις εκπαίδευσης υπολογίστηκαν τα παρακάτω:

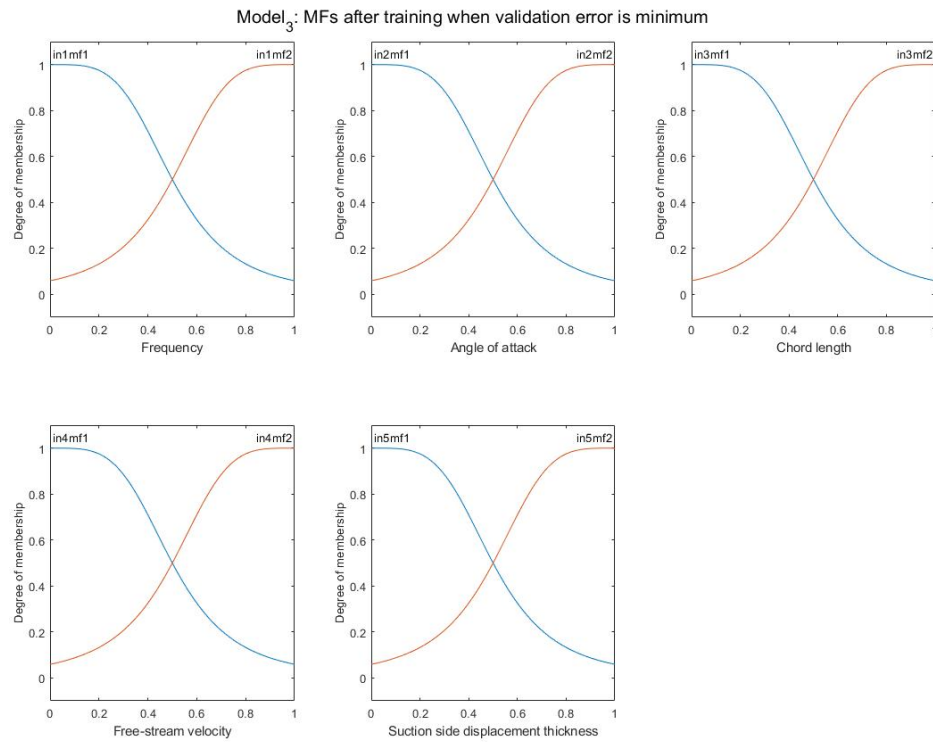
1. Διαγράμματα στα οποία απεικονίζονται οι τελικές μορφές των ασαφών συνόλων που προέκυψαν μέσω της διαδικασίας εκπαίδευσης.



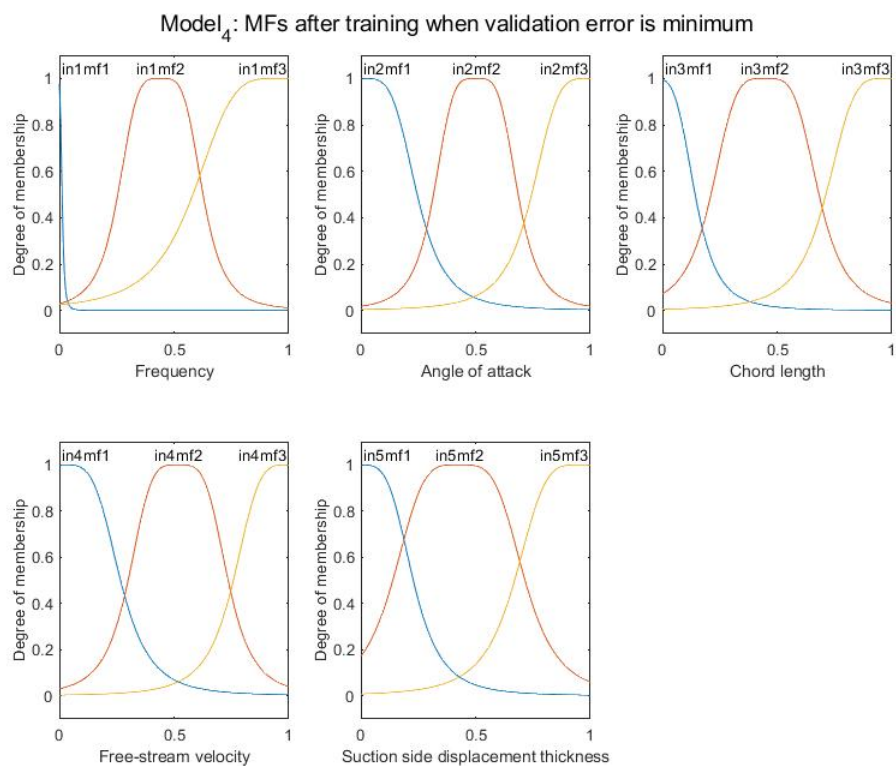
Σχήμα 1: Τελική μορφή ασαφών συνόλων του μοντέλου TSK_model_1.



Σχήμα 2: Τελική μορφή ασαφών συνόλων του μοντέλου TSK_model₂.

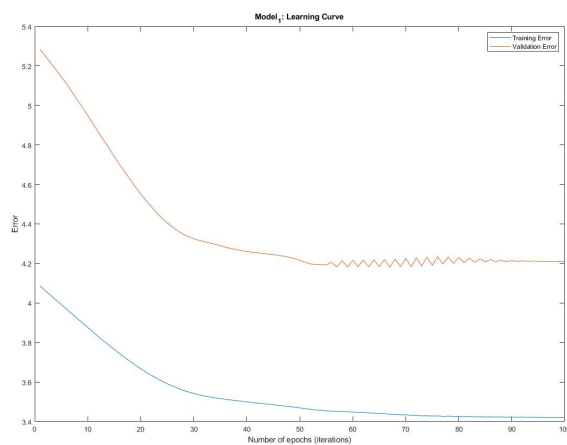


Σχήμα 3: Τελική μορφή ασαφών συνόλων του μοντέλου TSK_model₃.

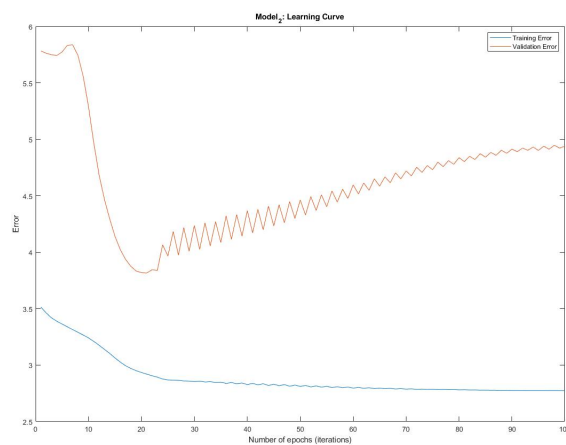


Σχήμα 4: Τελική μορφή ασαφών συνόλων του μοντέλου TSK_model4.

2. Διαγράμματα μάθησης (learning curves) όπου απεικονίζεται το σφάλμα του μοντέλου συναρτήσει του αριθμού των επαναλήψεων (iterations).

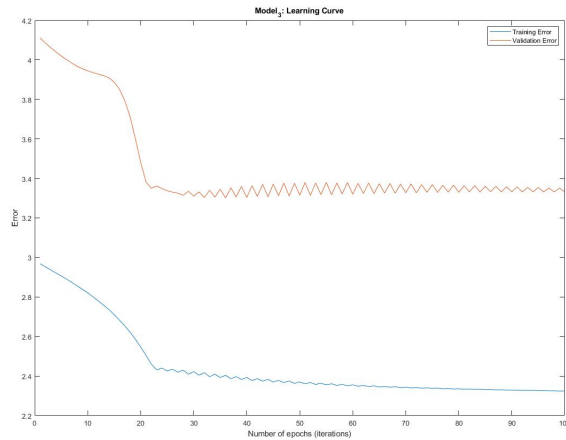


Μοντέλο TSK_model1.

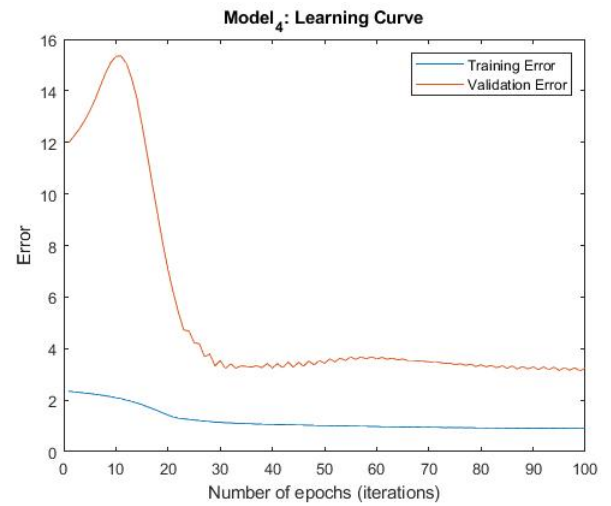


Μοντέλο TSK_model2.

Σχήμα 5: Διαγράμματα μάθησης σφάλματος συναρτήσει του αριθμού των επαναλήψεων του μοντέλου.



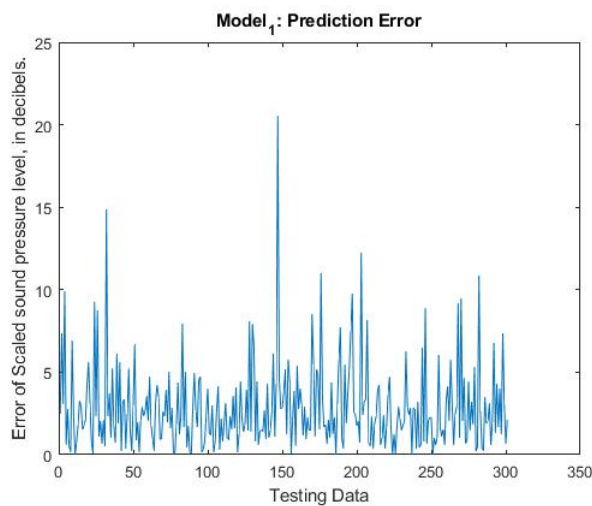
Μοντέλο TSK_model.3.



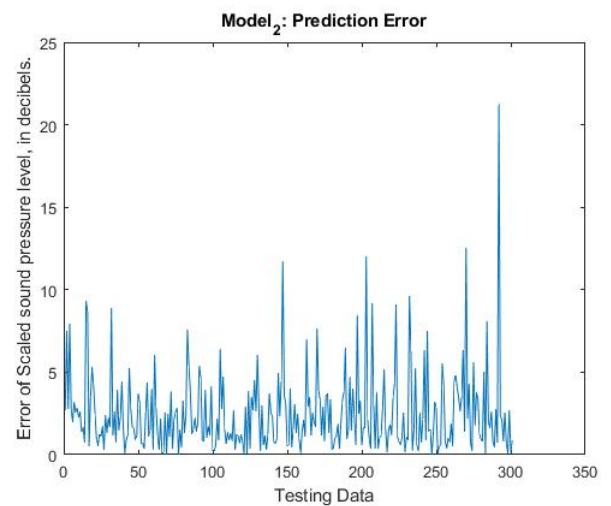
Μοντέλο TSK_model.4.

Σχήμα 6: Διαγράμματα μάθησης σφάλματος συναρτήσει του αριθμού των επαναλήψεων του μοντέλου.

3. Διαγράμματα αποτύπωσης των σφαλμάτων πρόβλεψης.

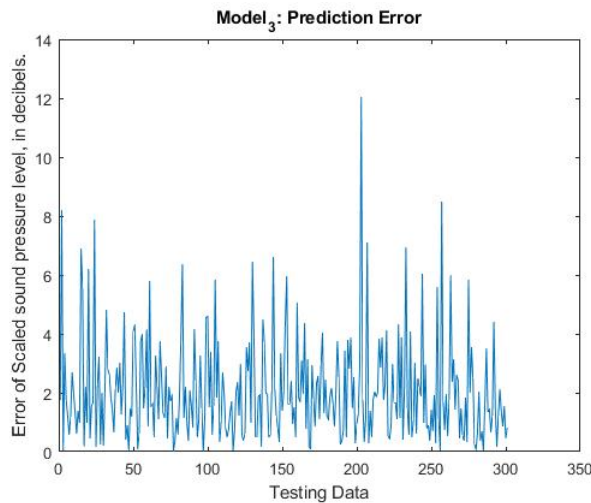


Μοντέλο TSK_model.1.

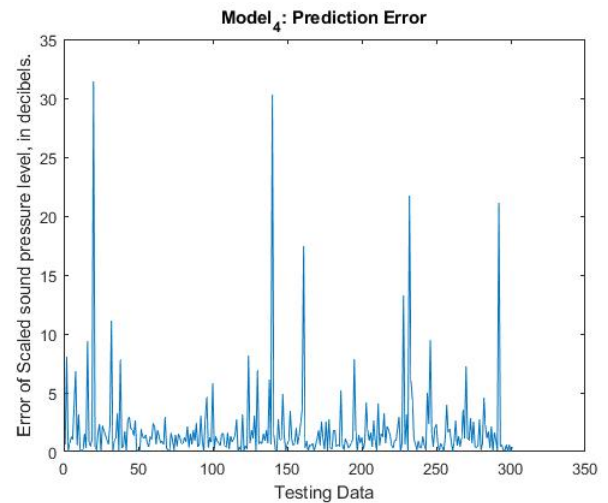


Μοντέλο TSK_model.2.

Σχήμα 7: Διαγράμματα αποτύπωσης των σφαλμάτων πρόβλεψης του μοντέλου.



Μοντέλο TSK_model.3.



Μοντέλο TSK_model.4.

Σχήμα 8: Διαγράμματα αποτύπωσης των σφαλμάτων πρόβλεψης του μοντέλου.

4. Δείκτες απόδοσης στον παρακάτω πίνακα.

Μοντέλο	RMSE	NMSE	NDEI	R^2
TSK_model.1	3.839718	0.337895	0.581287	0.662105
TSK_model.2	3.500550	0.280838	0.529941	0.719162
TSK_model.3	2.693165	0.166230	0.407713	0.833770
TSK_model.4	3.996081	0.365975	0.604959	0.634025

Πίνακας 2: Δείκτες απόδοσης.

1.3 Σχολιασμός αποτελεσμάτων

Με μια πρώτη ανάλυση διαπιστώνεται ότι το τρίτο μοντέλο είναι το καλύτερο καθώς είχε τον μεγαλύτερο δείκτη απόδοσης στα δεδομένα διάγνωσης αλλά και τις χαμηλότερες τιμές στις τρεις πρώτες μετρικές οι οποίες βασίζονται στο σφάλμα πρόβλεψης.

Επιπροσθέτως παρατηρώντας τις καμπύλες εκμάθησης, μπορεί να εξαχθεί το συμπέρασμα ότι πετυχαίνουμε πολύ γρηγορότερη σύγκλιση στα μοντέλα με τις τρεις συναρτήσεις συμμετοχής σε σχέση με αυτά που έχουν δύο.

2 Εφαρμογή σε dataset με υψηλή διαστασιμότητα

Το dataset που χρησιμοποιήθηκε στην παρούσα ενότητα ήταν το Superconductivity από το UCI Repository . Το μέγεθος του αριθμού των μεταβλητών των δεδομένων καθιστά αναγκαία τη χρήση μεθόδων μείωσης της διαστασιμότητας και του αριθμού των IF-THEN κανόνων καθώς με 81 όπως υπάρχουν στο dataset διαστάσεις και αν ο χώρος εισόδου κάθε μεταβλητής διαμεριζόταν με δύο ασαφή σύνολα, τότε θα καταλήγαμε με 2^{81} κανόνες.

Οπότε ο στόχος επίλυσης του προβλήματος θα επιτευχθεί μέσω της επιλογής χαρακτηριστικών και της χρήσης διαμέρισης διασχορπισμού. Οι δύο αυτές μέθοδοι εισάγουν στο πρόβλημα δύο ελεύθερες παραμέτρους, συγκεκριμένα, τον αριθμό των χαρακτηριστικών προς επιλογή και τον αριθμό των ομάδων που θα δημιουργηθούν.

Αναλυτικά έγινε η παρακάτω μοντελοποίηση του προβλήματος:

1. **Διαχωρίστηκε το σύνολο δεδομένων σε σύνολα εκπαίδευσης-επικύρωσης-ελέγχου:** όπως και στην πρώτη ενότητα.
2. **Επιλέχθηκαν βέλτιστες παράμετροι:** Η δημοφιλέστερη μέθοδος μέσω της οποίας επιτυγχάνεται αυτό είναι η αναζήτηση πλέγματος. Συγκεκριμένα αφού λήφθηκε ένα σύνολο τιμών για κάθε παράμετρο, δημιουργήθηκε ένα 2-διάστατο πλέγμα όπου κάθε σημείο αντιστοιχεί σε μία 2-άδα τιμών για τις εν λόγω παραμέτρους, και σε κάθε σημείο χρησιμοποιήθηκε μια μέθοδος αξιολόγησης για τον έλεγχο της ορθότητας των συγκεκριμένων τιμών.

Η μέθοδος που χρησιμοποιήθηκε ήταν η διασταυρωμένη επικύρωση (cross validation). Σύμφωνα με αυτή τη μέθοδο, και για επιλεγμένες τιμές των παραμέτρων, χωρίστηκε το σύνολο εκπαίδευσης σε δύο υποσύνολα, από τα οποία το ένα

χρησιμοποιήθηκε για την εκπαίδευση του μοντέλου και το άλλο για την αξιολόγησή του. Η διαδικασία αυτή επαναλήφθηκε πέντε φορές, με κάθε φορά διαφορετικό διαχωρισμό του συνόλου εκπαίδευσης και στο τέλος λήφθηκε ο μέσος όρος του σφάλματος του μοντέλου.

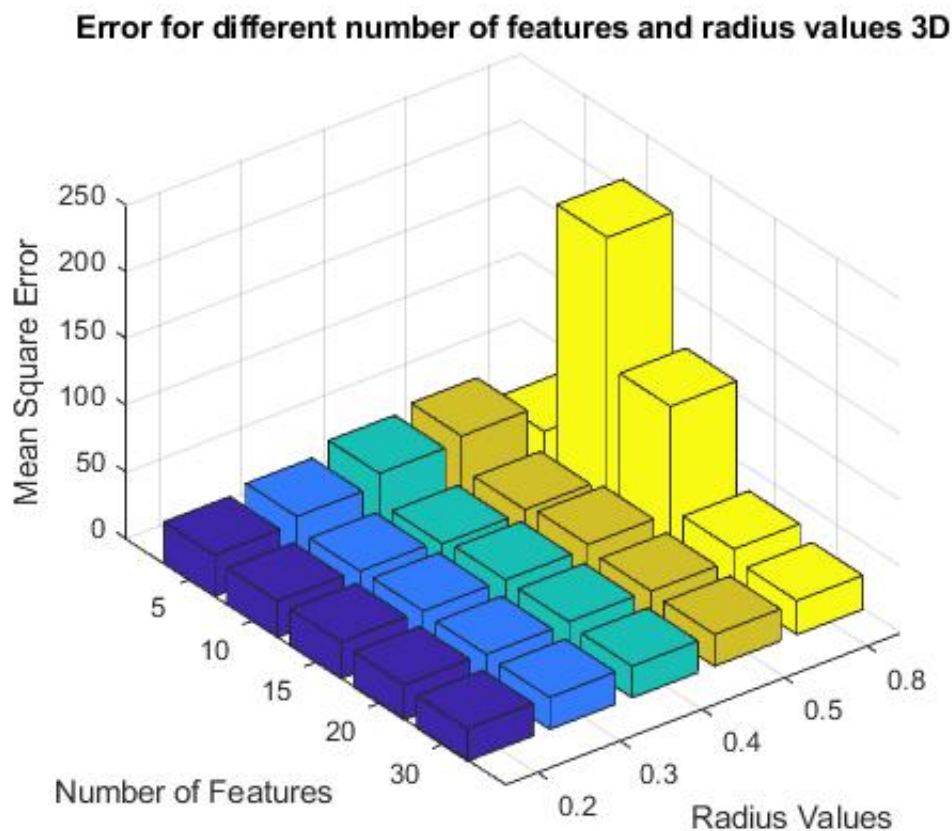
Η παραπάνω διαδικασία εκτελέστηκε για κάθε σημείο του πλέγματος, ελήφθησαν ως βέλτιστες τιμές των παραμέτρων οι τιμές που αντιστοιχούν στο μοντέλο που παρουσίασε το ελάχιστο μέσο σφάλμα. Οι τιμές αυτές χρησιμοποιήθηκαν για την εκπαίδευση του τελικού μοντέλου.

Ως μέθοδος ομαδοποίησης για τη δημιουργία των IF-ELSE κανόνων επιλέχθηκε ο αλγόριθμος Subtractive Clustering και η επιλογή χαρακτηριστικών εκτελέστηκε με τον αλγόριθμο RReliefF.

3. Με βάση τις βέλτιστες τιμές των παραμέτρων που επιλέχθηκαν από το προηγούμενο βήμα, εκπαιδεύτηκε ένα τελικό TSK μοντέλο και ελέγχθηκε η απόδοσή του στο σύνολο ελέγχου.

2.1 Ζητούμενα προβλήματος

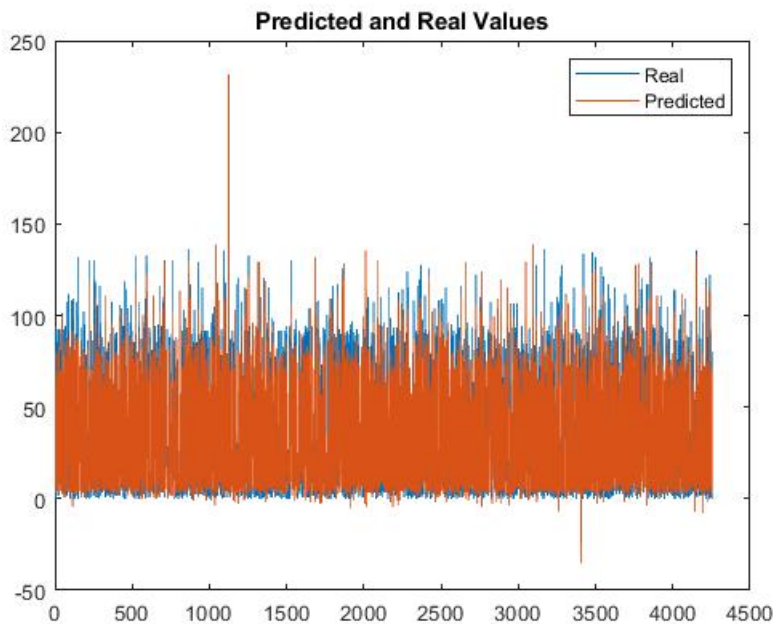
1. Για τα μοντέλα που παρακάτω δίνεται το διάγραμμα το οποίο απεικονίζει την σχέση του μέσου σφάλματος σε με τον αριθμό των επιλεγμένων χαρακτηριστικών και το μέγεθος της ακτίνας.



Σχήμα 9: Μέσο Σφάλμα ως συνάρτηση τον αριθμό των χαρακτηριστικών και το μέγεθος της ακτίνας

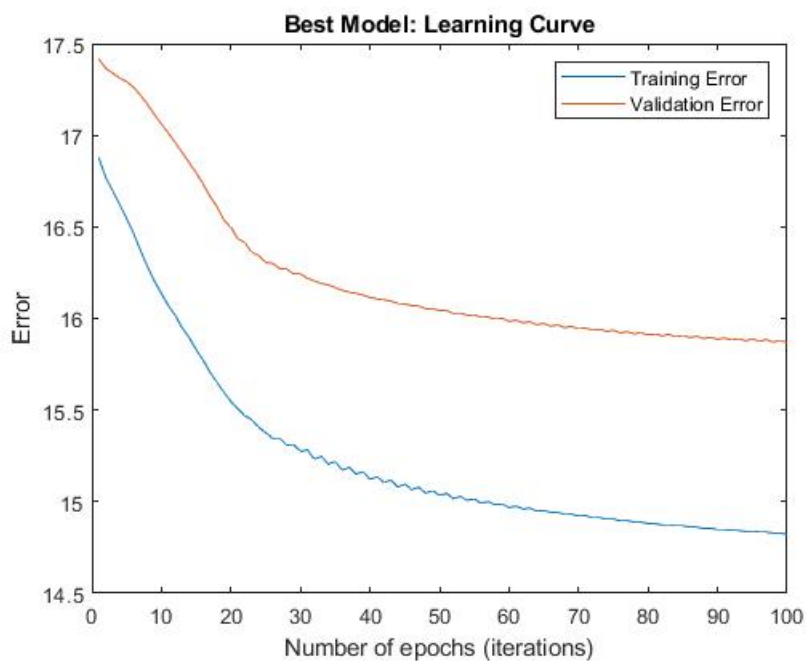
2. Για το τελικό μοντέλο ακόμα παράχθηκαν:

- Διάγραμμα όπου αποτυπώνονται οι προβλέψεις του τελικού μοντέλου καθώς και οι πραγματικές τιμές.



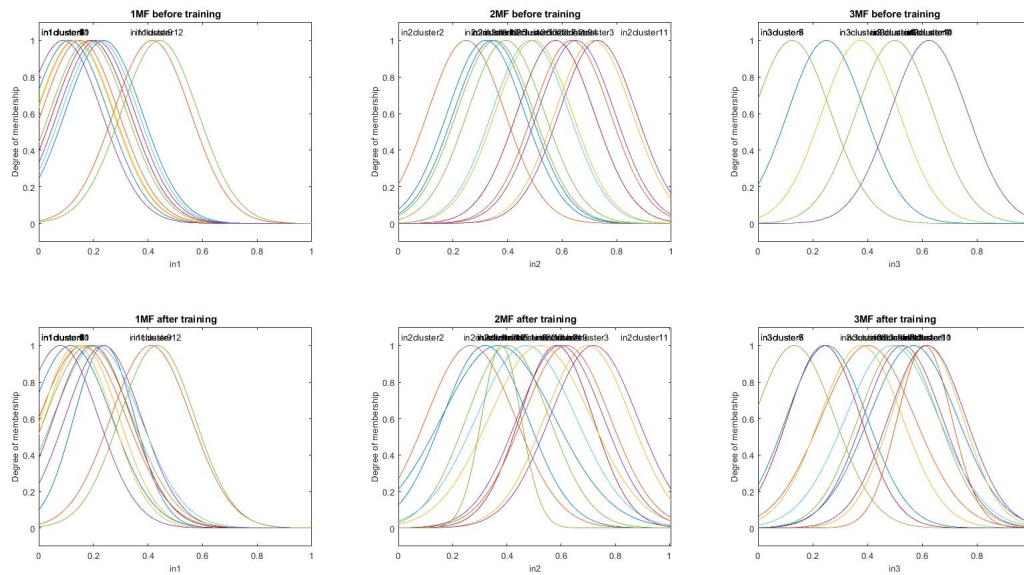
Σχήμα 10: Προβλέψεις και πραγματικές τιμές μοντέλου.

- Διάγραμμα εκμάθησης όπου απεικονίζεται το σφάλμα συναρτήσε του αριθμού επανάληψης.



Σχήμα 11: Καμπύλη εκμάθησης.

- Ενδεικτικά μερικά ασαφή σύνολα στην αρχική και τελική τους μορφή.



Σχήμα 12: Τρία ασαφή σύνολα στην αρχική και τελική τους μορφή (περιέχουν 12 κανόνες).

- Οι τιμές των δεικτών απόδοσης.

Μοντέλο	RMSE	NMSE	NDEI	R^2
Best_Model	15.554861	0.205941	0.453807	0.794059

Πίνακας 3: Δείκτες απόδοσης βέλτιστου μοντέλου

2.2 Σχολιασμός αποτελεσμάτων

Συγκρίνοντας τα αποτελέσματα των μετρικών των δύο πειραμάτων μπορεί να παρατηρηθεί πως τα αποτελέσματα του δεύτερου πειράματος δεν ήταν τα βέλτιστα καθώς έχουμε μεγαλύτερα νούμερα στις μετρικές που βασίζονται στο σφάλμα πρόβλεψης και μικρότερο δείκτη απόδοσης από τον καλύτερο του πρώτου πειράματος.

Αυτό το αποτέλεσμα ήταν αναμενόμενο καθώς έγινε μια προσπάθεια για την δημιουργία ενός μοντέλου το οποίο είχε 82 χαρακτηριστικά. Όμως τα αποτελέσματα δεν είναι εντελώς αποθαρρυντικά καθώς εκτός από το RMSE όλες οι άλλες μετρικές δεν έχουν πολύ μεγάλες διαφορές από το βέλτιστο μοντέλο του πρώτου πειράματος.