

Πανεπιστήμιο Δυτικής Μακεδονίας  
Τμήμα Ηλεκτρολόγων Μηχανικών & Μηχανικών  
Υπολογιστών

---

Αλγόριθμοι μηχανικής μάθησης στον  
εντοπισμό ανωμαλιών για πρόβλεψη  
συντήρησης σε βιομηχανικές εφαρμογές

---

Μπλεόνα Τσεκάνι (ΑΜ: 1086)  
Επιβλέπων Καθηγητής: Νικόλαος Πλόσκας

Εργαστήριο Ευφύων Συστημάτων & Βελτιστοποίησης  
22 Ιουνίου 2023



# Περίληψη

Ο όρος "Τέταρτη Βιομηχανική Επανάσταση" χρησιμοποιείται για να περιγράψει το συνεχιζόμενο μετασχηματισμό των βιομηχανιών με την ενσωμάτωση προηγμένων τεχνολογιών, όπως το Διαδίκτυο των πραγμάτων (IoT), την τεχνητή νοημοσύνη (AI), την υπολογιστική νέφους και την ανάλυση μεγάλων δεδομένων. Ο μετασχηματισμός αυτός οδηγεί στη δημιουργία έξυπνων εργοστασίων, όπου οι μηχανές και τα συστήματα επικοινωνούν και αλληλεπιδρούν μεταξύ τους, με αποτέλεσμα τη βελτίωση της παραγωγικότητας, της αποδοτικότητας και της κερδοφορίας. Το Industry 4.0 αντιπροσωπεύει μια σημαντική αλλαγή στον τρόπο λειτουργίας των επιχειρήσεων, καθώς επιτρέπει την αυτοματοποίηση πολλών διαδικασιών, τη δυνατότητα συλλογής και ανάλυσης τεράστιων ποσοτήτων δεδομένων σε πραγματικό χρόνο και τη δυνατότητα ανάπτυξης νέων προϊόντων και υπηρεσιών.

Η παρούσα διπλωματική εργασία επικεντρώνεται στην εφαρμογή αλγορίθμων επιβλεπόμενης μηχανικής μάθησης για την ανίχνευση ανωμαλιών σε βιομηχανικά συστήματα ηλεκτρικής παραγωγής. Χρησιμοποιήθηκαν δύο σύνολα δεδομένων για δυαδική και πολυταξική μέθοδο κατηγοριοποίησης, με την τελευταία να αντιμετωπίζεται και ως δυαδική για συγκριτικούς σκοπούς. Για τη βελτίωση της απόδοσης των μοντέλων χρησιμοποιήθηκαν μέθοδοι βελτιστοποίησης παραμέτρων. Τα αποτελέσματα των πειραμάτων καταδεικνύουν ότι οι προτεινόμενες μέθοδοι παρουσιάζουν ένα πολλά υποσχόμενο επίπεδο επιδόσεων, υποδεικνύοντας σημαντικές δυνατότητες για μελλοντικές εφαρμογές. Συνοψίζοντας, η παρούσα μελέτη παρέχει πολύτιμες πληροφορίες σχετικά με την αποτελεσματικότητα των αλγορίθμων μηχανικής μάθησης στο πλαίσιο της ανίχνευσης σφαλμάτων, με συνέπειες για τη μελλοντική έρευνα και την πρακτική εφαρμογή.

**Λέξεις κλειδιά:** Τέταρτη Βιομηχανική Επανάσταση, Μηχανική Μάθηση, Εντοπισμός Σφαλμάτων, Εντοπισμός Ανωμαλιών, Κατηγοριοποίηση

# Abstract

Industry 4.0 is a term used to describe the ongoing transformation of industries with the integration of advanced technologies such as the Internet of Things (IoT), artificial intelligence (AI), cloud computing, and big data analytics. This transformation is leading to the creation of smart factories, where machines and systems communicate and interact with each other, resulting in improved productivity, efficiency, and profitability. Industry 4.0 represents a significant shift in the way businesses operate, as it allows for the automation of many processes, the ability to collect and analyze vast amounts of data in real-time, and the potential for the development of new products and services. This concept is already transforming many industries, and its potential impact is far-reaching, with implications for businesses, governments, and society as a whole.

This diploma thesis focuses on the application of supervised machine learning algorithms for fault detection in electrical systems. Two datasets were employed, and binary and multiclass classification methods were tested, with the latter being treated as binary for comparative purposes. To enhance the performance of the models, tuning techniques were utilized. The results of the experiments demonstrate that the proposed methods exhibit a promising performance level, indicating significant potential for future application. In summary, this study provides valuable insights into the efficacy of machine learning algorithms in the context of fault detection, with implications for future research and practical implementation.

**Keywords:** Industry 4.0, Machine Learning, Fault Detection, Anomaly Detection, Classification

# Ευχαριστίες

Θα ήθελα καταρχήν να ευχαριστήσω τον καθηγητή κ. Νικόλαο Πλόσκα για την επίβλεψη αυτής της διπλωματικής εργασίας. Η καθοδήγηση του έπαιξε καθοριστικό ρόλο στη διαμόρφωση της διπλωματικής μου εργασίας και στον εμπλουτισμό της ακαδημαϊκής μου εμπειρίας.

Τέλος, το μεγαλύτερο ευχαριστώ ανήκει στην οικογένεια μου για τη συνεχή υποστήριξη και ενθάρρυνσή τους όλα αυτά τα φοιτητικά χρόνια. Η πίστη τους στις ικανότητές μου και οι θυσίες τους ήταν ανεκτίμητες σε αυτό το ταξίδι. Η αγάπη, η κατανόηση και η υπομονή τους μου προσέφεραν τη δύναμη να ξεπεράσω τις προκλήσεις και να επιμείνω. Είμαι πραγματικά ευγνώμων για την παρουσία τους στη ζωή μου και οφείλω την επιτυχία μου στην αμέριστη υποστήριξή τους.

# Δήλωση Πνευματικών Δικαιωμάτων

Δήλωση Πνευματικών Δικαιωμάτων Δηλώνω ρητά ότι, σύμφωνα με το άρθρο 8 του Ν. 1599/1986 και τα άρθρα 2,4,6 παρ. 3 του Ν. 1256/1982, η παρούσα Διπλωματική Εργασία με τίτλο "Αλγόριθμοι μηχανικής μάθησης στον εντοπισμό ανωμαλιών για πρόβλεψη συντήρησης σε βιομηχανικές εφαρμογές" καθώς και τα ηλεκτρονικά αρχεία και πηγαίοι κώδικες που αναπτύχθηκαν ή τροποποιήθηκαν στα πλαίσια αυτής της εργασίας και αναφέρονται ρητώς μέσα στο κείμενο που συνοδεύουν, και η οποία έχει εκπονηθεί στο Τμήμα Ηλεκτρολόγων Μηχανικών & Μηχανικών Υπολογιστών του Πανεπιστημίου Δυτικής Μακεδονίας, υπό την επίβλεψη του μέλους του Τμήματος κ. Νικόλαου Πλόσκα αποτελεί αποκλειστικά προϊόν προσωπικής εργασίας και δεν προσβάλλει κάθε μορφής πνευματικά δικαιώματα τρίτων και δεν είναι προϊόν μερικής ή ολικής αντιγραφής, οι πηγές δε που χρησιμοποιήθηκαν περιορίζονται στις βιβλιογραφικές αναφορές και μόνον. Τα σημεία όπου έχω χρησιμοποιήσει ιδέες, κείμενο, αρχεία ή / και πηγές άλλων συγγραφέων, αναφέρονται ευδιάκριτα στο κείμενο με την κατάλληλη παραπομπή και η σχετική αναφορά περιλαμβάνεται στο τμήμα των βιβλιογραφικών αναφορών με πλήρη περιγραφή.

Απαγορεύεται η αντιγραφή, αποθήκευση και διανομή της παρούσας εργασίας, εξ ολοκλήρου ή τμήματος αυτής, για εμπορικό σκοπό. Επιτρέπεται η ανατύπωση, αποθήκευση και διανομή για σκοπό μη κερδοσκοπικό, εκπαιδευτικής ή ερευνητικής φύσης, υπό την προϋπόθεση να αναφέρεται η πηγή προέλευσης και να διατηρείται το παρόν μήνυμα. Ερωτήματα που αφορούν τη χρήση της εργασίας για κερδοσκοπικό σκοπό πρέπει να απευθύνονται προς τον συγγραφέα. Οι απόψεις και τα συμπεράσματα που περιέχονται σε αυτό το έγγραφο εκφράζουν τον συγγραφέα και μόνο.

Copyright (C) Μπλεόνα Τσεκάνι & Νικόλαος Πλόσκας, 2023, Κοζάνη

Υπογραφή Φοιτητή



# Περιεχόμενα

<b>1</b>	<b>Εισαγωγή</b>	<b>14</b>
1.1	Ορισμός του προβλήματος . . . . .	14
1.2	Στόχοι της διπλωματικής εργασίας . . . . .	23
1.3	Διάρθρωση κειμένου . . . . .	24
<b>2</b>	<b>Θεωρητικό Υπόβαθρο</b>	<b>25</b>
2.1	Ορισμοί βασικών εννοιών . . . . .	25
2.2	Αλγόριθμοι επιβλεπόμενης μάθησης . . . . .	29
2.3	Μετρικές Αξιολόγησης Μοντέλων και Αποτελεσμάτων . . . . .	39
2.4	Μετρικές χρόνου εκτέλεσης . . . . .	42
<b>3</b>	<b>Βιβλιογραφική Ανασκόπηση</b>	<b>43</b>
3.1	Παραδείγματα εφαρμογής αλγορίθμων μηχανικής μάθησης . . . . .	43
3.1.1	Πρόβλεψη περιόδου βλάβης με χρήση τεχνητών νευρωνικών δικτύων . . . . .	43
3.1.2	Μια έξυπνη προσέγγιση για προεπεξεργασία και ανάλυση δεδομένων σε μια περίπτωση βιομηχανικής μελέτης . . . . .	46
3.1.3	Αρχιτεκτονική προβλεπτικής συντήρησης με μηχανική μάθηση σε πυρηνικές υποδομές . . . . .	48
3.1.4	SOPHIA: Μια αρχιτεκτονική IoT βασισμένη σε δεδομένα και μηχανική μάθηση για την προβλεπτική συντήρηση στη βιομηχανία 4.0 . . . . .	49
3.1.5	Προσέγγιση μηχανικής μάθησης για προβλεπτική συντήρηση στη βιομηχανία 4.0 . . . . .	50
3.1.6	Μηχανική μάθηση για προβλεπτική συντήρηση βιομηχανικών συστημάτων με χρήση δεδομένων αισθητήρων IoT. . . . .	51

3.1.7	Πρόβλεψη RUL στροβιλοκινητήρων μέσω της μηχανικής μάθησης	52
3.1.8	Ανάλυση βασισμένη σε αισθητήρες IoT, προεπεξεργασία Big Data και μοντέλων μηχανικής μάθησης για συστήματα παρακολούθησης πραγματικού χρόνου σε αυτοκινητοβιομηχανίες . . .	53
3.2	Ερευνητικά και αναπτυξιακά έργα στην Ευρωπαϊκή Ένωση . . . . .	54
3.3	Εμπορικές δραστηριότητες . . . . .	59
<b>4</b>	<b>Εργαλεία Υλοποίησης</b>	<b>62</b>
4.1	Εισαγωγή . . . . .	62
4.2	Jupyter Notebook . . . . .	62
4.3	Python . . . . .	63
4.4	Pandas . . . . .	63
4.5	Scikit-Learn . . . . .	63
4.6	Matplotlib . . . . .	64
<b>5</b>	<b>Ανάλυση Εφαρμογών και Προεπεξεργασία Δεδομένων</b>	<b>65</b>
5.1	Συστήματα ηλεκτρικής ενέργειας . . . . .	66
5.2	Ανάλυση δεδομένων . . . . .	67
5.3	Πλάνο μελέτης . . . . .	71
5.4	Τεχνικές αξιολόγησης . . . . .	72
<b>6</b>	<b>Πειραματικές Εφαρμογές και Αποτελέσματα</b>	<b>74</b>
6.1	Σχεδιασμός μοντέλου δυαδικής κατηγοριοποίησης . . . . .	75
6.1.1	Ανάλυση . . . . .	75
6.1.2	Σύνοψη . . . . .	78
6.2	Σχεδιασμός μοντέλου πολυταξικής κατηγοριοποίησης . . . . .	80
6.2.1	Ανάλυση . . . . .	80
6.2.2	Σύνοψη . . . . .	83
6.3	Χειρισμός πολυταξικής κατηγοριοποίησης ως δυαδική . . . . .	84
6.3.1	Ανάλυση . . . . .	84
6.3.2	Σύνοψη . . . . .	98
<b>7</b>	<b>Βελτιστοποίηση Παραμέτρων</b>	<b>102</b>
7.1	Παραμετροποίηση αλγορίθμων δυαδικής κατηγοριοποίησης . . . . .	103



---

7.1.1	Grid search CV . . . . .	103
7.1.2	Randomized search CV . . . . .	105
7.1.3	Σύνοψη . . . . .	107
7.2	Παραμετροποίηση αλγορίθμων πολυταξικής κατηγοριοποίησης . . . .	108
7.2.1	Grid search SV . . . . .	108
7.2.2	Randomized search CV . . . . .	110
7.2.3	Σύνοψη . . . . .	111
<b>8</b>	<b>Συμπεράσματα</b>	<b>113</b>

# Κατάλογος σχημάτων

1.1	Οι τέσσερις βιομηχανικές επαναστάσεις [37]	15
1.2	Στρατηγικές συντήρησης [23]	16
1.3	Ενισχυτικές τεχνολογίες και στοιχεία του Industry 4.0 [45]	19
2.1	Εξέλιξη Τεχνητής Νοημοσύνης, Μηχανικής Μάθησης και Βαθιάς Μάθησης [34]	26
2.2	Κατηγορίες Μηχανικής Μάθησης	28
2.3	Logistic Regression [40]	30
2.4	K-Nearest Neighbors[7]	31
2.5	Support Vector Machines [42]	32
2.6	Random Forest [47]	33
2.7	Gaussian Naive Bayes[41]	35
2.8	eXtreme Gradient Boosting[44]	36
2.9	Decision Tree [32]	37
2.10	Multi-layer Perceptron[15]	38
3.1	Έλεγχος αρχικού συνόλου [46]	45
3.2	Έλεγχος γενίκευσης ( $\alpha$ ) [46]	45
3.3	Έλεγχος γενίκευσης ( $\beta$ ) [46]	46
3.4	Αποτελέσματα K-means ομαδοποίησης για την πρώτη μηχανή [19]	47
3.5	Αποτελέσματα K-means ομαδοποίησης για τη δεύτερη μηχανή [19]	48
3.6	Οπτικοποίηση του Root Mean Square Error των RULs έναντι Machine Ids για διαφορετικούς αλγόριθμους σε τέσσερα διαφορετικά σύνολα δεδομένων [36]	53
5.1	Συστήματα ηλεκτρικής ενέργειας [4]	66

---

5.2	Ρεύμα και τάση στη γραμμή B σε φυσιολογικές και μη φυσιολογικές συνθήκες . . . . .	68
5.3	Ρεύμα και τάση στις γραμμές σε φυσιολογικές και μη φυσιολογικές συνθήκες . . . . .	69
5.4	Ρεύμα και τάση στις γραμμές A, B και C σε διάφορες συνθήκες . . . .	71
6.1	Διαδικασία μελέτης [13] . . . . .	74
6.2	Ακρίβεια και σφάλμα στη δυαδική κατηγοριοποίηση . . . . .	75
6.3	Μετρικές κατάστασης cross validation δυαδικής κατηγοριοποίησης . .	76
6.4	Σύγκριση των μοντέλων μέσω του χρόνου fit και score δυαδικής κατηγοριοποίησης . . . . .	77
6.5	Ακρίβεια και σφάλμα στην πολυταξική κατηγοριοποίηση . . . . .	80
6.6	Μετρικές κατάστασης cross validation στην πολυταξική κατηγοριοποίηση . . . . .	81
6.7	Confusion matrices . . . . .	82
6.8	Σύγκριση μοντέλων πολυταξικής κατηγοριοποίησης μέσω του χρόνου fit και score . . . . .	83
6.9	Ακρίβεια - σφάλμα σε 0000 . . . . .	85
6.10	Σύγκριση μοντέλων μέσω μετρικών κατηγοριοποίησης 0000 . . . . .	86
6.11	Χρόνοι εκπαίδευσης και αξιολόγησης σε 0000 . . . . .	86
6.12	Ακρίβεια - σφάλμα σε σφάλμα 1011 . . . . .	87
6.13	Σύγκριση μοντέλων μέσω μετρικών κατηγοριοποίησης σε σφάλμα 1011	88
6.14	Χρόνοι εκπαίδευσης και αξιολόγησης σε σφάλμα 1011 . . . . .	88
6.15	Ακρίβεια - σφάλμα σε 1111 . . . . .	89
6.16	Σύγκριση μοντέλων μέσω μετρικών κατηγοριοποίησης σε σφάλμα 1111	90
6.17	Χρόνοι εκπαίδευσης και αξιολόγησης σε σφάλμα 1111 . . . . .	90
6.18	Ακρίβεια - σφάλμα σε σφάλμα 1001 . . . . .	91
6.19	Σύγκριση μοντέλων μέσω μετρικών κατηγοριοποίησης σε σφάλμα 1001	92
6.20	Χρόνοι εκπαίδευσης αξιολόγησης σε σφάλμα 1001 . . . . .	92
6.21	Ακρίβεια - σφάλμα σε σφάλμα 0111 . . . . .	93
6.22	Σύγκριση μοντέλων μέσω μετρικών κατηγοριοποίησης σε σφάλμα 0111	94
6.23	Χρόνοι εκπαίδευσης και αξιολόγησης σε σφάλμα 0111 . . . . .	94
6.24	Ακρίβεια - σφάλμα σε σφάλμα 0110 . . . . .	95

---

6.25 Σύγκριση μοντέλων μέσω μετρικών κατηγοριοποίησης σε σφάλμα 0110	96
6.26 Χρόνοι εκπαίδευσης και αξιολόγησης σε σφάλμα 0110 . . . . .	96

# Κατάλογος πινάκων

3.1	RMSE τιμές των υποσυνόλων που χρησιμοποιήθηκαν στον έλεγχο των τεχνικών μηχανικής μάθησης . . . . .	44
3.2	Τιμές RMSE για τον έλεγχο γενίκευσης των τεχνικών μηχανικής μάθησης	46
3.3	Αριθμός περιπτώσεων της κάθε συστάδας για το κάθε μηχάνημα . . . . .	48
3.4	Σύγκριση αποτελεσμάτων . . . . .	49
3.5	Μετρικές ορθότητας, μνήμης και ακρίβεια των GBM, DRF και XGBoost	50
3.6	Αποτελέσματα ομαδοποίησης . . . . .	51
3.7	Σύγκριση διαφορετικών επιβλεπόμενων αλγορίθμων . . . . .	52
3.8	Σύγκριση αποδόσεων διάφορων συγκριτικών μοντέλων . . . . .	54
6.1	Πίνακας τιμών ακρίβειας - σφάλματος στη δυαδική κατηγοριοποίηση	76
6.2	Πίνακας μετρικών κατάστασης cross validation στη δυαδική κατηγοριοποίηση . . . . .	76
6.3	Πίνακας χρονικών τιμών fit - score δυαδικής κατηγοριοποίησης . . . . .	77
6.4	Πίνακας τιμών ακρίβειας - σφάλματος στην πολυταξική κατηγοριοποίηση . . . . .	80
6.5	Πίνακας μετρικών κατάστασης cross validation στην πολυταξική κατηγοριοποίηση . . . . .	81
6.6	Πίνακας χρονικών τιμών fit - score . . . . .	83
6.7	Πίνακας τιμών ακρίβειας - σφάλματος σε 0000 . . . . .	85
6.8	Πίνακας μετρικών κατάστασης cross validation σε 0000 . . . . .	86
6.9	Πίνακας χρονικών τιμών εκπαίδευσης - αξιολόγησης σε 0000 . . . . .	87
6.10	Πίνακας τιμών ακρίβειας - σφάλματος σε σφάλμα 1011 . . . . .	87
6.11	Πίνακας μετρικών κατάστασης cross validation σε σφάλμα 1011 . . . . .	88
6.12	Πίνακας χρονικών τιμών εκπαίδευσης - αξιολόγησης σε 1011 . . . . .	89
6.13	Πίνακας τιμών ακρίβειας - σφάλματος σε σφάλμα 1111 . . . . .	89

6.14	Πίνακας μετρικών κατάστασης cross validation σε σφάλμα 1111 . . . . .	90
6.15	Πίνακας χρονικών τιμών εκπαίδευσης - αξιολόγησης . . . . .	91
6.16	Πίνακας τιμών ακρίβειας - σφάλματος σε σφάλμα 1001 . . . . .	91
6.17	Πίνακας μετρικών κατάστασης cross validation σε σφάλμα 1001 . . . . .	92
6.18	Πίνακας χρονικών τιμών εκπαίδευσης - αξιολόγησης . . . . .	93
6.19	Πίνακας τιμών ακρίβειας - σφάλματος σε σφάλμα 0111 . . . . .	93
6.20	Πίνακας μετρικών κατάστασης cross validation σε σφάλμα 0111 . . . . .	94
6.21	Πίνακας χρονικών τιμών εκπαίδευσης - αξιολόγησης . . . . .	95
6.22	Πίνακας τιμών ακρίβειας - σφάλματος σε σφάλμα 0111 . . . . .	95
6.23	Πίνακας μετρικών κατάστασης cross validation σε σφάλμα 0110 . . . . .	96
6.24	Πίνακας Χρονικών Τιμών Εκπαίδευσης - Αξιολόγησης . . . . .	97
6.25	Συγκριτικός πίνακας τιμών ακρίβειας σε πολυταξική και δυαδική κα- τηγοριοποίηση . . . . .	98
6.26	Συγκριτικός πίνακας τιμών σφάλματος σε πολυταξική και δυαδική κατηγοριοποίηση . . . . .	98
6.27	Πίνακας μετρικών κατάστασης cross validation σε δυαδική και πολυ- ταξική κατηγοριοποίηση . . . . .	99
6.28	Πίνακας χρονικών τιμών fit - score σε πολυταξική και δυαδική κατη- γοριοποίηση . . . . .	100
7.1	Αποτελέσματα μετρικών στη δυαδική κατηγοριοποίηση . . . . .	103
7.2	Τιμές παραμέτρων D-Tree στη δυαδική κατηγοριοποίηση . . . . .	103
7.3	Αποτελέσματα μετρικών D-Tree στη δυαδική κατηγοριοποίηση . . . . .	104
7.4	Τιμές παραμέτρων Random Forest στη δυαδική κατηγοριοποίηση . . . . .	104
7.5	Αποτελέσματα μετρικών Random Forest στη δυαδική κατηγοριοποίηση	104
7.6	Τιμές παραμέτρων XGBoost στη δυαδική κατηγοριοποίηση . . . . .	105
7.7	Αποτελέσματα μετρικών XGBoost στη δυαδική κατηγοριοποίηση . . . . .	105
7.8	Τιμές παραμέτρων D-Tree στη δυαδική κατηγοριοποίηση . . . . .	105
7.9	Αποτελέσματα μετρικών D-Tree στη δυαδική κατηγοριοποίηση . . . . .	106
7.10	Τιμές παραμέτρων Random Forest στη δυαδική κατηγοριοποίηση . . . . .	106
7.11	Αποτελέσματα μετρικών Random Forest στη δυαδική κατηγοριοποίηση	106
7.12	Τιμές παραμέτρων XGBoost στη δυαδική κατηγοριοποίηση . . . . .	107
7.13	Αποτελέσματα μετρικών XGBoost στη δυαδική κατηγοριοποίηση . . . . .	107

---

7.14 Πίνακας μετρικών κατάστασης Cross Validation στην πολυταξική κατηγοριοποίηση . . . . .	108
7.15 Τιμές παραμέτρων D-Tree στην πολυταξική κατηγοριοποίηση . . . . .	108
7.16 Αποτελέσματα μετρικών D-Tree στην πολυταξική κατηγοριοποίηση . .	109
7.17 Τιμές παραμέτρων Random Forest στην πολυταξική κατηγοριοποίηση	109
7.18 Αποτελέσματα μετρικών Random Forest στην πολυταξική κατηγοριοποίηση . . . . .	109
7.19 Τιμές παραμέτρων XGBoost στην πολυταξική κατηγοριοποίηση . . . .	109
7.20 Αποτελέσματα μετρικών XGBoost στην πολυταξική κατηγοριοποίηση .	110
7.21 Τιμές παραμέτρων D-Tree στην πολυταξική κατηγοριοποίηση . . . . .	110
7.22 Αποτελέσματα μετρικών D-Tree στην πολυταξική κατηγοριοποίηση . .	110
7.23 Τιμές παραμέτρων Random Forest στην πολυταξική κατηγοριοποίηση	110
7.24 Αποτελέσματα μετρικών Random Forest στην πολυταξική κατηγοριοποίηση . . . . .	111
7.25 Τιμές παραμέτρων XGBoost στην πολυταξική κατηγοριοποίηση . . . .	111
7.26 Αποτελέσματα μετρικών XGoost στην πολυταξική κατηγοριοποίηση .	111





# Κεφάλαιο 1

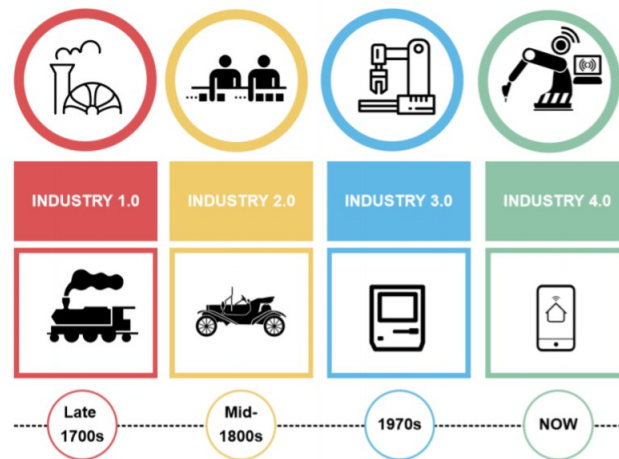
## Εισαγωγή

### 1.1 Ορισμός του προβλήματος

Επί του παρόντος, η βιομηχανία περνά από αυτό που οι ειδικοί ονομάζουν "Τέταρτη Βιομηχανική Επανάσταση" ή αλλιώς Industry 4.0. Σύμφωνα με το Σχήμα 1.1 η πρώτη βιομηχανική επανάσταση ξεκινάει από το 1750 έως το 1830, με την εφεύρεση της ατμομηχανής και των σιδηροδρόμων. Η δεύτερη βιομηχανική επανάσταση διήρκεσε από το 1870 έως το 1900, εκείνη την περίοδο κατέστη δυνατή η μαζική παραγωγή μέσω της χρήσης ηλεκτρικής ενέργειας, κινητήρων εσωτερικής καύσης, χημικών ουσιών και πετρελαίου. Η τρίτη βιομηχανική επανάσταση ξεκίνησε τη δεκαετία του 1960 και επικεντρώθηκε στη διάδοση της μικροηλεκτρονικής, βασιζόμενη στην ανακάλυψη του τρανζίστορ (1947), των ολοκληρωμένων κυκλωμάτων (1957), της επίπεδης διαδικασίας (1959), των ημιαγωγών (1960), του μικροεπεξεργαστή (1971), των προσωπικών υπολογιστών (1970-1980), του Διαδικτύου (1990s) και των smartphones (2000s). Η τέταρτη βιομηχανική επανάσταση σχετίζεται άμεσα με την ενοποίηση των φυσικών και των ψηφιακών συστημάτων στο περιβάλλοντων παραγωγής. Η ενοποίηση επιτρέπει τη συλλογή μεγάλου όγκου δεδομένων από διαφορετικά συστήματα, τοποθετημένα σε διαφορετικά μέρη του εργοστασίου. Επιπλέον, οι τεχνολογίες που σχετίζονται με το Industry 4.0 ενσωματώνουν ανθρώπινο δυναμικό, μηχανήματα και προϊόντα επιτρέποντας ταχύτερη και πιο στοχευόμενη ανταλλαγή πληροφοριών.

Ο μεγάλος όγκος δεδομένων που συλλέγεται από τα βιομηχανικά συστήματα περιέχει πληροφορίες σχετικά με τις διεργασίες, τα γεγονότα και τους συναγερμούς που προκύπτουν σε μια βιομηχανική γραμμή παραγωγής. Επιπλέον, όταν επεξεργα-

Σχήμα 1.1: Οι τέσσερις βιομηχανικές επαναστάσεις [37]



στούν και αναλυθούν, είναι πιθανόν να εξαχθούν από τα δεδομένα αυτά πολύτιμες πληροφορίες για τη διαδικασία της κατασκευής και της δυναμικής του συστήματος. Εφαρμόζοντας αναλυτικές προσεγγίσεις με βάση τα δεδομένα, είναι πιθανόν να βρεθούν ερμηνευτικά αποτελέσματα σχετικά με τη στρατηγική λήψη αποφάσεων, να παρέχονται πλεονεκτήματα σχετικά με τη μείωση του κόστους συντήρησης, τη μείωση των σφαλμάτων στις μηχανές, τη μείωση των διακοπών συντήρησης, την αύξηση της διάρκειας ζωής των ανταλλακτικών, την αύξηση της παραγωγής, τη βελτίωση της ασφάλειας των χειριστών και το συνολικό κέρδος μεταξύ άλλων.

Τα πλεονεκτήματα που αναφέρθηκαν προηγουμένως συνδέονται στενά με τις διαδικασίες συντήρησης. Στις βιομηχανίες, η συντήρηση των μηχανημάτων είναι ένα βασικό ζήτημα καθώς επηρεάζει τη χρονική διάρκεια λειτουργίας τους και την αποτελεσματικότητά τους. Για αυτόν το λόγο για να αποφευχθεί ο τερματισμός των διαδικασιών παραγωγής τα σφάλματα θα πρέπει να εντοπίζονται και να λύνονται.

Υπάρχουν διάφορες κατηγορίες στρατηγικών διαχείρισης συντήρησης, στη συνέχεια παρουσιάζονται οι πιο βασικές [8] οι οποίες απεικονίζονται και στο Σχήμα 1.2:

- Προβλεπτική Συντήρηση (Predictive Maintenance - PdM)

Χρησιμοποιεί προγνωστικά εργαλεία για να καθορίσει το πότε θα είναι απαραίτητες οι ενέργειες συντήρησης. Βασίζεται στη συνεχή παρακολούθηση ενός μηχανήματος ή της ακεραιότητας μιας διεργασίας με σκοπό την εφαρμογή της συντήρησης όταν χρειαστεί. Επιπλέον, δίνει τη δυνατότητα της έγκαιρης ανίχνευσης των βλαβών χάρη στα εργαλεία πρόβλεψης που βασίζονται στη χρήση

---

προηγούμενων δεδομένων, σε παράγοντες ακεραιότητας, σε στατιστικές μεθόδους και σε μηχανικές προσεγγίσεις.

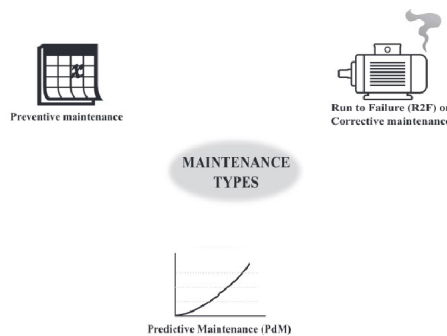
- Προληπτική Συντήρηση (Preventive Maintenance - PnM)

Ονομάζεται και αλλιώς χρονική συντήρηση ή προγραμματισμένη. Είναι μια τεχνική συντήρησης που εκτελείται περιοδικά με βάση κάποιο προγραμματισμένο χρονοδιάγραμμα ή με διαδοχικές προσεγγίσεις διεργασιών για την πρόβλεψη αποτυχιών σε κάποια διεργασία ή μηχανήμα. Γενικά, είναι μια ικανοποιητική προσέγγιση για να αποφευχθούν βλάβες. Ωστόσο, πραγματοποιούνται περιττές διορθωτικές ενέργειες οι οποίες οδηγούν στην αύξηση του λειτουργικού κόστους.

- Λειτουργία μέχρι τη βλάβη (Run-to-Failure - R2F) ή αλλιώς Αντιδραστική Συντήρηση (Reactive Maintenance)

Συμβαίνει μόνο όταν διακοπεί η εργασία κάποιας μηχανής. Αποτελεί την πιο απλή στρατηγική συντήρησης και η διακοπή της παραγωγής σε συνδυασμό με την επισκευή των εξαρτημάτων που πρέπει να αντικατασταθούν προσθέτουν ένα κόστος στη διαδικασία.

Σχήμα 1.2: Στρατηγικές συντήρησης [23]



Γενικά, η στρατηγική R2F επιβραδύνει τις ενέργειες συντήρησης των βιομηχανιών και από την άλλη πλευρά η στρατηγική PnM προβλέπει τις παρεμβάσεις συντήρησης με αποτέλεσμα την πρόωρη αντικατάσταση των ανταλλακτικών. Μια καλή στρατηγική συντήρησης θα πρέπει να βελτιώνει την κατάσταση των μηχανημάτων, να μειώνει τα ποσοστά βλάβης τους και να ελαχιστοποιεί το κόστος συντήρησης ενώ ταυτόχρονα να μεγιστοποιεί τη διάρκεια ζωής τους. Σύμφωνα με τα παραπάνω, η

---

στρατηγική PdM είναι αυτή που ξεχωρίζει περισσότερο μεταξύ των άλλων στρατηγικών και χρησιμοποιείται στην εποχή της Τέταρτης Βιομηχανικής Επανάστασης λόγω της ικανότητας της να βελτιώνει την χρήση και τη διαχείριση των μηχανημάτων.

Όπως περιγράφηκε πιο πάνω, η PdM ασχολείται με την πρόβλεψη σφαλμάτων ή αστοχιών πριν ακόμη συμβούν. Οι προσεγγιστικές τεχνικές συντήρησης που είναι σε θέση να παρακολουθούν τις καταστάσεις των μηχανημάτων για διαγνωστικούς και προγνωστικούς σκοπούς ομαδοποιούνται σε τρεις βασικές κατηγορίες: στατιστικές, τεχνητής νοημοσύνης και βασισμένες σε μοντέλα. Οι προσεγγίσεις που βασίζονται σε μοντέλα χρειάζονται μηχανιστικές και θεωρητικές γνώσεις για το μηχανήμα που πρόκειται να παρακολουθηθεί και η στατιστική προσέγγιση προαπαιτεί ένα μαθηματικό υπόβαθρο με αποτέλεσμα η προσέγγιση τεχνητής νοημοσύνης να εφαρμόζεται όλο και περισσότερο σε εφαρμογές PdM.

Η μηχανική μάθηση αποτελεί κλάδο της τεχνητής νοημοσύνης και έχει αναδειχθεί ως ένα ισχυρό εργαλείο για την ανάπτυξη έξυπνων προβλεπτικών αλγορίθμων. Μέσω της μηχανικής μάθησης δύναται η δυνατότητα χειρισμού υψηλών διαστάσεων και πολυμεταβλητών δεδομένων καθώς και η εξαγωγή κρυφών σχέσεων μεταξύ των δεδομένων σε ένα πολύπλοκο και δυναμικό περιβάλλον. Επομένως, η μηχανική μάθηση προσφέρει μια ισχυρή προβλεπτική προσέγγιση σε εφαρμογές PdM.

Το πανεπιστήμιο Rmit Australia έκανε μια αναφορά [31] και εντόπισε έξι βασικές αρχές καίριας σημασίας που θα πρέπει να ασπαστούν από τις βιομηχανίες έτσι ώστε αυτές να επωφεληθούν από την τεχνολογία του Industry 4.0:

#### 1. Διαλειτουργικότητα

Αναφέρεται στην ικανότητα αντικειμένων, μηχανών και ατόμων σε μια επιχείρηση να επικοινωνούν, να ανταλλάσσουν δεδομένα και να συντονίζουν δραστηριότητες. Η ικανότητα αυτή συμβάλει σημαντικά στην αποτελεσματικότητα και στη βελτιστοποίηση των διαδικασιών μέσω των δεδομένων και των πληροφοριών που παρέχονται. Σημαντικό βήμα για τις επιχειρήσεις είναι η ψηφιοποίηση των λειτουργιών τους με χρήση cloud computing ή ενοποίηση πλατφορμών και λογισμικού ανοιχτού κώδικα, όπως Linux, Android, Apache, OpenOffice, GnuCash, ADempiere, SugarCRM, Drupal, Wordpress και OpenCart.

---

## 2. Εικονοποίηση

Η αρχή αυτή αναφέρεται σε δύο διαφορετικά σενάρια. Στο πρώτο σενάριο, ένας εικονικός πόρος δημιουργείται από πολλούς φυσικούς πόρους και το περιβάλλον μιας βιομηχανίας προσομοιώνεται δημιουργώντας "ψηφιακά δίδυμα" φυσικών στοιχείων με χρήση δεδομένων από αισθητήρες. Στο δεύτερο σενάριο, πολλοί εικονικοί πόροι δημιουργούνται από έναν ή περισσότερους φυσικούς πόρους, το λογισμικό χρησιμοποιείται για να διαιρέσει έναν φυσικό διακομιστή σε πολλούς εικονικούς διακομιστές που λειτουργούν σαν μοναδικές φυσικές συσκευές. Μέσω της εικονικοποίησης οι εφαρμογές, οι υπολογιστές, οι διακομιστές και τα δεδομένα δεν εξαρτώνται πλέον από μια φυσική συσκευή και έτσι βελτιώνεται η αξιοπιστία και ενεργοποίηση προσθηκών όταν απαιτείται.

## 3. Αποκέντρωση

Η κεντρική ιδέα της έννοιας αυτής είναι η μετακίνηση συστημάτων σε εξαρτήματα και όχι σε έναν κεντρικό υπολογιστή για απεριόριστη επεκτασιμότητα και ευελιξία. Με το Industry 4.0, όλη η τεχνολογία είναι αποκεντρωμένη, διευκολύνοντας τη δημιουργία αποκεντρωμένων συστημάτων σε όλες τις βιομηχανίες σε παγκόσμια κλίμακα.

## 4. Δυνατότητα Πραγματικού Χρόνου

Συλλογή και ανάλυση δεδομένων σε πραγματικό χρόνο, επιτρέποντας την άμεση λήψη αποφάσεων κάθε στιγμή. Χάρη στον πολλαπλασιασμό των αισθητήρων και των συνδεδεμένων συσκευών στο διαδίκτυο, οι περισσότερες επιχειρήσεις έχουν ήδη πρόσβαση σε μεγάλες ποσότητες δεδομένων πραγματικού χρόνου. Η πρόκληση έγκειται στην επεξεργασία και ανάλυση τους, έτσι ώστε οι δραστηριότητες να μπορούν να βελτιώνονται συνεχώς.

## 5. Προσανατολισμένη Λειτουργία

Η δυνατότητα πραγματικού χρόνου που παρέχεται χάρη στα μεγάλα δεδομένα και στην ελεύθερη ροή πληροφοριών επιτρέπει στις επιχειρήσεις να ικανοποιούν καλύτερα τις ανάγκες των πελατών. Οι επιχειρήσεις προσαρμόζονται στις μεταβαλλόμενες ανάγκες και προσδοκίες των πελατών παρέχοντας εξατομικευμένες υπηρεσίες. Ως αποτέλεσμα, παρατηρείται μια μετατόπιση της

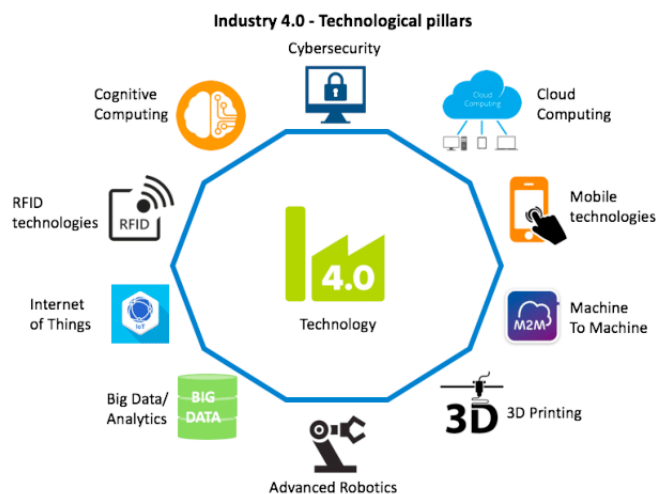
προσοχής προς τους πελάτες και όχι προς τα προϊόντα, καθώς και σε υπηρεσίες που προσαρμόζονται στις ανάγκες και όχι στη μαζική παραγωγή.

## 6. Αρθρωτή οργάνωση

Αναφέρεται στην ικανότητα των επιχειρήσεων να προσαρμόζονται με ευελιξία στις μεταβαλλόμενες απαιτήσεις και στις ανάγκες της βιομηχανίας. Δίνεται στις επιχειρήσεις η δυνατότητα να εστιάζουν σε αυτό που κάνουν καλά. Το προσωπικό γίνεται εξειδικευμένο στις βασικές δραστηριότητες, οι οποίες επεκτείνονται καθώς αναπτύσσεται η επιχείρηση, αυξάνοντας την ευελιξία και την προσαρμοστικότητα της επιχείρησης, καθώς και την ικανότητα τους να ανταποκρίνονται πιο γρήγορα στις μεταβαλλόμενες συνθήκες της αγοράς.

Το Industry 4.0 αναφέρεται στην ολοκλήρωση πολλαπλών τεχνολογιών και παραγόντων με κοινό στόχο τη βελτίωση της αποτελεσματικότητας και της απόκρισης ενός συστήματος παραγωγής. Στο Σχήμα 1.3 παρουσιάζονται οι κύριες τεχνολογίες:

Σχήμα 1.3: Ενισχυτικές τεχνολογίες και στοιχεία του Industry 4.0 [45]



Παρακάτω, δίνεται ένας σύντομος ορισμός και γίνεται ανάλυση των τεχνολογιών του Industry 4.0:

- Τεχνητή Νοημοσύνη (Artificial Intelligence - AI)

Αναφέρεται στην ανάπτυξη υπολογιστικών συστημάτων τα οποία είναι ικανά να εκτελούν εργασίες που συνήθως απαιτούν ανθρώπινη νοημοσύνη. Αποσκοπεί στην προσομοίωση της ευφυούς συμπεριφοράς, συμπεριλαμβανομένων εργασιών όπως επίλυσης προβλημάτων, λήψης αποφάσεων, αναγνώρισης ομιλίας,

---

μετάφρασης και αναγνώρισης εικόνων. Περιλαμβάνει διάφορες τεχνικές και προσεγγίσεις που επιτρέπουν στις μηχανές να αντιλαμβάνονται, να συλλογίζονται, να μαθαίνουν και να αλληλεπιδρούν με το περιβάλλον τους.

- Μηχανική Μάθηση (Machine Learning)

Η μηχανική μάθηση είναι ένα υποσύνολο της τεχνητής νοημοσύνης που επικεντρώνεται στην ανάπτυξη αλγορίθμων και μοντέλων που επιτρέπουν στους υπολογιστές να μαθαίνουν από δεδομένα και να κάνουν προβλέψεις ή αποφάσεις χωρίς να είναι ρητά προγραμματισμένοι. Οι αλγόριθμοι Μηχανικής Μάθησης μαθαίνουν μοτίβα και σχέσεις στα δεδομένα μέσω της εκπαίδευσης και βελτιώνουν επαναληπτικά την απόδοσή τους με την πάροδο του χρόνου. Χρησιμοποιείται ευρέως σε διάφορες εφαρμογές, συμπεριλαμβανομένης της ανάλυσης δεδομένων, της αναγνώρισης προτύπων και των αυτόνομων οχημάτων.

- Κυβερνο-φυσικά συστήματα (Cyber-Physical Systems - CPS)

Αντιπροσωπεύουν μια νέα γενιά ψηφιακών συστημάτων, τα οποία αποτελούνται από δύο κύρια λειτουργικά στοιχεία:

1. προηγμένη συνδεσιμότητα που εξασφαλίζει τη συλλογή δεδομένων σε πραγματικό χρόνο από το φυσικό κόσμο και την ανατροφοδότηση πληροφοριών από τον κυβερνοχώρο.
2. ευφυής διαχείριση δεδομένων, ανάλυση και υπολογιστική ικανότητα που δομούν τον κυβερνοχώρο.

Η χρήση του CPS στοχεύει στην υλοποίηση συστημάτων μεγάλης κλίμακας βελτιώνοντας τη δυνατότητα προσαρμογής, την αυτονομία, την αποδοτικότητα, τη λειτουργικότητα, την αξιοπιστία, την ασφάλεια και τη χρηστικότητα τους.

- Διαδίκτυο των Πραγμάτων (Internet of Things - IoT)

Περιγράφει το δίκτυο των φυσικών αντικειμένων που είναι ενσωματωμένα με αισθητήρες, λογισμικό και άλλες τεχνολογίες, με σκοπό τη σύνδεση και την ανταλλαγή δεδομένων με άλλες συσκευές και συστήματα μέσω του διαδικτύου. Οι συσκευές αυτές μπορεί να είναι οικιακές συσκευές ή εξελιγμένα βιομηχανικά εργαλεία.

---

- Βιομηχανικό Διαδίκτυο Πραγμάτων (Industrial Internet of Things - IIoT)

Αναφέρεται σε διασυνδεδεμένους αισθητήρες, εργαλεία, άλλες συσκευές που είναι δικτυωμένες μαζί με υπολογιστικές βιομηχανικές εφαρμογές, στη βιομηχανοποίηση και στη διαχείριση ενέργειας. Η συνδεσιμότητα αυτή επιτρέπει τη συλλογή, ανταλλαγή και ανάλυση δεδομένων, διευκολύνοντας ενδεχόμενες βελτιώσεις στην παραγωγικότητα και την αποδοτικότητα καθώς παρέχει και άλλα οικονομικά οφέλη. Το IIoT είναι μια εξέλιξη ενός κατακεντρωμένου συστήματος ελέγχου (Distributed Control System – DCS) που επιτρέπει υψηλότερο βαθμό αυτοματοποίησης με τη χρήση υπολογιστικού νέφους και βελτιστοποιημένου ελέγχου των διαδικασιών.

- Έξυπνη Κατασκευαστική (Smart Manufacture)

Είναι μια ευρεία κατηγορία κατασκευαστικής που χρησιμοποιεί υπολογιστική κατασκευαστική, υψηλά επίπεδα προσαρμοστικότητας, ραγδαίες αλλαγές στο σχεδιασμό, τεχνολογίες ψηφιακών πληροφοριών και πιο ευέλικτες τεχνικές εκπαίδευσης του εργατικού δυναμικού. Άλλοι στόχοι που προσπαθεί να πετύχει η έξυπνη κατασκευαστική είναι οι γρήγορες αλλαγές στα επίπεδα παραγωγής με βάση τη ζήτηση, βελτιστοποίηση της αλυσίδας εφοδιασμού, αποδοτική παραγωγή και δυνατότητα ανακύκλωσης. Ο ευρύς ορισμός της έξυπνης κατασκευαστικής καλύπτει πολλές διαφορετικές τεχνολογίες. Μερικές από αυτές είναι η δυνατότητα επεξεργασίας μεγάλων δεδομένων (big data), συσκευές και υπηρεσίες βιομηχανικής συνδεσιμότητας και προηγμένη ρομποτική.

- Υπολογιστική Νέφους (Cloud Computing)

Είναι η κατά παραγγελία διαθεσιμότητα υπολογιστικών πόρων, όπως χώρου αποθήκευσης δεδομένων (cloud storage) και υπολογιστικής ισχύος.

- Ψηφιακό Δίδυμο (Digital Twin)

Είναι ένα εικονικό αντίγραφο ενός φυσικού αντικειμένου ή συστήματος. Καταγράφει σε πραγματικό χρόνο τα δεδομένα και τη συμπεριφορά του φυσικού ομολόγου. Επιτρέπει προσομοιώσεις, προβλέψεις και βελτιστοποιήσεις του φυσικού αντικειμένου ή συστήματος.

- Προηγμένα Ρομποτικά Συστήματα (Advanced Robotics)



---

Σε σύγκριση με τα συμβατικά ρομπότ, τα προηγμένα ρομπότ έχουν καλύτερη αντίληψη, δυνατότητα ενσωμάτωσης, προσαρμοστικότητα και φορητότητα. Αυτές οι βελτιώσεις επιτρέπουν την ταχύτερη εγκατάσταση, ανάθεση, αναδιαμόρφωση καθώς και πιο αποδοτικές και σταθερές λειτουργίες. Το κόστος αυτού του εξελιγμένου εξοπλισμού θα μειωθεί καθώς θα μειώνονται οι τιμές των αισθητήρων και της υπολογιστικής ισχύς και καθώς το λογισμικό θα αντικαθιστά όλο και περισσότερο το υλικό έτσι ώστε να καθοδηγεί τη λειτουργικότητα. Ως εκ τούτου, τα προηγμένα ρομπότ θα μπορούν να εκτελούν πολλές εργασίες πολύ πιο οικονομικά από την προηγούμενη γενιά των αυτοματοποιημένων συστημάτων.

- **Μεγάλα Δεδομένα (Big Data)**

Συνήθως περιλαμβάνουν σύνολα δεδομένων με μεγέθη που καθιστούν δύσκολη την καταγραφή, τη διαχείριση και την επεξεργασία δεδομένων σε ανεκτό χρονικό διάστημα από τα συχνά χρησιμοποιούμενα εργαλεία λογισμικού. Η φιλοσοφία των Big Data περιλαμβάνει μη δομημένα, ημιδομημένα και δομημένα δεδομένα, ωστόσο εστιάζουν πιο πολύ στα μη δομημένα δεδομένα. Τα Big Data απαιτούν ένα σύνολο τεχνικών και τεχνολογιών με νέες μορφές ενοποίησης για να εξαχθούν πληροφορίες από τα σύνολα δεδομένων που είναι ποικιλόμορφα και σύνθετα.

- **Γνωστική Πληροφορική (Cognitive Computing)**

Αναφέρεται σε τεχνολογικές πλατφόρμες οι οποίες σε γενικές γραμμές βασίζονται στις αρχές της τεχνητής νοημοσύνης και της επεξεργασίας σημάτων. Αυτές οι πλατφόρμες περιλαμβάνουν μηχανική μάθηση, συλλογιστική, επεξεργασία φυσικής γλώσσας, αναγνώριση ομιλίας και όρασης (αναγνώριση αντικειμένου), αλληλεπίδραση ανθρώπου-υπολογιστή, διάλογο και αφήγηση μεταξύ άλλων τεχνολογιών.

- **Επαυξημένη και Εικονική Πραγματικότητα (Augmented and Virtual Reality)**

Η επαυξημένη πραγματικότητα – AR πρέπει να περιλαμβάνει ένα συνδυασμό μεταξύ του εικονικού και του πραγματικού περιβάλλοντος, να είναι σε θέση να αλληλεπιδράσει με το άμεσο περιβάλλον, να καταγράφει και να συνδέει πραγματικά και εικονικά αντικείμενα.

---

Τα χαρακτηριστικά της εικονικής πραγματικότητας – VR παρέχουν τόσο σωματική όσο και ψυχολογική παρουσία. Στα πλαίσια αυτά, ο χρήστης απομονώνεται από τον πραγματικό κόσμο σε κάποιο βαθμό, που κυμαίνεται από μερική έως πλήρη απομόνωση στην οποία δεν υπάρχει αλληλεπίδραση με τον έξω κόσμο.

## 1.2 Στόχοι της διπλωματικής εργασίας

Ο σκοπός της παρούσας διπλωματικής εργασίας είναι η εφαρμογή αλγορίθμων μηχανικής μάθησης με στόχο τον εντοπισμό ανωμαλιών σε ένα σύνολο δεδομένων βιομηχανικού συστήματος. Επομένως, γίνεται ανάλυση των δεδομένων με κύριο στόχο την εξαγωγή πληροφοριών που υποβοηθούν τις βιομηχανικές διαδικασίες σε όλα τα επίπεδα καθώς και στη λήψη σωστών και κρίσιμων αποφάσεων.

Πιο συγκεκριμένα, τα δεδομένα που αναλύονται προέρχονται από μια γραμμή μεταφοράς ενός συστήματος ηλεκτρικής ενέργειας. Οι ανάγκες ηλεκτρικής ενέργειας έχουν αυξηθεί εκθετικά κατά τη σύγχρονη εποχή και ο κύριος ρόλος μιας γραμμής μεταφοράς είναι να μεταφέρει ενέργεια από την περιοχή της πηγής στο δίκτυο διανομής. Το σύστημα ηλεκτρικής ενέργειας αποτελείται από τόσα πολλά σύνθετα δυναμικά και αλληλεπιδρώντα στοιχεία που είναι πάντα επιρρεπή σε διαταραχές και ηλεκτρικά σφάλματα.

Εξετάστηκαν δύο σύνολα δεδομένων, το πρώτο σύνολο χρησιμοποιείται για να εξαχθούν πληροφορίες σχετικά με το αν θα εμφανίσει το σύστημα σφάλμα ή όχι. Ενώ, το δεύτερο σύνολο χρησιμοποιείται για να εξαχθούν πληροφορίες σχετικά με τον τύπο του σφάλματος.

Συνολικά, εκτελέστηκαν τρεις διαδικασίες κατηγοριοποίησης και δύο διαδικασίες παραμετροποίησης:

1. Δυαδική κατηγοριοποίηση στο πρώτο σύνολο δεδομένων
2. Πολυταξική κατηγοριοποίηση στο δεύτερο σύνολο δεδομένων
3. Χειρισμός της πολυταξικής κατηγοριοποίησης ως δυαδική
4. Παραμετροποίηση των αλγορίθμων με στόχο τη βελτιστοποίησή τους

---

Οι σταθμοί παραγωγής ηλεκτρικής ενέργειας υψηλής δυναμικότητας και η έννοια του δικτύου απαιτούν την ανίχνευση σφαλμάτων και τη λειτουργία του εξοπλισμού προστασίας στον ελάχιστο δυνατό χρόνο για να παραμείνουν σταθεροί. Τα σφάλματα στις γραμμές μεταφοράς του συστήματος ηλεκτρικής ενέργειας υποτίθεται ότι πρέπει πρώτα να ανιχνεύονται και να κατηγοριοποιούνται σωστά και να αποκαθίστανται στον ελάχιστο δυνατό χρόνο. Επομένως, οι επιτυχείς προβλέψεις θα μπορούσαν να βοηθήσουν στη διάκριση μεταξύ ελαττωματικών και υγιών συστημάτων ηλεκτρικής ενέργειας.

### 1.3 Διάρθρωση κειμένου

Τα επόμενα τρία κεφάλαια έχουν ως στόχο να καλύψουν το θεωρητικό κομμάτι της εργασίας. Πιο συγκεκριμένα στο δεύτερο κεφάλαιο, αναλύονται γενικές έννοιες σχετικές με την Τεχνητή Νοημοσύνη, αλγόριθμοι που θα χρησιμοποιηθούν και μετρικές αξιολόγησης. Έπειτα, στο τρίτο κεφάλαιο παρουσιάζεται μια βιβλιογραφική ανασκόπηση των μεθόδων μηχανικής μάθησης που εφαρμόζονται στην προβλεπτική μηχανική. Στο τέταρτο κεφάλαιο γίνεται μια εισαγωγή και παρουσίαση των τεχνολογιών που χρησιμοποιήθηκαν.

Το βασικό κομμάτι της εργασίας αποτελείται από τα επόμενα τρία κεφάλαια. Στο πέμπτο κεφάλαιο αναλύεται η εφαρμογή και τα δεδομένα που χρησιμοποιήθηκαν. Το έκτο κεφάλαιο σχετίζεται με τις αναλύσεις που υλοποιήθηκαν και τα αποτελέσματα που προκύπτουν. Ενώ στο έβδομο κεφάλαιο εξετάστηκε η περίπτωση παραμετροποίησης των αλγορίθμων με τις καλύτερες επιδόσεις, με κύριο σκοπό τη βελτίωσή τους.

Τέλος, το όγδοο κεφάλαιο αποτελεί τον επίλογο της διπλωματικής εργασίας. Παρουσιάζονται τα βασικά συμπεράσματα της εργασίας και οι μελλοντικές επεκτάσεις.

# Κεφάλαιο 2

## Θεωρητικό Υπόβαθρο

Το κεφάλαιο αυτό παρουσιάζει και αναλύει το θεωρητικό κομμάτι της εργασίας. Πιο συγκεκριμένα, παρέχεται μια εννοιολογική διάκριση μεταξύ σχετικών όρων και εννοιών καθώς και η περιγραφή των αλγορίθμων μηχανικής μάθησης που χρησιμοποιήθηκαν για την εξαγωγή των αποτελεσμάτων.

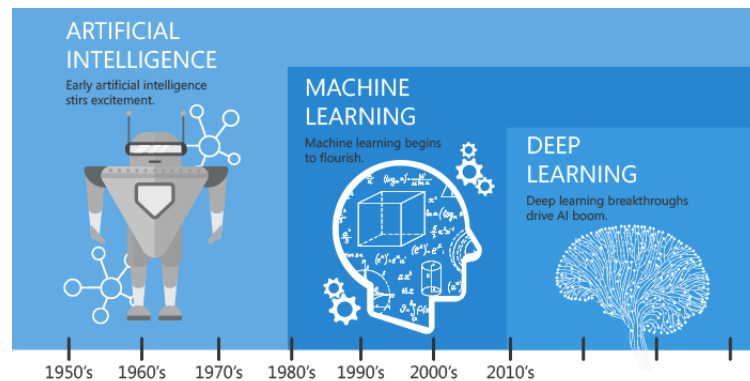
### 2.1 Ορισμοί βασικών εννοιών

Τα ευφυή συστήματα που προσφέρουν δυνατότητες τεχνητής νοημοσύνης συχνά βασίζονται στη μηχανική μάθηση. Με τον όρο μηχανική μάθηση περιγράφεται η ικανότητα των συστημάτων να μαθαίνουν από τα δεδομένα εκπαίδευσης έτσι ώστε να αυτοματοποιηθεί η διαδικασία ανάπτυξης αναλυτικών μοντέλων με τελικό στόχο την επίλυση σχετικών εργασιών. Η βαθιά μάθηση (deep learning) είναι μια έννοια μηχανικής μάθησης που βασίζεται στα τεχνητά νευρωνικά δίκτυα (artificial neural networks). Σε πολλές εφαρμογές, τα μοντέλα βαθιάς μάθησης ξεπερνούν τα ρηχά μοντέλα μηχανικής μάθησης (shallow machine learning models) και τις παραδοσιακές προσεγγίσεις ανάλυσης δεδομένων.

Σε αυτή την ενότητα, συνοψίζονται οι βασικές αρχές της τεχνητής νοημοσύνης, της μηχανικής και της βαθιάς μάθησης για να δημιουργηθεί μια ευρύτερη κατανόηση της μεθοδολογίας που χρησιμοποιείται για την υποστήριξη των σημερινών ευφυών συστημάτων. Δίνεται μεγαλύτερη έμφαση στο κομμάτι της μηχανικής μάθησης το οποίο είναι και το αντικείμενο αυτής της εργασίας.

Σύμφωνα με το Σχήμα 2.1, η τεχνητή νοημοσύνη αποτελείται από τεχνικές που επιτρέπουν στους υπολογιστές να μιμούνται την ανθρώπινη νοημοσύνη. Η μηχανική μάθηση είναι ένα υποσύνολο της τεχνητής νοημοσύνης που περιλαμβάνει αφη-

Σχήμα 2.1: Εξέλιξη Τεχνητής Νοημοσύνης, Μηχανικής Μάθησης και Βαθιάς Μάθησης [34]



ρημένες στατιστικές τεχνικές, οι οποίες επιτρέπουν στις μηχανές να αποκτήσουν εμπειρία σε διάφορες εργασίες και να τις βελτιώνουν. Η βαθιά μάθηση είναι ένα υποσύνολο της μηχανικής μάθησης που αποτελείται από αλγορίθμους οι οποίοι επιτρέπουν στο λογισμικό να εκπαιδευτεί ώστε να εκτελεί εργασίες όπως αναγνώριση λόγου και εικόνας. Αυτό επιτυγχάνεται με την έκθεση πολυεπίπεδων νευρωνικών δικτύων σε τεράστιες ποσότητες δεδομένων.

- Τεχνητή Νοημοσύνη

Στόχος της Τεχνητής Νοημοσύνης είναι η δημιουργία μιας μηχανής η οποία θα συμπεριφέρεται όπως ένας συνηθισμένος άνθρωπος και θα αποτελεί μια βελτιωμένη εκδοχή της τρέχουσας συμπεριφοράς των μηχανημάτων για την αντιμετώπιση σύνθετων εργασιών. Μέσω της έρευνας στο πεδίο της ΤΝ ο άνθρωπος μπόρεσε να καταλάβει την ευφυή συμπεριφορά του. Το ανθρώπινο γένος έχει μια ενδιαφέρουσα προσέγγιση στην επίλυση προβλημάτων η οποία βασίζεται στην αφηρημένη σκέψη, σε υψηλού επιπέδου αναλυτικό σκεπτικό και αναγνώριση προτύπων. Η τεχνητή νοημοσύνη βοηθάει στην κατανόηση των διεργασιών μέσω της αναδημιουργίας τους, επιτρέποντας παράλληλα την ενίσχυση των ανθρώπινων ικανοτήτων. Παρόλο που είναι ένα σχετικά καινούργιο ανεξάρτητο πεδίο μελέτης, έχει ρίζες στο παρελθόν. Ξεκίνησε περίπου πριν 2400 χρόνια όταν ο Έλληνας φιλόσοφος Αριστοτέλης εφηύρε την έννοια του λογικού συλλογισμού και συνεχίστηκε με τους Leibniz και Newton. Ακόμη, ο George Boole ανέπτυξε την άλγεβρα Boole τον 19ο αιώνα, η οποία έβαλε τις βάσεις των υπολογιστικών κυκλωμάτων. Παρόλα αυτά, η κεντρική ιδέα μιας συλλογιστικής μηχανής ήρθε από τον Alan Turing ο οποίος πρότεινε το τεστ

---

Turing το 1950, το οποίο δίνει τον ορισμό της νοημοσύνης σε μια μηχανή. Ο όρος Τεχνητή Νοημοσύνη επινοήθηκε από τον John MacCarthy το 1956. Στην απλούστερη μορφή της η τεχνητή νοημοσύνη είναι ένα πεδίο το οποίο συνδυάζει την επιστήμη των υπολογιστών και ισχυρά σύνολα δεδομένων, ώστε να επιλύσει προβλήματα. Παρακάτω αναλύονται οι τύποι Τεχνητής Νοημοσύνης:

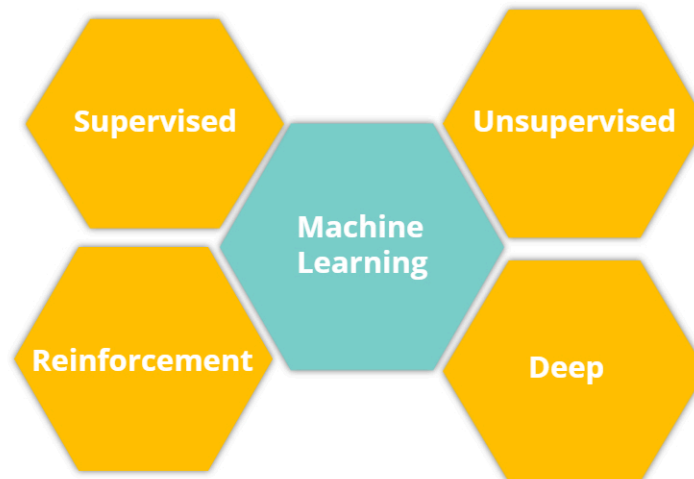
1. Αδύναμη Τεχνητή Νοημοσύνη – επίσης ονομάζεται και ρηχή (Narrow AI ή Artificial Narrow Intelligence – ANI). Ο τύπος αυτός εκπαιδεύεται και στοχεύει στην επίτευξη συγκεκριμένων εργασιών. Κινεί το μεγαλύτερο μέρος των εφαρμογών τεχνητής νοημοσύνης που υπάρχουν σήμερα όπως Siri, Alexa, IBM Watson και αυτόνομα οχήματα.
  2. Ισχυρή Τεχνητή Νοημοσύνη – είναι δημιούργημα της Τεχνητής Γενικής Νοημοσύνης (Artificial General Intelligence – AGI) και της Τεχνητής Υπερνοημοσύνης (Artificial Super Intelligence – ASI). Η AGI είναι μια θεωρητική μορφή τεχνητής νοημοσύνης σύμφωνα με την οποία μια μηχανή θα έχει ισάξια νοημοσύνη με την ανθρώπινη. Αυτό θα επιτευχθεί μέσω της συνειδητής αυτογνωσίας η οποία θα δύναται να λύνει προβλήματα, να μαθαίνει και να σχεδιάζει μελλοντικά σχέδια. Η ASI θα υπερβαίνει τη νοημοσύνη και την ικανότητα του ανθρώπινου εγκεφάλου. Παρόλα αυτά, εξακολουθεί να είναι εντελώς θεωρητική χωρίς πρακτικά παραδείγματα.
- Μηχανική Μάθηση

Είναι ένας κλάδος της Τεχνητής Νοημοσύνης και της επιστήμης των υπολογιστών που επικεντρώνεται στη χρήση δεδομένων και αλγορίθμων έτσι ώστε να μιμηθεί τον τρόπο με τον οποίο μαθαίνει ο άνθρωπος, βελτιώνοντας σταδιακά την ακρίβεια του. Μέσω της χρήσης στατιστικών μεθόδων, οι αλγόριθμοι εκπαιδεύονται να κάνουν κατηγοριοποιήσεις ή προβλέψεις, εξάγοντας βασικές γνώσεις μέσω της εξόρυξης δεδομένων. Αυτές οι γνώσεις καθορίζουν στη συνέχεια τις αποφάσεις που λαμβάνονται μέσα σε εφαρμογές και επιχειρήσεις έχοντας ως κύριο στόχο την επιρροή των βασικών μέτρων ανάπτυξης. Τα συστήματα μηχανικής μάθησης μπορούν να κατηγοριοποιηθούν ανάλογα με το πλήθος και τον τύπο επίβλεψης που λαμβάνουν καθ' όλη τη διάρκεια της εκπαίδευσης. Υπάρχουν τέσσερες κύριες κατηγορίες όπως απεικονίζονται και

---

στο Σχήμα 2.2: επιβλεπόμενη, μη-επιβλεπόμενη, ενισχυτική και βαθιά μάθηση. Η κάθε κατηγορία περιγράφεται εν συντομία παρακάτω[18].

Σχήμα 2.2: Κατηγορίες Μηχανικής Μάθησης



### 1. Επιβλεπόμενη Μάθηση (Supervised Learning)

Είναι μια μέθοδος με την οποία ένας ειδικός εισάγει γνωστές εξόδους για συγκεκριμένες εισόδους έτσι ώστε να εκπαιδευτεί ο αλγόριθμος και χρησιμοποιείται ευρέως για κατηγοριοποίηση και παλινδρόμηση. Έτσι, η εποπτευόμενη μάθηση χρησιμοποιείται συνήθως σε περιπτώσεις όπου υπάρχουν δείγματα στα οποία έχουν επισημανθεί μία ή περισσότερες ετικέτες. Κάποιοι από τους πιο δημοφιλείς αλγόριθμους επιβλεπόμενης μάθησης είναι τα Νευρωτικά Δίκτυα (Artificial neural network - ANN) και οι Μηχανές Υποστήριξης Διανυσμάτων (Support Vector Machine - SVM).

### 2. Μη-Επιβλεπόμενη Μάθηση (Unsupervised learning)

Δεν παρέχεται ανατροφοδότηση από κανέναν και ο αλγόριθμος βρίσκει μοτίβα σε άγνωστα σύνολα δεδομένων. Με αυτόν τον τρόπο γίνεται χρήση δεδομένων που δεν έχουν κάποια ετικέτα για εκπαίδευση του αλγορίθμου. Ο πιο δημοφιλής και γνωστός αλγόριθμος μη-επιβλεπόμενης μάθησης είναι η Ανάλυση Κύριων Συνιστωσών (Principal Component Analysis

---

- PCA) και χρησιμοποιείται κυρίως για σκοπούς παρακολούθησης.

### 3. Ενισχυτική Μάθηση (Reinforcement learning)

Η ενισχυτική μάθηση απαιτεί διαδοχικές ενέργειες στις οποίες δοκιμάζει το αποτέλεσμα τους, έτσι ώστε να επιλεχθούν αυτές που ταιριάζουν καλύτερα στο πρόβλημα που αντιμετωπίζετε. Ως εκ τούτου, η ενισχυτική μάθηση διαφέρει σημαντικά από τις υπόλοιπες κατηγορίες οι οποίες βασίζονται στην αξιοποίηση παλαιότερων δεδομένων, δημιουργώντας πληροφορία από προηγούμενες αποφάσεις και ανταμοιβές.

### 4. Βαθιά Μάθηση (Deep learning)

Χρησιμοποιούνται πολλαπλά επίπεδα για να δημιουργηθεί ένα ANN, το οποίο είναι σε θέση να λαμβάνει έξυπνες αποφάσεις, χειρίζεται μεγάλες ποσότητες δεδομένων υψηλής πολυπλοκότητας, χωρίς καμία ανθρώπινη παρέμβαση. Μερικοί αλγόριθμοι βαθιάς μάθησης είναι τα Συνελικτικά Νευρωνικά Δίκτυα (Convolutional Neural Network - CNN), η Περιορισμένη Μηχανή Boltzmann (Restricted Boltzmann Machine - RBM) και οι Αυτό – κωδικοποιητές (autoencoder - AE).

## 2.2 Αλγόριθμοι επιβλεπόμενης μάθησης

Στην ενότητα αυτή περιγράφονται και αναλύονται οι αλγόριθμοι επιβλεπόμενης μάθησης για το πως λειτουργούν και πότε χρησιμοποιούνται. Πιο συγκεκριμένα, οι αλγόριθμοι αυτοί ανήκουν στις κατηγορίες της κατηγοριοποίησης και παλινδρόμησης και επιπρόσθετα σε επόμενο κεφάλαιο εφαρμόζονται και εξετάζονται τα μοντέλα σε διάφορες εφαρμογές.

Η Κατηγοριοποίηση (Classification) και η Παλινδρόμηση (Regression) είναι δύο από τις πιο γνωστές εφαρμογές στη μηχανική μάθηση. Οι αλγόριθμοι κατηγοριοποίησης χρησιμοποιούνται στην πρόβλεψη μιας κατηγορηματικής ετικέτας ή κλάσης για μια δοσμένη είσοδο, ενώ οι αλγόριθμοι παλινδρόμησης χρησιμοποιούνται στην πρόβλεψη συνεχών τιμών για δεδομένη είσοδο.

Οι αλγόριθμοι κατηγοριοποίησης χρησιμοποιούνται σε ένα ευρύ φάσμα εφαρμογών, από την αναγνώριση ανεπιθύμητης αλληλογραφίας μέχρι σε ιατρικές διαγνώσεις. Μερικοί από τους πιο γνωστούς αλγόριθμους είναι οι: Λογιστική Παλινδρό-



---

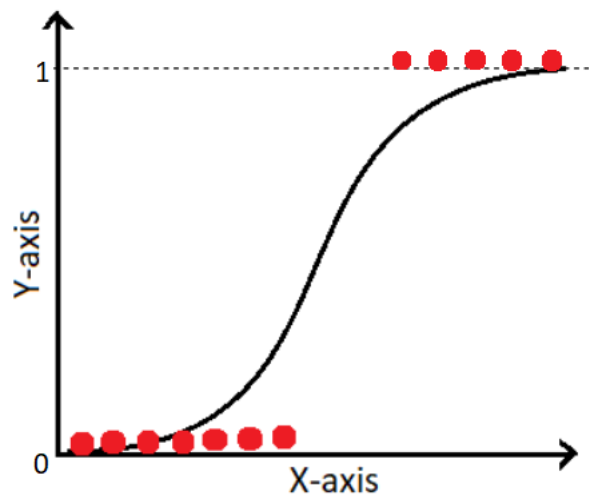
μηση (Logistic Regression), Δένδρα Απόφασης (Decision Trees) και Κ-Πλησιέστεροι Γείτονες (K-Nearest Neighbors).

Από την άλλη πλευρά, οι αλγόριθμοι παλινδρόμησης χρησιμοποιούνται στην πρόβλεψη αριθμητικών τιμών όπως μετοχών, θερμοκρασίας ή βάρους. Ορισμένοι από τους πιο γνωστούς αλγόριθμους παλινδρόμησης είναι οι Γραμμική Παλινδρόμηση (Linear Regression), Πολυωνυμική Παλινδρόμηση (Polynomial Regression) και Παλινδρόμηση με Διανύσματα Υποστήριξης (Support Vector Regression).

- Λογιστική Παλινδρόμηση (Logistic Regression)

Ο Logistic Regression είναι μία ισχυρή στατιστική μέθοδος που χρησιμοποιείται ευρέως στη μηχανική μάθηση και στην ανάλυση δεδομένων κυρίως σε εφαρμογές κατηγοριοποίησης[30]. Ο αλγόριθμος χρησιμοποιεί μια λογιστική συνάρτηση, την Sigmoid, για να μοντελοποιήσει την πιθανότητα να προκύψει μια συγκεκριμένη κλάση ή ένα γεγονός.

Σχήμα 2.3: Logistic Regression [40]



Η λογιστική συνάρτηση είναι μια καμπύλη σχήματος S όπως παρατηρείται και στο Σχήμα 2.3 που μπορεί να λάβει ως είσοδο οποιοδήποτε πραγματικό αριθμό και να τον αντιστοιχήσει σε μια τιμή μεταξύ του 0 και 1. Σε επόμενο στάδιο, η προβλεπόμενη πιθανότητα αντιστοιχίζεται σε κλίμακα 0.5 για να προβλεφθεί η κατηγορία. Ο Logistic Regression βρίσκει τις καλύτερες παραμέτρους ή αλλιώς συντελεστές, για τη λογιστική συνάρτηση, μειώνοντας το σφάλμα μεταξύ της προβλεπόμενης και της πραγματικής τιμής. Η μείωση αυτή επιτυγχάνεται με τη χρήση του Gradient Descent, ο οποίος ανήκει στην

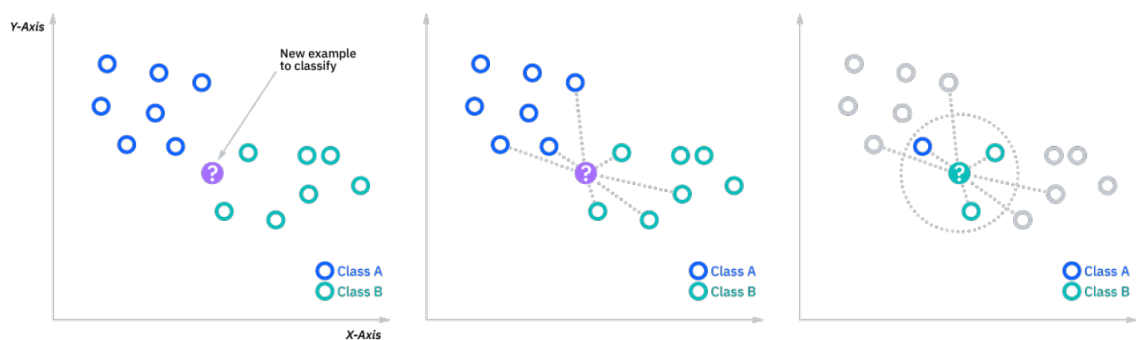
κατηγορία των αλγορίθμων βελτιστοποίησης. Οι συντελεστές της λογιστικής συνάρτησης αντιπροσωπεύουν τη σημασία του κάθε χαρακτηριστικού και την επίδραση τους στο αποτέλεσμα.

Ο Logistic Regression χρησιμοποιείται ευρέως σε διάφορους τομείς όπως σε βαθμολόγηση πιστώσεων, σε ιατρικές διαγνώσεις και έρευνες αγοράς. Οι λόγοι για τους οποίους ο αλγόριθμος είναι δημοφιλής έχουν να κάνουν με την απλότητα του ως προς τη χρήση αλλά και ως προς την ερμηνεία καθώς μπορεί να χειριστεί τόσο γραμμικές όσο και μη γραμμικές σχέσεις μεταξύ των ανεξάρτητων και εξαρτημένων μεταβλητών. Επιπλέον, είναι ανθεκτικός σε θόρυβο και ακραίες τιμές και μπορεί να διαχειριστεί δεδομένα υψηλών διαστάσεων.

- K-Πλησιέστεροι Γείτονες (K-Nearest Neighbors)

Ο KNeighbors χρησιμοποιείται για κατηγοριοποίηση και παλινδρόμηση. Η βασική ιδέα πίσω από τον KNeighbors είναι να προβλέπει τη κλάση ή την τιμή ενός στοιχείου του συνόλου δεδομένων βασισμένο στους k κοντινούς γείτονες.

Σχήμα 2.4: K-Nearest Neighbors[7]



Αρχικά, ο αλγόριθμος αποθηκεύει όλα τα δεδομένα εκπαίδευσης με τις αντίστοιχες κατηγορίες τους. Στη συνέχεια, όταν ζητηθεί να κατηγοριοποιηθεί ένα καινούργιο στοιχείο ο αλγόριθμος υπολογίζει την απόσταση μεταξύ του καινούργιου στοιχείου και όλων των αποθηκευμένων δεδομένων που χρησιμοποι-

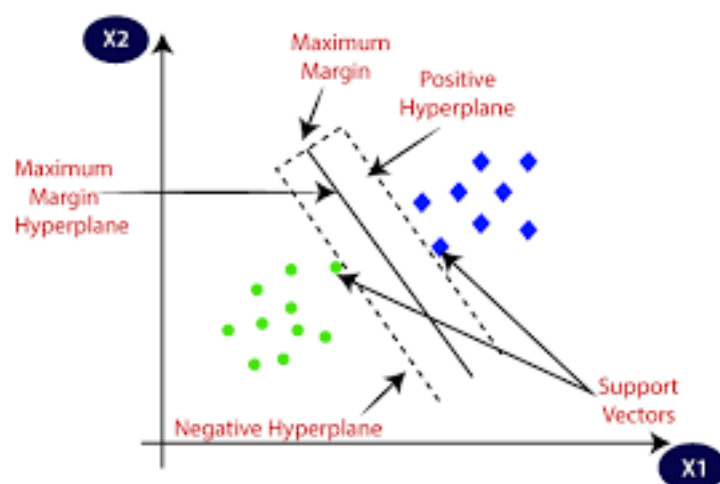
ήθηκαν για εκπαίδευση. Έπειτα, επιλέγει τα  $k$  στοιχεία εκπαίδευσης που είναι πιο κοντά στο νέο στοιχείο και χρησιμοποιεί τις αντίστοιχες κατηγορίες τους για να προβλέψει τη κατηγορία του καινούργιου στοιχείου[16]. Η διαδικασία που αναφέρθηκε απεικονίζεται στο Σχήμα 2.4. Η τιμή της παραμέτρου  $k$  επιλέγεται από το χρήστη, συνήθως η επιλογή γίνεται με βάση το συγκεκριμένο πρόβλημα και το μέγεθος του συνόλου.

Ο κατηγοριοποιητής KNeighbors είναι απλός και εύκολος στη κατανόηση. Δεν κάνει κάποια υπόθεση σχετικά με την κατανομή των δεδομένων και δεν είναι ευαίσθητος ως προς την κλίμακα των χαρακτηριστικών. Ωστόσο, το κόστος υπολογισμού μπορεί να είναι αρκετά ακριβό και μπορεί επίσης να γίνει πιο ευαίσθητος στην επιλογή της παραμέτρου  $k$ . Άλλο ένα μειονέκτημα του είναι η τάση του να είναι ευαίσθητο σε θορυβώδη δεδομένα και άσχετα χαρακτηριστικά.

- Παλινδρόμηση με Διανύσματα Υποστήριξης (Support Vector Machines)

Ο SVM χρησιμοποιείται σε προβλήματα κατηγοριοποίησης. Είναι ένας ισχυρός και ευέλικτος αλγόριθμος που χειρίζεται τόσο γραμμικά όσο και μη γραμμικά δεδομένα. Η βασική ιδέα πίσω από τον SVM είναι η εύρεση μιας γραμμής ή ενός επίπεδου που θα διαχωρίζει καλύτερα τα σημεία διαφορετικών κλάσεων.

Σχήμα 2.5: Support Vector Machines [42]



Ο SVM βρίσκει το υπερεπίπεδο μεγιστοποιώντας το περιθώριο όπως φαίνεται στο Σχήμα 2.5, το οποίο ορίζεται ως η απόσταση μεταξύ του υπερεπίπεδου και των πλησιέστερων σημείων της κάθε κλάσης[20]. Τα πλησιέστερα σημεία

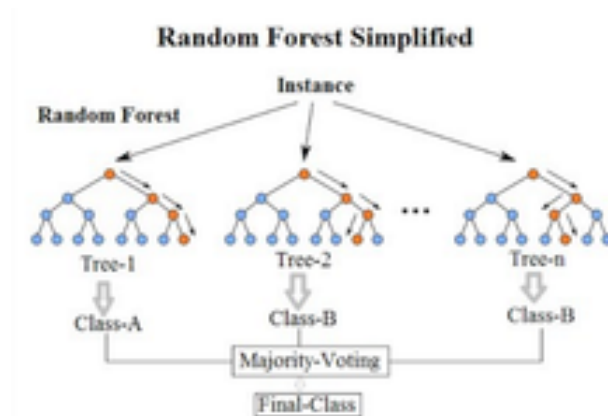
είναι γνωστά επίσης και ως διανύσματα υποστήριξης. Με τη μεγιστοποίηση του περιθωρίου, ο SVM διασφαλίζει ότι το όριο απόφασης απέχει αρκετά από τα κοντινά σημεία δεδομένων, παρέχοντας καλές αποδόσεις γενίκευσης του μοντέλου. Επιτρέπει επιπλέον στον αλγόριθμο να είναι ανθεκτικός στην παρουσία ακραίων τιμών.

Ο SVM είναι ένας ισχυρός αλγόριθμος ο οποίος μπορεί να χειριστεί μη γραμμικά όρια αποφάσεων χρησιμοποιώντας το τέχνασμα του πυρήνα. Η τεχνική αυτή επιτρέπει στον SVM να αντιστοιχεί τα δεδομένα εισόδου σε ένα χώρο υψηλότερων διαστάσεων όπου τα δεδομένα γίνονται γραμμικά διαχωρίσιμα. Οι πιο συχνά χρησιμοποιούμενες συναρτήσεις πυρήνα είναι οι γραμμικοί, πολυωνυμικοί και οι Πυρήνες Ακτινικής Συνάρτησης Βάσης (Radial basis function kernel - RBF). Παρόλα αυτά, ο SVM μπορεί να είναι ευαίσθητος στην επιλογή της συνάρτησης πυρήνα και μπορεί επίσης ο υπολογισμός να είναι ιδιαίτερα δαπανηρός για μεγάλα σύνολα δεδομένων.

- Δάση τυχαίας απόφασης (Random Forest)

Ο Random Forest χρησιμοποιείται σε προβλήματα κατηγοριοποίησης και παλινδρόμησης. Λειτουργεί συνδυάζοντας πολλαπλά δέντρα αποφάσεων για να φτάσει στην τελική πρόβλεψη. Τα δέντρα απόφασης δημιουργούνται από ένα τυχαίο υποσύνολο των δεδομένων εκπαίδευσης και ένα τυχαίο υποσύνολο των χαρακτηριστικών και κάθε δέντρο από το δάσος ψηφίζει για την τελική πρόβλεψη, Σχήμα 2.6. Ο αλγόριθμος συγκεντρώνει τα αποτελέσματα όλων των δέντρων έτσι ώστε να υπολογίσει το τελικό αποτέλεσμα[25].

Σχήμα 2.6: Random Forest [47]



---

Ένα από τα βασικά πλεονεκτήματα του Random Forest είναι η ικανότητα του να διαχειρίζεται μεγάλα και πολύπλοκα σύνολα δεδομένων, καθώς και να διαχειρίζεται δεδομένα που λείπουν ή που είναι θορυβώδη. Αυτό έγκειται στο γεγονός ότι ο αλγόριθμος είναι ικανός να συλλάβει μη γραμμικές σχέσεις στα δεδομένα και να κάνει ακριβείς προβλέψεις ακόμα και με περιορισμένα δεδομένα. Επιπρόσθετα, ο αλγόριθμος είναι λιγότερο επιρρεπής στην υπερπροσαρμογή (overfitting) από τα κλασσικά δέντρα απόφασης, η οποία μπορεί να συμβεί όταν ένα μοναδικό δέντρο εκπαιδεύεται σε ένα μεγάλο σύνολο δεδομένων και γίνεται υπερβολικά πολύπλοκο.

Ο Random Forest παρέχει επίσης μια βαθμολόγηση της σημαντικότητας των χαρακτηριστικών, η οποία επιτρέπει στους χρήστες να καθορίσουν τη σχετική σημασία του κάθε χαρακτηριστικού στο σύνολο δεδομένων. Αυτό μπορεί να είναι χρήσιμο στην επιλογή των χαρακτηριστικών όπως και στην εμφάνιση των χαρακτηριστικών που έχουν τη μεγαλύτερη επίδραση στη μεταβλητή στόχο. Ο Random Forest είναι ένα πολύτιμο εργαλείο τόσο για την εξερεύνηση δεδομένων όσο και για την προγνωστική μοντελοποίηση, και έχει χρησιμοποιηθεί ευρέως σε διάφορους κλάδους, όπως η χρηματοδότηση, η υγειονομική περίθαλψη και το μάρκετινγκ. Συνολικά, ο Random Forest είναι ένας ισχυρός και εύελικτος αλγόριθμος μηχανικής μάθησης που μπορεί να προσφέρει υψηλή ακρίβεια και ισχυρά αποτελέσματα για πολλούς τύπους προβλημάτων.

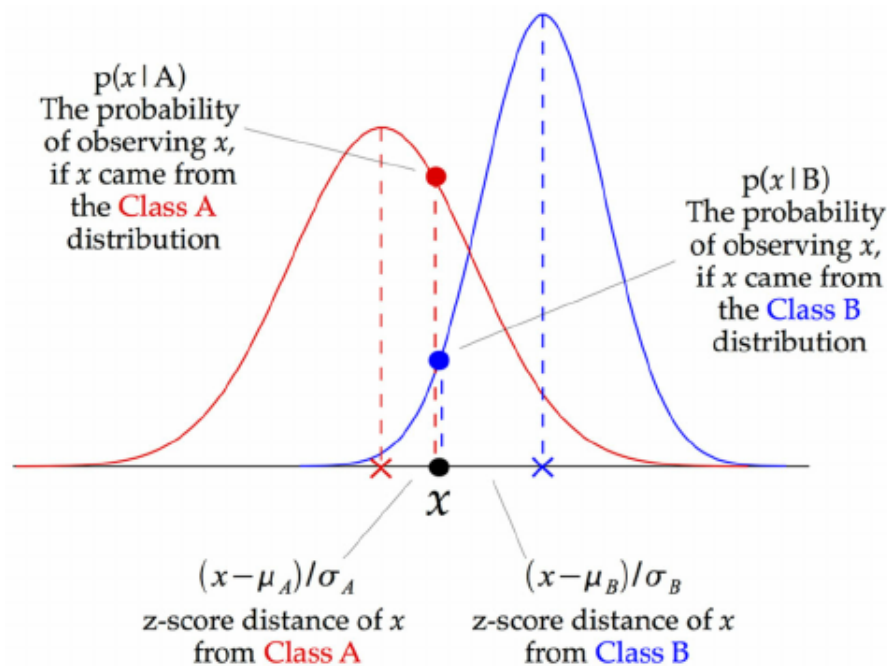
- Gaussian Naive Bayes

Ο GaussianNB είναι ένας δημοφιλής αλγόριθμος που χρησιμοποιείται για κατηγοριοποίηση στη μηχανική μάθηση. Στηρίζεται στην υπόθεση ότι τα δεδομένα παράγονται από μια κατανομή Gauss. Ο αλγόριθμος χρησιμοποιεί αυτήν την υπόθεση για να εκτιμήσει την πιθανότητα ένα δεδομένο δείγμα να ανήκει σε μια συγκεκριμένη κατηγορία, όπως παρουσιάζεται και στο Σχήμα 2.7. Ο GaussianNB είναι ιδιαίτερα χρήσιμος για σύνολα δεδομένων όπου τα χαρακτηριστικά έχουν συνεχείς τιμές και οι κλάσεις είναι καλά διαχωρισμένες από ένα γραμμικό όριο[43].

Ένα από τα κύρια πλεονεκτήματα του GaussianNB είναι η απλότητα του. Ο αλγόριθμος απαιτεί ελάχιστο συντονισμό και είναι εύκολος στην υλοποίηση,

πράγμα που τον κάνει μια πολύ καλή επιλογή για αρχάριους στη μηχανική μάθηση. Επιπρόσθετα, είναι υπολογιστικά αποδοτικός και χρειάζεται ένα σχετικά μικρό αριθμό δεδομένων για την παραγωγή ακριβών αποτελεσμάτων. Αυτό τον κάνει μια αρκετά καλή επιλογή σε περιπτώσεις για δεδομένα με μεγάλο αριθμό χαρακτηριστικών ή σε περιπτώσεις όπου οι υπολογιστικοί πόροι είναι περιορισμένοι.

Σχήμα 2.7: Gaussian Naive Bayes[41]



Παρόλα αυτά, ο GaussianNB έχει επίσης ορισμένους περιορισμούς. Υποθέτει ότι τα δεδομένα ακολουθούν μια κατανομή Gauss, κάτι που μπορεί να μην ισχύει πάντα. Στην περίπτωση όπου τα δεδομένα δεν είναι γκαουσιανά, ο αλγόριθμος ενδέχεται να μην αποδίδει τόσο καλά όσο άλλοι αλγόριθμοι κατηγοριοποίησης. Επιπρόσθετα, ο GaussianNB δεν είναι κατάλληλος για σύνολα δεδομένων όπου οι κλάσεις δεν είναι καλά διαχωρισμένες μέσω ενός γραμμικού ορίου, καθώς κάνει κάποιες υποθέσεις σχετικά με την υποκείμενη κατανομή των δεδομένων που μπορεί να μην πληρούνται.

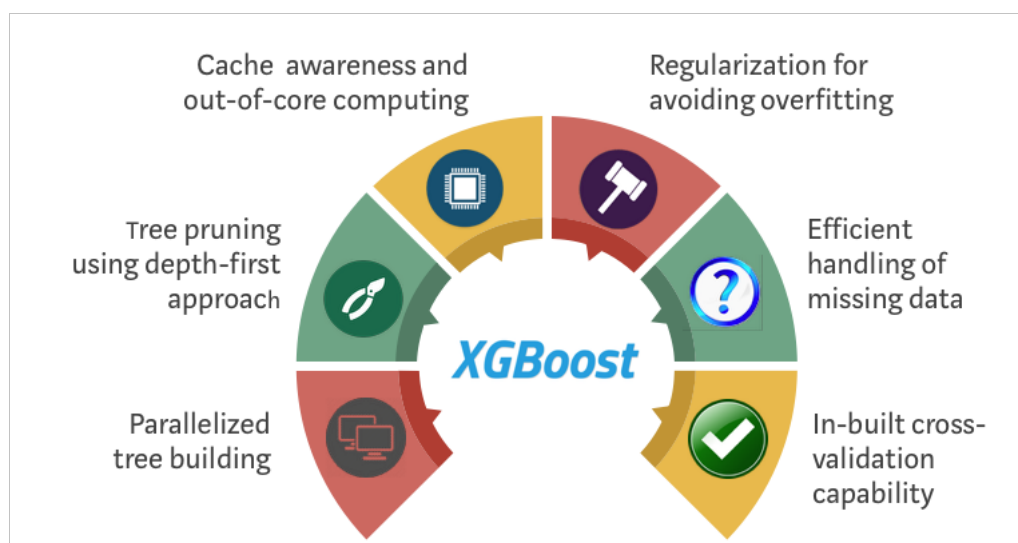
- eXtreme Gradient Boosting

Ο XGBoost είναι ένας ισχυρός και δημοφιλής αλγόριθμος για κατηγοριοποίηση και παλινδρόμηση στη μηχανική μάθηση. Πρόκειται για μια εφαρμογή του Gradient Boosting, μιας συνολικής μεθόδου μάθησης που συνδυάζει πολλαπλά

αδύναμα μοντέλα για να παραχθεί ένα ισχυρό μοντέλο [24]. Ο XGB χρησιμοποιεί δέντρα απόφασης στη βάση της εκμάθησης του και ο αλγόριθμος επαναληπτικά βελτιώνει το μοντέλο προσαρμόζοντας τα βάρη των δέντρων απόφασης για να μειώσει το σφάλμα.

Ένα από τα κύρια πλεονεκτήματα του XGB είναι η υψηλή ακρίβεια και η ικανότητα να διαχειρίζεται μεγάλα σύνολα δεδομένων με μεγάλο αριθμό χαρακτηριστικών. Είναι επίσης ανθεκτικός στο θόρυβο και στις ακραίες τιμές κάτι που τον κάνει μια αρκετά καλή επιλογή για σύνολα δεδομένων με θόρυβο και κενές μεταβλητές. Επιπλέον, ο κατηγοριοποιητής διαθέτει έναν ενσωματωμένο μηχανισμό επιλογής χαρακτηριστικών ο οποίος βοηθά στον εντοπισμό των πιο σημαντικών χαρακτηριστικών για τη διαδικασία της κατηγοριοποίησης, κάτι που μπορεί να είναι χρήσιμο όταν πρόκειται για σύνολα δεδομένων υψηλών διαστάσεων. Στο Σχήμα 2.8 παρουσιάζονται τα πλεονεκτήματα του αλγορίθμου.

Σχήμα 2.8: eXtreme Gradient Boosting[44]



Ωστόσο, ο XGB έχει επίσης ορισμένους περιορισμούς. Είναι ευαίσθητος στην κλίμακα των χαρακτηριστικών, επομένως είναι σημαντική η προ-επεξεργασία των δεδομένων έτσι ώστε να είναι βέβαιο ότι τα δεδομένα βρίσκονται στην ίδια κλίμακα. Επιπρόσθετα μπορεί να είναι υπολογιστικά δαπανηρός και να απαιτεί αρκετή μνήμη ειδικά όταν πρόκειται για μεγάλα σύνολα δεδομένων. Τέλος, ο XGBoost δεν είναι κατάλληλος για σύνολα δεδομένα με πολύ λίγες παρατηρήσεις ή σύνολα δεδομένων όπου οι κλάσεις δε διαχωρίζονται καλά

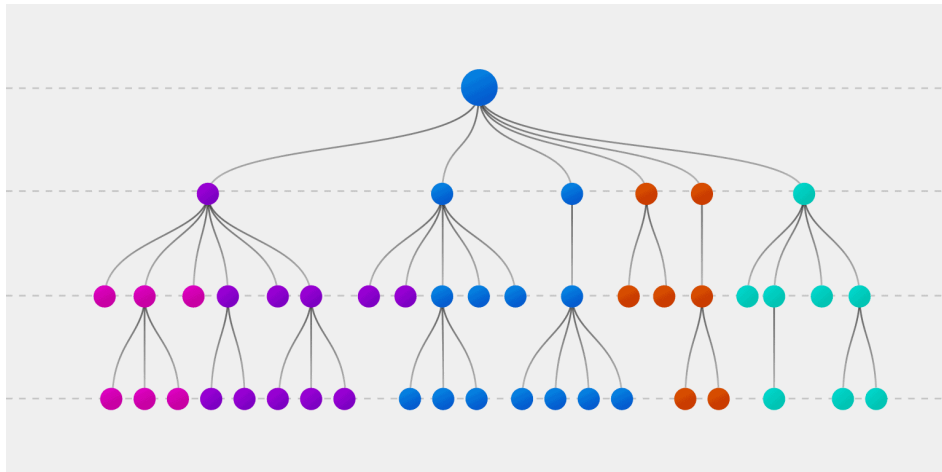
---

από ένα γραμμικό όριο.

- Δέντρα Απόφασης (Decision Tree)

Ο Decision Tree είναι ένας αλγόριθμος μηχανικής μάθησης με επίβλεψη που χρησιμοποιείται σε προβλήματα κατηγοριοποίησης και παλινδρόμησης στη μηχανική μάθηση. Ανήκει στην κατηγορία των αλγορίθμων που βασίζονται στα δέντρα. Ξεκινάει με έναν κόμβο, που ονομάζεται ρίζα, και διαχωρίζει τα δεδομένα σε υποσύνολα με βάση κάποια συγκεκριμένα κριτήρια [39]. Τα υποσύνολα και αυτά διαδοχικά διαχωρίζονται σε περισσότερα υποσύνολα έως ότου κατηγοριοποιηθούν τα δεδομένα ή ο αλγόριθμος φτάσει κάποιο τερματικό κριτήριο, Σχήμα 2.9.

Σχήμα 2.9: Decision Tree [32]



Τα δέντρα απόφασης χρησιμοποιούν μια ευρετική προσέγγιση για να επιλέξουν το καλύτερο χαρακτηριστικό που θα διαχωρίσει τα δεδομένα σε κάθε κόμβο. Η ευρετική προσέγγιση βασίζεται συνήθως στο κέρδος πληροφορίας, το οποίο μετράει τη μείωση της φασαρίας σε κάθε διαχωρισμό. Το χαρακτηριστικό με το μεγαλύτερο πληροφοριακό κέρδος επιλέγεται ως το χαρακτηριστικό για το διαχωρισμό των δεδομένων.

Τα δέντρα απόφασης είναι εύκολα, ερμηνεύσιμα και μπορούν να διαχειριστούν κατηγορηματικές και αριθμητικές τιμές. Μπορούν επίσης να διαχειριστούν κενά δεδομένα και οριακές τιμές. Είναι επίσης εύκολα στην οπτικοποίηση και στην ερμηνεία, γεγονός που τα καθιστά χρήσιμα στην εξήγηση πολύπλοκων μοντέλων. Παρόλα αυτά, μπορεί να είναι επιρρεπείς στην υπερπροσαρμογή

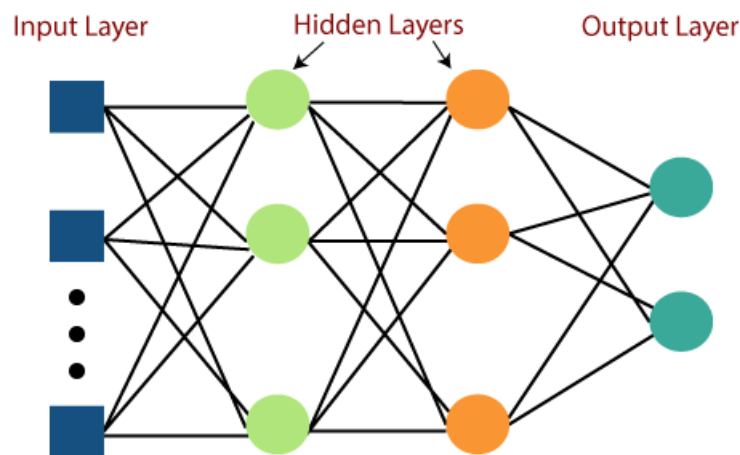


(overfitting) και μπορεί να γίνουν πολύπλοκα όταν έχουν να κάνουν με μεγάλα και πολύπλοκα σύνολα δεδομένων. Για να αποφευχθεί η υπερπροσαρμογή (overfitting), χρησιμοποιούνται τεχνικές όπως το κλάδεμα (pruning) θέτοντας ένα μέγιστο όριο και προσθέτοντας έναν όρο κανονικοποίησης.

- Multi-layer Perceptron

Ο MLP χρησιμοποιείται σε εφαρμογές κατηγοριοποίησης και παλινδρόμησης. Είναι ένα είδος τεχνητών νευρωνικών δικτύων που έχει πολλαπλά επίπεδα, με τουλάχιστον ένα επίπεδο εισόδου, ένα ή περισσότερα κρυφά επίπεδα και ένα επίπεδο εξόδου, Σχήμα 2.10. Τα επίπεδα είναι πλήρως συνδεδεμένα και ο αλγόριθμος εκπαιδεύεται με τη χρήση οπισθοδιάδοσης (backpropagation).

Σχήμα 2.10: Multi-layer Perceptron[15]



Είναι ικανός στο να μαθαίνει μη-γραμμικές σχέσεις μεταξύ εισόδου και εξόδου. Κάθε κόμβος στα κρυμμένα επίπεδα εφαρμόζει μια μη γραμμική συνάρτηση στις εισόδους του, η οποία επιτρέπει το δίκτυο να μαθαίνει σύνθετες αναπαραστάσεις των δεδομένων. Ο αριθμός των κρυφών επιπέδων και ο αριθμός των κόμβων σε κάθε επίπεδο είναι υπερπαραμέτροι που μπορούν να ρυθμιστούν για να βελτιωθεί η απόδοση του μοντέλου [29].

Ο MLP είναι ένας ισχυρός αλγόριθμος που μπορεί να επιτύχει κορυφαίες επιδόσεις σε ένα ευρύ φάσμα προβλημάτων. Μπορεί να διαχειριστεί μεγάλα και υψηλών διαστάσεων σύνολα δεδομένων καθώς και να μάθει μη γραμμικές σχέσεις μεταξύ εισόδου και εξόδου. Ωστόσο, ο MLP μπορεί να είναι ευαίσθητος στην επιλογή των υπερπαραμέτρων και μπορεί να είναι υπολογιστικά κοστο-

---

βόρος στην εκπαίδευση. Επίσης, μπορεί να είναι ευαίσθητος στην κλιμάκωση των χαρακτηριστικών εισόδου και να χρειαστεί ένα μεγάλο αριθμό δεδομένων για να εκπαιδευτεί. Για την αντιμετώπιση αυτών των προβλημάτων χρησιμοποιούνται τεχνικές κανονικοποίησης όπως η εγκατάλειψη (dropout), πρόωρη διακοπή (early stopping) και κανονικοποίηση L1, L2.

## 2.3 Μετρικές Αξιολόγησης Μοντέλων και Αποτελεσμάτων

Πριν αναλυθούν οι μετρικές, παρατίθεται η ερμηνεία όρων που θα χρησιμοποιηθούν:

1. Αληθές Θετικό – True Positive (TP): αναφέρεται σε περιπτώσεις όπου η πραγματική κλάση της εισόδου και η προβλεπόμενη είναι θετικές. Δηλαδή το μοντέλο σωστά αναγνώρισε τη θετική κλάση.
2. Αληθές Αρνητικό – True Negative (TN): στη περίπτωση αυτή η πραγματική κλάση της εισόδου και η προβλεπόμενη είναι αρνητικές. Με άλλα λόγια, το μοντέλο αναγνώρισε σωστά την αρνητική κλάση.
3. Ψευδές Αρνητικό – False Negative (FN): έχει να κάνει με περιπτώσεις όπου η πραγματική κλάση της εισόδου είναι θετική αλλά η προβλεπόμενη κλάση είναι αρνητική. Το μοντέλο δεν αναγνώρισε σωστά τη θετική κλάση.
4. Ψευδές Θετικό – False Positive (FP): στην περίπτωση αυτή η πραγματική κλάση εισόδου είναι αρνητική αλλά η προβλεπόμενη κλάση είναι θετική. Το μοντέλο δηλαδή δεν αναγνώρισε σωστά την αρνητική είσοδο.

Για να αξιολογηθεί η απόδοση των μοντέλων μηχανικής μάθησης, χρησιμοποιούνται διάφορες μετρικές για την αξιολόγηση της ποιότητας των προβλέψεων του μοντέλου. Στην υποενότητα αυτή θα συζητηθούν ορισμένες ευρέως χρησιμοποιούμενες μετρικές στη μηχανική μάθηση, και θα διερευνηθούν τα πλεονεκτήματα και οι αδυναμίες τους [27].

- Ακρίβεια (Accuracy): είναι η πιο κοινή μετρική που χρησιμοποιείται για την αξιολόγηση της απόδοσης του μοντέλου. Μετράει την αναλογία των σωστών προβλέψεων προς το συνολικό αριθμό περιπτώσεων του συνόλου δεδομένων.

---

Η μετρική αυτή είναι καταλληλότερη για ισορροπημένα σύνολα δεδομένων.

- Τύπος:  $Accuracy = \frac{correct\ predictions}{total\ predictions}$
- Πλεονεκτήματα:
  - \* Είναι απλή και εύκολη στην κατανόηση
  - \* Παρέχει μια πλήρη εικόνα της απόδοσης του μοντέλου
- Μειονεκτήματα
  - \* Δεν είναι κατάλληλη μετρική για ανισόρροπα σύνολα δεδομένων καθώς μπορεί να δώσει μια παραπλανητική αξιολόγηση της απόδοσης
- Ανάκληση (Recall): γνωστή και ως ευαισθησία ή αλλιώς ποσοστό των αληθώς θετικών, είναι η μετρική της ικανότητας του μοντέλου να αναγνωρίζει σωστά όλες τις θετικές περιπτώσεις.
  - Τύπος:  $Recall = \frac{true\ positives}{true\ positives + false\ negatives}$
  - Πλεονεκτήματα:
    - \* Είναι μια χρήσιμη τεχνική σε περιπτώσεις όπου το κόστος των ψευδή αρνητικών είναι υψηλό
    - \* Βοηθάει στον προσδιορισμό της ικανότητας του μοντέλου να διακρίνει τα αληθή θετικά
  - Μειονεκτήματα:
    - \* Δε λαμβάνει υπόψιν τα ψευδή θετικά
- Ακρίβεια θετικών προβλέψεων (Precision): είναι η μετρικής της ικανότητας του μοντέλου να αναγνωρίζει σωστά τις θετικές περιπτώσεις
  - Τύπος:  $Precision = \frac{true\ positives}{true\ positives + false\ positives}$
  - Πλεονεκτήματα:
    - \* Είναι μια χρήσιμη τεχνική σε περιπτώσεις όπου το κόστος των ψευδή θετικών είναι υψηλό
    - \* Βοηθάει στον προσδιορισμό της ικανότητας του μοντέλου να διακρίνει τα αληθή θετικά

- 
- Μειονεκτήματα
    - \* Δε λαμβάνει υπόψιν του τα ψευδή αρνητικά
  - Σκορ F1 (F1 score): είναι ένα μέτρο της ακρίβειας του μοντέλου που λαμβάνει υπόψιν του την ακρίβεια θετικών προβλέψεων και την ανάκληση. Είναι ο αρμονικός μέσος του precision και του recall και παρέχει μια ενιαία βαθμολογία που αντιπροσωπεύει την ισορροπία μεταξύ των δύο μετρικών.
    - Τύπος:  $F1score = 2 * \frac{Precision * Recall}{Precision + Recall}$
    - Πλεονεκτήματα:
      - \* Είναι μια χρήσιμη τεχνική για ανισόροπα σύνολα δεδομένων, καθώς λαμβάνει υπόψιν το precision και το recall
    - Μειονεκτήματα:
      - \* Δεν παρέχει μια ξεκάθαρη ερμηνεία της απόδοσης του μοντέλου
  - Περιοχή κάτω από την καμπύλη ROC (ROC AUC): η μετρική αυτή χρησιμοποιείται για να αξιολογήσει την απόδοση των μοντέλων δυαδικής κατηγοριοποίησης. Μετρά την περιοχή κάτω από την καμπύλη ROC, η οποία απεικονίζει το ποσοστό των αληθών θετικών (TPR) έναντι του ποσοστού των ψευδών θετικών (FPR) σε διάφορα όρια.
    - Πλεονεκτήματα:
      - \* Είναι ανθεκτικό σε ανισόροπα σύνολα δεδομένων
    - Μειονεκτήματα:
      - \* Δεν είναι κατάλληλο για προβλήματα κατηγοριοποίησης πολλών κατηγοριών

Σε μερικές μετρικές χρησιμοποιήθηκε η σταθμισμένη μορφή τους η οποία λαμβάνει υπόψιν τη σημασία ή τη συσχέτιση των διάφορων κλάσεων ή δειγμάτων στο σύνολο δεδομένων. Αποδίδει ένα βάρος σε κάθε κλάση ή δείγμα βασισμένο σε κάποια προκαθορισμένα κριτήρια, όπως τη συχνότητα τους ή τη σημαντικότητα τους στον τομέα. Στη συνέχεια, η μετρική υπολογίζεται δίνοντας βάση στα δείγματα ή στη κλάση που θεωρείται πιο σημαντική ή δύσκολη να προβλεφθεί.

---

Σε γενικές γραμμές, οι σταθμισμένες μετρικές μπορεί να είναι χρήσιμες όταν πρόκειται για ανισόρροπα σύνολα δεδομένων. Ωστόσο, είναι σημαντικό να επιλέγονται τα βάρη προσεκτικά, με βάση τη γνώση του τομέα και εμπειρικών στοιχείων προκειμένου να αποφευχθεί η εισαγωγή διάφορων τάσεων ή υπερβολική προσαρμογή στα δεδομένα εκπαίδευσης.

## 2.4 Μετρικές χρόνου εκτέλεσης

Εκτός από την ακρίβεια και το σφάλμα παρατίθενται και κάποιες άλλες μετρικές κατηγοριοποίησης έτσι ώστε να υπάρξει μια πλήρης εικόνα για την επιλογή του μοντέλου [21].

- Χρόνος εκπαίδευσης (fit time): ο χρόνος που απαιτείται έτσι ώστε να εκπαιδευτεί το μοντέλο σε ένα σύνολο δεδομένων. Αυτό συμπεριλαμβάνει το χρόνο που απαιτείται για την επεξεργασία των δεδομένων, την προσαρμογή των παραμέτρων του μοντέλου στα δεδομένα καθώς και οποιαδήποτε προ-επεξεργασία και εξαγωγή των χαρακτηριστικών που απαιτείται.
- Χρόνος αξιολόγησης (score time): έχει να κάνει με το χρόνο που απαιτείται για να γίνουν οι προβλέψεις με τη χρήση του εκπαιδευμένου μοντέλου σε ένα νέο σύνολο δεδομένων. Αυτό συμπεριλαμβάνει το χρόνο που απαιτείται για την επεξεργασία των δεδομένων και την πρόβλεψη καθώς και τη μετ-επεξεργασία ή αξιολόγηση που απαιτείται.

Οι δύο χρονικές μετρήσεις είναι σημαντικοί παράγοντες που πρέπει να ληφθούν υπόψιν όταν επιλέγεται ένα μοντέλο ή όταν ρυθμίζονται οι παράμετροι του. Ένα μοντέλο που απαιτεί πολύ χρόνο για την προσαρμογή μπορεί να μην είναι κατάλληλο για εφαρμογές πραγματικού χρόνου και ένα μοντέλο που χρειάζεται μεγάλο χρόνο αποτελεσμάτων μπορεί να μην είναι κατάλληλο για εφαρμογές μεγάλης κλίμακας. Οι χρόνοι αυτοί μπορεί να επηρεαστούν από διάφορους παράγοντες, όπως η πολυπλοκότητα του μοντέλου, το μέγεθος του συνόλου δεδομένων και το υλικό που χρησιμοποιήθηκε.

# Κεφάλαιο 3

## Βιβλιογραφική Ανασκόπηση

### 3.1 Παραδείγματα εφαρμογής αλγορίθμων μηχανικής μάθησης

Ο σκοπός αυτού του κεφαλαίου είναι να παρουσιάσει μια ανασκόπηση της βιβλιογραφίας των μεθόδων μηχανικής μάθησης που εφαρμόζονται στην προβλεπτική μηχανική. Η έρευνα αυτή πραγματοποιήθηκε στα αρχικά στάδια της διπλωματικής εργασίας με σκοπό την καλύτερη κατανόηση του θεωρητικού αλλά και του πρακτικού υπόβαθρου.

#### 3.1.1 Πρόβλεψη περιόδου βλάβης με χρήση τεχνητών νευρωνικών δικτύων

Οι Sampraiο et al. [46] προτείνουν μια μεθοδολογία για την επεξεργασία και το μετασχηματισμό των συνθετικών δεδομένων που συλλέγονται από ένα σύστημα δονήσεων το οποίο προσομοιώνει έναν κινητήρα. Στη συνέχεια δημιουργείται μια βάση δεδομένων που εκπαιδεύει και ελέγχει ένα Τεχνητό Νευρωνικό Δίκτυο το οποίο θα παρέχει τη δυνατότητα της μελλοντικής πρόβλεψης των καταστάσεων του εξαρτήματος, προσδιορίζοντας τη χρονική περίοδο που θα συμβεί κάποια βλάβη. Για την επίτευξη αυτού του σκοπού, κατασκευάστηκε ένα μοντέλο συσκευής για να προσομοιώσει αντιπροσωπευτικές δονήσεις κινητήρα, αποτελούμενο από έναν ανεμιστήρα φύξης υπολογιστή και μερικούς μαγνήτες. Οι μετρήσεις πραγματοποιήθηκαν χρησιμοποιώντας ένα επιταχυνσιόμετρο και τα δεδομένα συλλέχθηκαν και υποβλήθηκαν σε επεξεργασία για την παραγωγή ενός δομημένου συνόλου δεδομένων. Η εκπαίδευση των νευρωνικών δικτύων με το σύνολο δεδομένων συγκλίνει γρήγορα και σταθερά ενώ οι έλεγχοι που πραγματοποιήθηκαν ήταν η k-πλή σταυρωτή αξιολόγηση και η γενίκευση του μοντέλου. Οι έλεγχοι χρησιμοποίησαν ως δείκτη απόδοσης

τη ρίζα της μέσης τετραγωνικής απόκλισης (Root-mean-square deviation) RMSE και παρουσίασαν εξαιρετικές επιδόσεις. Οι ίδιοι έλεγχοι πραγματοποιήθηκαν και με άλλες τεχνικές μηχανικής μάθησης όπως τα Regression Tree, Random Forest, Support Vector Machine, για να αποδειχθεί η αποτελεσματικότητα των νευρωνικών δικτύων κυρίως στη γενικευσιμότητα. Μέσω του υπολογισμού του RMSE, κάνοντας χρήση των εκτιμώμενων και προβλεπόμενων τιμών, ήταν δυνατόν να υπολογιστεί η ακρίβεια του μοντέλου για το πόσο αντανακλά την πραγματικότητα του συστήματος που μελετάτε. Σε αυτή τη μελέτη η τιμή του  $k$  ορίστηκε ως 5.

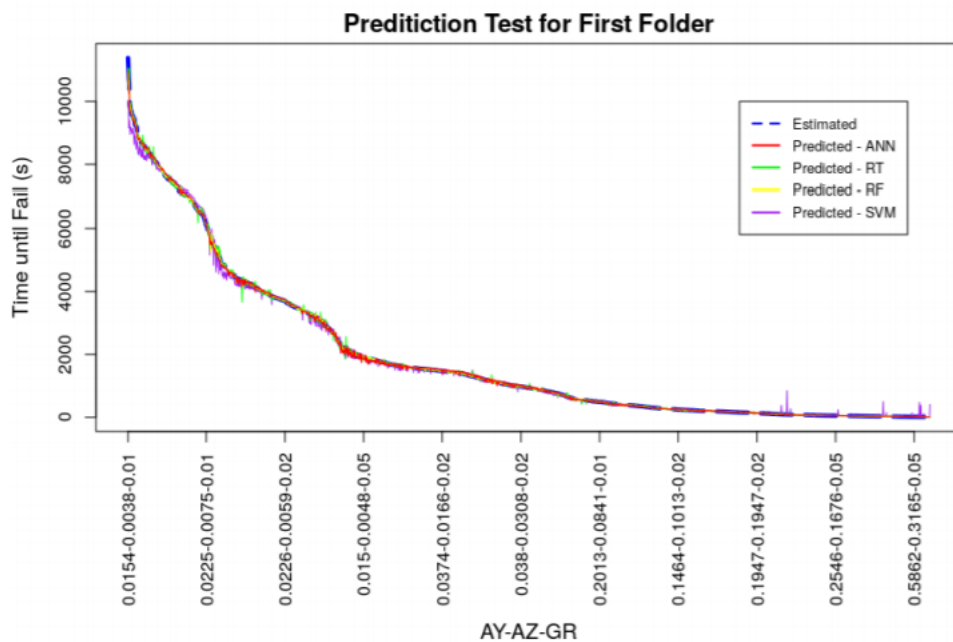
Αρχείο	ANN	RT	RF	SVM
1	0.0039	0.0047	0.0025	0.0106
2	0.0035	0.0054	0.0035	0.0129
3	0.0028	0.0051	0.0022	0.0105
4	0.0041	0.0052	0.0024	0.0120
5	0.0049	0.0052	0.0026	0.0123
Μέσος Όρος	0.0038	0.0051	0.0026	0.0117

Πίνακας 3.1: RMSE τιμές των υποσυνόλων που χρησιμοποιήθηκαν στον έλεγχο των τεχνικών

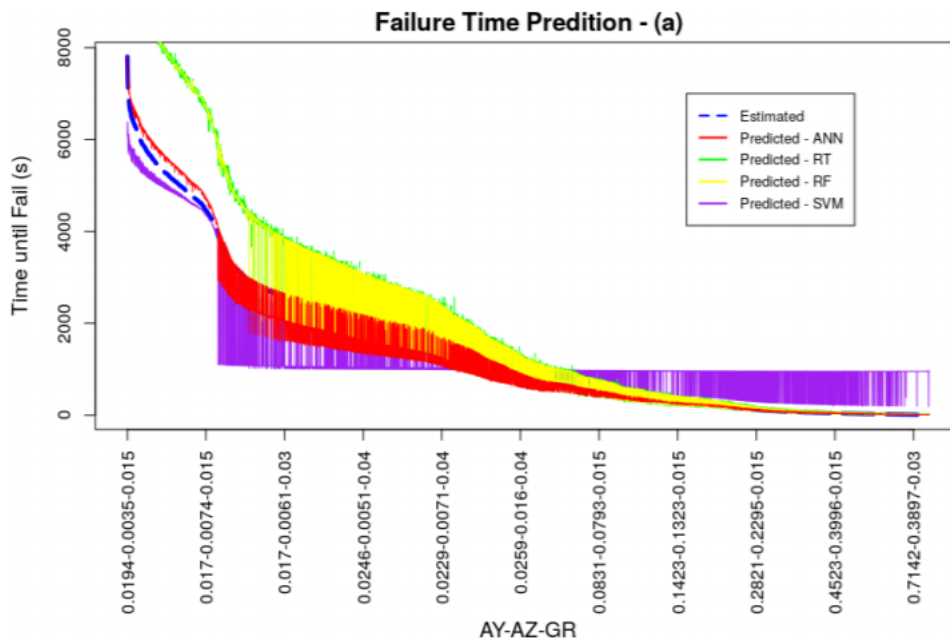
Όλα τα μοντέλα είχαν κάποια σύγκλιση στη φάση της εκπαίδευσης και ήταν σε θέση να αποδώσουν αρκετά καλά. Ο μέσος όρος RMSE, όπως φαίνεται στον Πίνακα 3.1, για το μοντέλο ANN ήταν 0.0038 και ως εκ τούτου βρίσκεται σε καλύτερη θέση από τα μοντέλα που εκπαιδεύτηκαν με τεχνικές RF και SVM ενώ η τεχνική RT έδωσε την καλύτερη απόδοση της  $k$ -πλή σταυρωτής αξιολόγησης. Στα σύνολα δεδομένων που χρησιμοποιήθηκαν για την εκπαίδευση και τον έλεγχο των μοντέλων διπλασιάστηκε ο ρυθμός αύξησης του πλάτους και των σημάτων δόνησης έτσι ώστε να προσομοιωθούν διαφορετικά σενάρια δονητικής συμπεριφοράς δημιουργώντας βέλτιστες δοκιμές. Έτσι, δημιουργήθηκαν δύο νέα σύνολα σε σχέση με το αρχικό σύνολο δεδομένων εκπαίδευσης προκειμένου να επικυρωθεί η γενικευσιμότητα του μοντέλου. Το γενικευμένο σύνολο δεδομένων (α) δημιουργήθηκε με το ίδιο πλάτος σημάτων δόνησης αλλά με διαφορετικό ρυθμό αύξησης της τάξης του 0.015, 0.03 και 0.04 οδηγώντας σε ένα σύνολο 9.180 παρατηρήσεων. Το γενικευμένο σύνολο δεδομένων (β) δημιουργήθηκε από τις μέσες τιμές πλάτους της περιστροφικής ταχύτητας του κινητήρα ψύξης, διπλασιασμένο με τους ρυθμούς αύξησης των σημάτων δόνησης οδηγώντας σε 153 παρατηρήσεις. Ακολουθούν συγκριτικά γραφήματα στα Σχήματα 3.1, 3.2 και 3.3 μεταξύ των εκτιμώμενων και των προβλεπόμενων τιμών

του αρχικού και των δύο συνόλων δεδομένων γενίκευσης, αντίστοιχα.

Σχήμα 3.1: Έλεγχος αρχικού συνόλου [46]



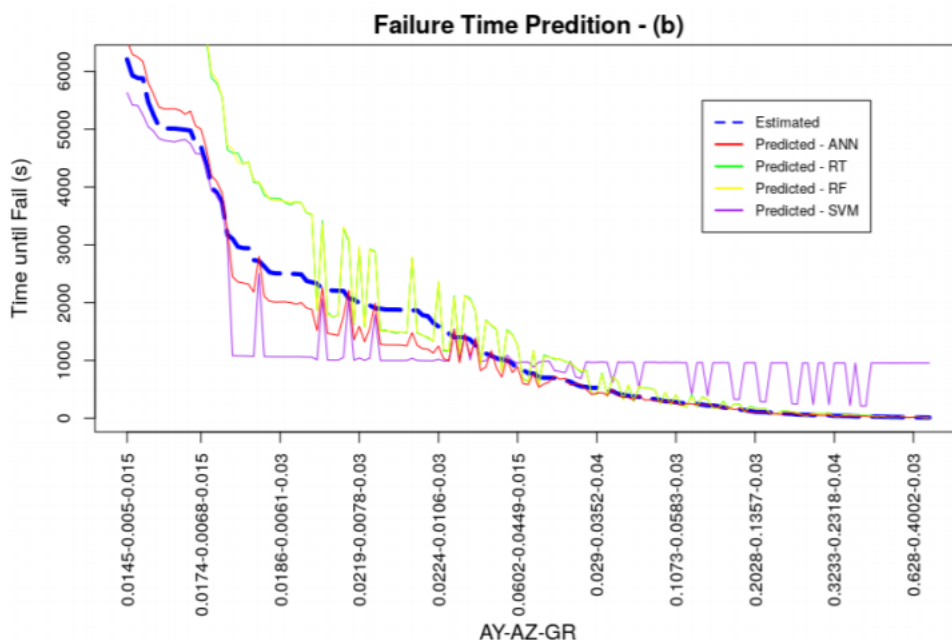
Σχήμα 3.2: Έλεγχος γενίκευσης ( $\alpha$ ) [46]



Οι τιμές των συνόλων δεδομένων ( $\alpha$ ) και ( $\beta$ ) δεν ξεπερνούν τις μέγιστες και ελάχιστες τιμές του συνόλου δεδομένων των εκπαιδευόμενων μοντέλων έτσι ώστε να μην επηρεαστεί η απόδοση των τεχνικών που βασίζονται σε δέντρα. Από τα Σχήματα 3.2, 3.3 και τον Πίνακα 3.2 μπορεί να παρατηρηθεί ότι η πρόβλεψη ακολούθησε την ίδια συμπεριφορά με εκείνη που παρουσιάστηκε στη k-πλή σταυρωτή αξιολόγηση



Σχήμα 3.3: Έλεγχος γενίκευσης ( $\beta$ ) [46]



αλλά με πιο έντονο τρόπο. Οι βραχυπρόθεσμες προβλέψεις ήταν ακριβείς για τις τεχνικές ANN, RT και RF ωστόσο, για μεσοπρόθεσμες και μακροπρόθεσμες προβλέψεις η τεχνική ANN απέδωσε καλύτερα από τις υπόλοιπες. Γενικά ο ANN είχε το καλύτερο αποτέλεσμα δείκτη RMSE για τους ελέγχους γενίκευσης με τιμή 0.0313 για το σύνολο (α) και 0.1184 για το σύνολο (β). Οι τιμές που παρουσιάζονται είναι χαμηλές αποδεικνύοντας την καλή γενικευσιμότητα του μοντέλου. Οι τεχνικές RT και RF είχαν καλή γενικευσιμότητα για μεσοπρόθεσμες προβλέψεις και ακόμη και με χαμηλούς δείκτες RMSE ο SVM δεν είχε πολύ καλή απόδοση.

Γενίκευση	ANN	RT	RF	SVM
a	0.0313	0.0922	0.0920	0.0696
b	0.1184	0.1417	0.1430	0.1237

Πίνακας 3.2: Τιμές RMSE για τον έλεγχο γενίκευσης των τεχνικών μηχανικής μάθησης

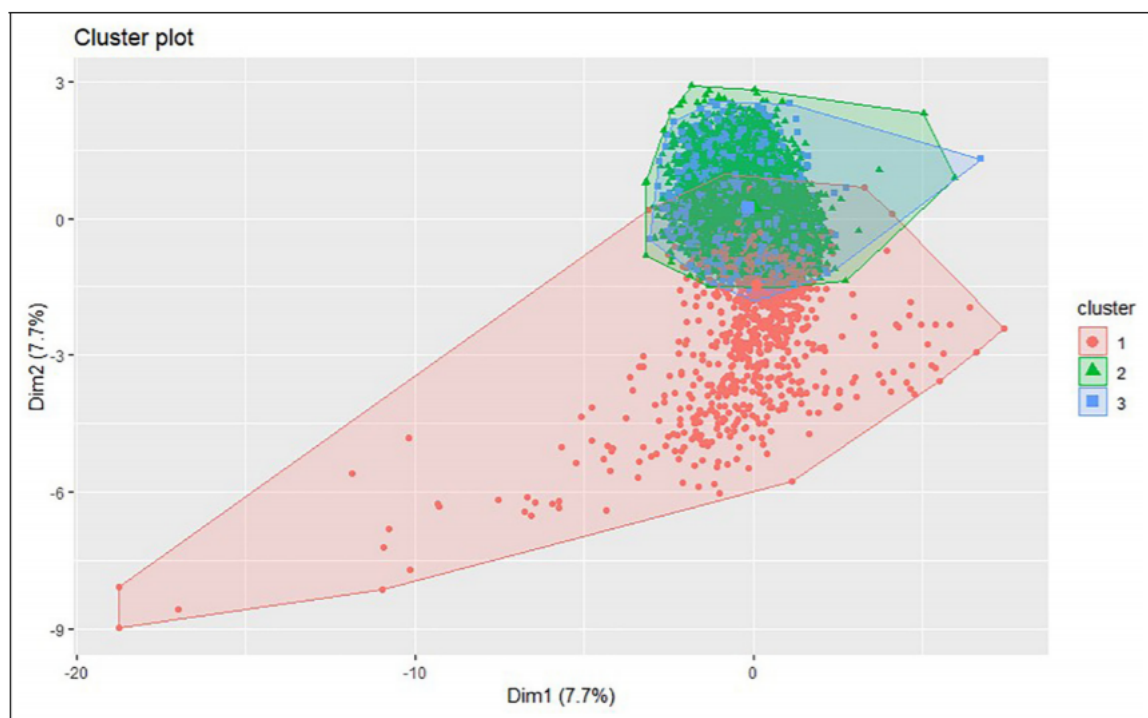
### 3.1.2 Μια έξυπνη προσέγγιση για προεπεξεργασία και ανάλυση δεδομένων σε μια περίπτωση βιομηχανικής μελέτης

Οι Bekar et al. [19] εξετάζουν κυρίως την πρόβλεψη αποτυχιών καθώς και τις βασικές αιτίες δύο μηχανών συμφόρησης σε μια από τις κορυφαίες κατασκευαστικές εταιρίες που βρίσκονται στην Σουηδία. Χρησιμοποιούνται δεδομένα υψηλής διάστασης τα οποία συλλέγονται μέσω εξελιγμένων αισθητήριων συστημάτων, PLC μηχανών

νημάτων, συστημάτων παρακολούθησης παραγωγής και συστημάτων συντήρησης. Τα δεδομένα έχουν να κάνουν με την ισχύ, τη ροπή, τη δόνηση και τη θερμοκρασία. Συνολικά υπάρχουν 13 σύνολα δεδομένων για το πρώτο μηχάνημα και 9 σύνολα δεδομένων για το δεύτερο. Αρχικά, τα δεδομένα υπόκεινται σε κλιμάκωση έτσι ώστε να προετοιμαστούν για τη διεργασία της κατηγοριοποίησης. Έπειτα, κανονικοποιούνται με τη μέθοδο z-score. Σε αυτό το σημείο εφαρμόζεται ο PCA αλγόριθμος στα δεδομένα των μηχανών και ύστερα αυτά εισάγονται στον αλγόριθμο ομαδοποίησης K-means. Με τη μέθοδο του αγκώνα υπολογίστηκε ότι το k θα είναι ίσο με 3 και τα αποτελέσματα του K-means φαίνονται στα γραφήματα των Σχημάτων 3.4 και 3.5 για την πρώτη και τη δεύτερη μηχανή, αντίστοιχα.

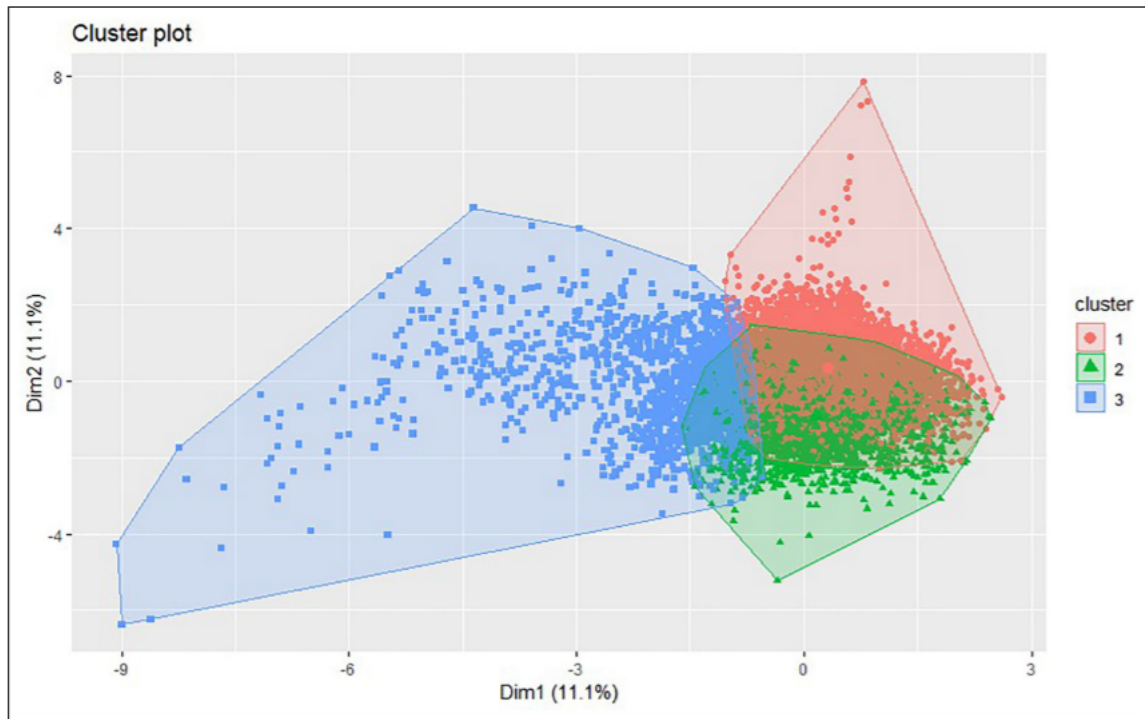
Τα αποτελέσματα της ομαδοποίησης για την πρώτη μηχανή υποδεικνύουν ότι υπάρχουν κανονικά δείγματα μέσα στο σύνολο δεδομένων τα οποία εντοπίζονται στις χαμηλές συχνότητες και τα οποία μπορεί να υποδεικνύουν κάποια βλάβη στην κατάσταση της μηχανής.

Σχήμα 3.4: Αποτελέσματα K-means ομαδοποίησης για την πρώτη μηχανή [19]



Μπορεί επίσης να παρατηρηθεί ότι υπάρχουν αρκετά κανονικά δείγματα που πρέπει να αναλυθούν και για το δεύτερο μηχάνημα. Είναι περιπτώσεις όπου η συχνότητα με την οποία εμφανίζονται είναι μικρή δεδομένου του συνολικού συνόλου

Σχήμα 3.5: Αποτελέσματα K-means ομαδοποίησης για τη δεύτερη μηχανή [19]



όπως φαίνεται στον Πίνακα 3.3 και μπορεί και αυτά να υποδεικνύουν κάποιο πρόβλημα στο μηχάνημα.

Μηχάνημα	Συστάδα	Περιπτώσεις
Machine_1	1	248
	2	7488
	3	5414
Machine_2	1	2888
	2	8343
	3	1743

Πίνακας 3.3: Αριθμός περιπτώσεων της κάθε συστάδας για το κάθε μηχάνημα

### 3.1.3 Αρχιτεκτονική προβλεπτικής συντήρησης με μηχανική μάθηση σε πυρηνικές υποδομές

Οι Gohel et al. [28] είχαν στόχο το σχεδιασμό και την ανάλυση ενός προηγμένου συστήματος ανάλυσης προβλεπτικής συντήρησης σε συνδυασμό με αλγορίθμους μηχανικής μάθησης. Το σύστημα αυτό θα χρησιμοποιηθεί για να προβλέπει βλάβες σε πυρηνικούς σταθμούς με κύριο μέλημα την προστασία του περιβάλλοντος. Θα συλλεχθούν και θα αναλυθούν ετερογενή δεδομένα από διάφορους αισθητήρες θερ-

μότητας, δόνησης, πίεσης ,επιταχυνσιόμετρων καθώς και άλλων υποσυστημάτων του σταθμού. Έπειτα, θα εφαρμοστούν οι αλγόριθμοι μηχανικής μάθησης SVM και Logistic Regression. Τα δεδομένα της εργασίας είναι συνθετικά καθώς βρίσκονται ακόμα στο στάδιο συλλογής πραγματικών δεδομένων. Στον Πίνακα 3.4 που ακολουθεί παρουσιάζεται μια σύγκριση των αποτελεσμάτων της συγκεκριμένης μελέτης με άλλες που σχετίζονται με πυρηνικούς σταθμούς και όπως παρατηρείται έχει την καλύτερη επίδοση.

Υπάρχουσες Εργασίες	Τύποι Χαρακτηριστικών	Αλγόριθμοι Μηχανικής Μάθησης	Ακρίβεια
24	PCA	Linear	80%
25	RFE	LR, RF, MLP	67-75%
26	Multiple	SVN,ANN	Vary
Proposed Work	No Variance & Co-related	SVM, LR	95%

Πίνακας 3.4: Σύγκριση αποτελεσμάτων

Η τεχνικές μηχανικής μάθησης SVM και Linear Regression χρησιμοποιούνται για να διερευνήσουν και να συγκρίνουν σπάνια γεγονότα που θα μπορούσαν να συμβούν σε πυρηνικούς σταθμούς. Τέλος, ο SVM παρέχει καλύτερες μετρήσεις απόδοσης.

### **3.1.4 SOPHIA: Μια αρχιτεκτονική IoT βασισμένη σε δεδομένα και μηχανική μάθηση για την προβλεπτική συντήρηση στη βιομηχανία 4.0**

Οι Calabrese et al. [22] εφάρμοσαν μια μέθοδο που βασίζεται στα δεδομένα και τη μηχανική μάθηση σε μηχανήματα βιομηχανιών ξυλουργίας. Πιο συγκεκριμένα, παρουσιάζεται μια εφαρμογή μηχανικής μάθησης για PdM στο ρουλεμάν της μηχανής Electros spindle(ES). Οι προβλεπόμενες βλάβες υπολογίζονται μέσω μοντέλων κατηγοριοποίησης δέντρων (Gradient Boosting, Random Forest και Extreme Gradient Boosting). Προτείνεται μια PdM μεθοδολογία η οποία εκμεταλλεύεται την πληροφορία των μεγάλων δεδομένων που παρέχονται από το σύστημα καταγραφής δεδομένων, χωρίς να χρειαστεί να τοποθετηθούν αισθητήρες στο μηχανήμα έτσι ώστε να συλλέξει δεδομένα κατάστασης. Εξετάστηκε κυρίως η περίπτωση του σπασμένου ρουλεμάν το οποίο είναι υπεύθυνο για την υπολειτουργία του κινητήρα στα Rover μηχανήματα ξυλουργικής. Τα αρχεία καταγραφής συλλέχθηκαν σε μια χρονική περίοδο πέντε μηνών και αναφέρονταν σε 14 μηχανήματα Rover, 5 ES με σπασμένο ρουλεμάν και τα υπόλοιπα 9 ES χωρίς ρουλεμάν. Κατά τη διάρκεια

της αξιολόγησης πραγματοποιήθηκαν δοκιμές και με άλλες μεθοδολογίες μηχανικής μάθησης (SVM με γραμμικό και Γκαουσιανό πυρήνα, Nearest-Neighbor (NN) και Decision Tree(DT)). Σύμφωνα με τα αποτελέσματα, το Extreme Gradient Boosting (XGBoost), το Distributed Random Forest (DRF) και το Gradient Boosting Machine (GBM) απέδωσαν καλύτερα στην ακρίβεια, την ανάκληση και στην ακρίβεια των θετικών προβλέψεων από ότι ο SVM και ο DT. Για να ληφθούν αμερόληπτες εκτιμήσεις το σύνολο δεδομένων χωρίζεται σε δύο υποσύνολα, το σύνολο εκπαίδευσης και το σύνολο ελέγχου με ποσοστό 70% και 30% αντίστοιχα.

Μοντέλο	RUL	Ακρίβεια		Ανάκληση		Ακρίβεια Θετικών Προβλέψεων	
		Εκπαίδευση	Αξιολόγηση	Εκπαίδευση	Αξιολόγηση	Εκπαίδευση	Αξιολόγηση
Distributed Random Forest(DRF)	30	97.4(0.4)	96.8	98.8(0.4)	96.5	98.3(0.3)	100
	20	97.8(0.2)	96.5	99.3(0.2)	96.5	98.4(0.4)	99.3
	10	98.2(0.5)	96.5	99.6(0.2)	96.9	98.5(0.4)	99.3
Gradient Boosting Machine(GBM)	30	98.2(0.2)	98.9	99.6(0.1)	99.6	98.3(0.3)	99.1
	20	98.7(0.4)	97.8	99.5(0.2)	98.7	99.1(0.2)	98.8
	10	98.5(0.8)	97.8	93.3(0.2)	97.9	99.2(0.6)	99.8
Extreme Gradient Boosting (XGB)	30	98.1(0.5)	98.8	99.5(0.1)	100	98.4(0.7)	98.7
	20	98.3(0.7)	96.3	99.7(0.01)	100	98.4(0.7)	96.2
	10	98.4(0.6)	98.8	99.6(0.1)	99.8	98.8(0.7)	99.8

Πίνακας 3.5: Μετρικές ορθότητας, μνήμης και ακρίβεια των GBM, DRF και XGBoost

Η μεθοδολογία αυτή επέτρεψε τον ταυτόχρονο έλεγχο πολλαπλών συνδεδεμένων μηχανών, καθιστώντας δυνατή την καθημερινή παρακολούθηση, που μπορεί να υιοθετηθεί σαν τεχνική μαζί με τη διαχείριση συντήρησης. Η αποτελεσματικότητα της μεθόδου αποδεικνύεται ελέγχοντας ένα ανεξάρτητο δείγμα ξυλουργικών μηχανών. Ο GBM ομαδοποίησε σωστά 78 από τις 81 καταγραφές λαμβάνοντας υπόψη την περίοδο διακοπής λειτουργίας της μηχανής κατά τη χρονική περίοδο των 30 ημερών και 605 από τις 610 σε περίοδο κανονικής λειτουργίας αποδίδοντας 98.2% στην ακρίβεια, 98.6% στην ανάκληση και 98.3% στην ακρίβεια θετικών προβλέψεων όπως φαίνεται και στον Πίνακα 3.5.

### 3.1.5 Προσέγγιση μηχανικής μάθησης για προβλεπτική συντήρηση στη βιομηχανία 4.0

Οι Paolanti et al. [38] παρουσιάζουν μία νέα μεθοδολογία PdM βασισμένη σε μια προσέγγιση μηχανικής μάθησης η οποία εφαρμόζεται σε μηχανήματα κοπής ξυλουργικών εξαρτημάτων, το οποίο αποτελεί το κέντρο επεξεργασίας σε βιομηχανίες ξύλου. Τα δεδομένα έχουν συλλεχθεί από διάφορους αισθητήρες, PLC μηχανήματα

και πρωτόκολλα επικοινωνίας τα οποία θα επεξεργαστούν με εργαλεία ανάλυσης δεδομένων (Data Analysis Tool). Η προτεινόμενη PdM μεθοδολογία επιτρέπει την εφαρμογή δυναμικών κανόνων απόφασης με στόχο τη διαχείριση συντήρησης, το οποίο επιτυγχάνεται με την εκμάθηση του Random Forest στο Azure Machine Learning Studio. Το σύνολο δεδομένων περιέχει 530.731 στοιχεία που έχουν καταγραφεί σε πραγματικό χρόνο όπου το καθένα αποτελείται από 15 διαφορετικά χαρακτηριστικά. Τα αρχικά αποτελέσματα όπως φαίνονται στον Πίνακα 3.6 δείχνουν μία σωστή συμπεριφορά στην πρόβλεψη διάφορων καταστάσεων μηχανών με μια υψηλή ακρίβεια της τάξης του 95%.

Μοντέλα	Αποτελέσματα
Overall Accuracy	0.95
Average Accuracy	0.92
Micro-Averaged Precision	0.94
Macro-Averaged Precision	0.93
Micro-Averaged Recall	0.95
Macro-Averaged Recall	0.94

Πίνακας 3.6: Αποτελέσματα ομαδοποίησης

### 3.1.6 Μηχανική μάθηση για προβλεπτική συντήρηση βιομηχανικών συστημάτων με χρήση δεδομένων αισθητήρων IoT.

Οι Kanawaday et al. [33] εξέτασαν την περίπτωση ενός μηχανήματος κοπής που αποτελείται από 14 βραχίονες οι οποίοι παράγουν συσκευασίες μεταβλητού μεγέθους. Τα δεδομένα συλλέχθηκαν από τους αισθητήρες σε μια χρονική περίοδο ενός μήνα με ρυθμό δειγματοληψίας ανά δευτερόλεπτο. Έπειτα, αποθηκεύτηκαν με μορφή CSV στο σύστημα. Το αρχείο αποτελείται από 5 στήλες: Χρονική σφραγίδα, ένταση, πίεση, πλάτος και διάμετρος. Διάφορα επιβλεπόμενα μοντέλα όπως τα Νευρωτικά Δίκτυα, Μηχανές Διανυσμάτων Υποστήριξης SVM και το CART χρησιμοποιήθηκαν στο σύνολο δεδομένων για να εκπαιδεύσουν έναν κατηγοριοποιητή έτσι ώστε να εντοπίσει ποιοτικά λάθη στους κύκλους παραγωγής. Τα επιβλεπόμενα μοντέλα εκπαιδεύτηκαν σε ένα σύνολο δεδομένων το οποίο χωρίζεται σε 70% εκπαίδευση, 10% επικύρωση και 20% έλεγχο, τα αποτελέσματα των οποίων απεικονίζονται στον Πίνακα 3.7. Συμπεραίνεται ότι τα Βαθιά Νευρωτικά Δίκτυα είναι πιο αποτελεσματικά.

Μοντέλα Επιβλεπόμενης Μάθησης	Ακρίβεια Πρόβλεψης(%)
Naive Bayes	96.61
Support Vector Machine	95.52
CART	94.46
Deep Neural Network	98.69

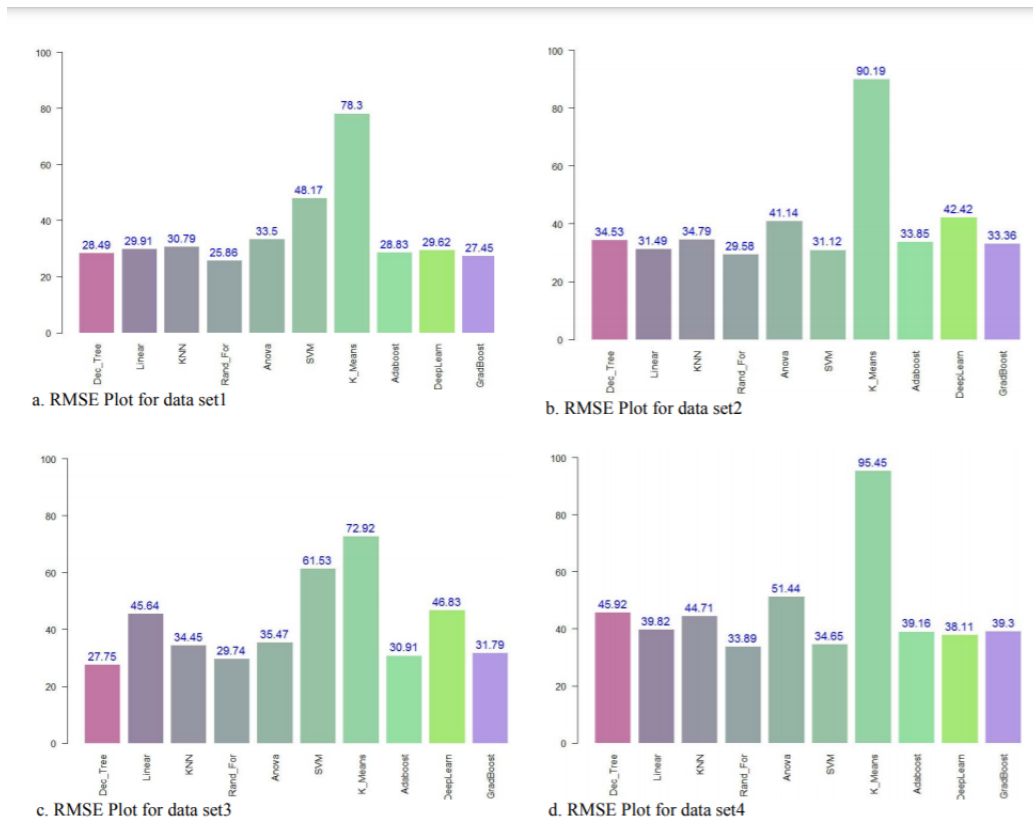
Πίνακας 3.7: Σύγκριση διαφορετικών επιβλεπόμενων αλγορίθμων

### 3.1.7 Πρόβλεψη RUL στροβιλοκινητήρων μέσω της μηχανικής μάθησης

Οι Mathew et al. [36] πραγματοποίησαν μια συγκριτική μελέτη αλγορίθμων μηχανικής μάθησης με σκοπό την πρόβλεψη του ωφέλιμου χρόνου ζωής ενός ανεμιστήρα στροβιλοκινητήρων αεροσκαφών. Το σύνολο δεδομένων προέρχεται από την αποθήκη δεδομένων της NASA. Το επιλεγμένο σύνολο δεδομένων περιλαμβάνει μετρήσεις από τους αισθητήρες φθαρμένων στροβιλοκινητήρων με ανεμιστήρα. Παρόλο που οι κινητήρες είναι του ίδιου τύπου, κάθε κινητήρας ξεκινάει με διαφορετικό βαθμό αρχικών συνθηκών και υπάρχουν παραλλαγές στη διαδικασία κατασκευής τους, οι οποίες δεν είναι γνωστές στο χρήστη. Για τους υπο εξέταση στροβιλοκινητήρες η απόδοση τους μπορεί να αλλάξει προσαρμόζοντας τρεις λειτουργικές ρυθμίσεις. Κάθε κινητήρας έχει 21 αισθητήρες συλλογής διαφορετικών μετρήσεων που σχετίζονται με την κατάσταση του κινητήρα και το χρόνο εκτέλεσης. Το κύριο χαρακτηριστικό του συνόλου δεδομένων είναι ότι πρόκειται για χρονοσειρές. Παρακάτω αναφέρονται οι αλγόριθμοι που χρησιμοποιήθηκαν για τη μελέτη αυτή: Linear Regression, Decision Tree, Support Vector Machine (SVM), Random Forest, K-Nearest Neighbors, K Means, Gradient Boosting Method, AdaBoost, Deep Learning και Anova. Η εκπαίδευση και η αξιολόγηση του μοντέλου Μηχανικής Μάθησης πραγματοποιήθηκε με τη χρήση τεσσάρων συνόλων δεδομένων όπου το κάθε ένα περιείχε δεδομένα από 100 – 250 κινητήρες. Οι δέκα αλγόριθμοι εξετάστηκαν με τη χρήση και των τεσσάρων συνόλων δεδομένων. Τα γραφήματα του Σχήματος 3.6 παρουσιάζουν τη μέση τετραγωνική ρίζα του σφάλματος και παρατηρείται ότι τα καλύτερα αποτελέσματα καταγράφηκαν από τον αλγόριθμο Random Forest ο οποίος καταγράφει την απόκλιση πολλών ταυτόχρονων μεταβλητών εισόδου και επιτρέπει υψηλό βαθμό παρατηρήσεων να λάβουν μέρος στην πρόβλεψη. Επίσης, παρατηρήθηκε ότι η απόδοση και των δέκα αλγορίθμων ήταν σύμφωνη στα τέσσερα σύνολα δεδομένων δημιουργώντας αναλογική ακρίβεια για τους διαφορετικούς αλ-

γορίθμους που δοκιμάστηκαν.

Σχήμα 3.6: Οπτικοποίηση του Root Mean Square Error των RULs έναντι Machine Ids για διαφορετικούς αλγόριθμους σε τέσσερα διαφορετικά σύνολα δεδομένων [36]



### 3.1.8 Ανάλυση βασισμένη σε αισθητήρες IoT, προεπεξεργασία Big Data και μοντέλων μηχανικής μάθησης για συστήματα παρακολούθησης πραγματικού χρόνου σε αυτοκινητοβιομηχανίες

Οι Syafrudin et al. [48] πρότειναν ένα σύστημα παρακολούθησης σε πραγματικό χρόνο που αξιοποιεί αισθητήρες βασισμένους σε IoT, επεξεργασία μεγάλων δεδομένων και ένα υβριδικό μοντέλο πρόβλεψης. Αρχικά, αναπτύχθηκαν IoT βασισμένοι αισθητήρες που συλλέγουν δεδομένα θερμοτήτας, υγρασίας, επιτάχυνσης και γυροσκοπίου. Έπειτα, το προτεινόμενο υβριδικό μοντέλο πρόβλεψης που χρησιμοποιείται είναι η Χωρική Ομαδοποίηση Εφαρμογών με Θόρυβο με Βάση την Πυκνότητα (Density-Based Spatial Clustering of Applications with Noise - DBSCAN) σε συνδυασμό με τον Random Forest που χρησιμοποιούνται για να αφαιρέσουν τις ακραίες τιμές από τα δεδομένα και να παρέχουν ανίχνευση σφαλμάτων κατά τη διαδικασία κατασκευής. Το προτεινόμενο μοντέλο αξιολογήθηκε και ελέγχθηκε



στη γραμμική συναρμολόγηση κατασκευής σε μια Κορεάτικη αυτοκινητοβιομηχανία. Στον Πίνακα 3.8 παρουσιάζονται τα αποτελέσματα σύγκρισης των επιδόσεων για τα διάφορα μοντέλα κατηγοριοποίησης.

Μοντέλο	Ακρίβεια Θετικών Προβλέψεων(%)	Ανάκληση(%)	Ακρίβεια(%)
Naive Bayes (NB)	94.1	93.6	93.567
Logistics Regression (LR)	98	98	97.953
Multilayer Perceptron(MLP)	96.8	96.8	96.784
Random Forest (RF)	98.5	98.5	98.538
DBSCAN + NB	96.8	96.7	96.74
DBSCAN + LR	98.6	98.5	98.52
DBSCAN + MLP	98.8	98.8	98.81
Hybrid Prediction Model (DBSCAN + RF)	100	100	100

Πίνακας 3.8: Σύγκριση αποδόσεων διάφορων συγκριτικών μοντέλων

Αρκετά συμβατικά μοντέλα κατηγοριοποίησης όπως το Naïve Bayes (NB), ο Logistic Regression (LR), ο Multilayer Perceptron (MLP) και ο Random Forest (RF) συγκρίθηκαν με το υβριδικό μοντέλο πρόβλεψης έτσι ώστε να αναγνωρίσει και να προβλέψει ασυνήθιστα συμβάντα. Το προτεινόμενο μοντέλο συγκέντρωσε την υψηλότερη ακρίβεια 100% σε σύγκριση με τα υπόλοιπα μοντέλα κατηγοριοποίησης. Υπήρξε μια μικρή βελτίωση στην ακρίβεια του μοντέλου μετά την ενσωμάτωση του DBSCAN για τον εντοπισμό των ακραίων τιμών. Ενσωματώνοντας το DBSCAN στο Random Forest αύξησε την ακρίβεια κατά 1.462% σε σύγκριση με το συμβατικό μοντέλο Random Forest. Επιπλέον, παρατηρήθηκε βελτίωση στην ακρίβεια και σε άλλα συμβατικά μοντέλα κατηγοριοποίησης μετά την εφαρμογή του DBSCAN της τάξης του 3.173%, 0.567% και 2.026% για το Naive Bayes, το Logistic Regression και το Multilayer Perceptron, αντίστοιχα.

## 3.2 Ερευνητικά και αναπτυξιακά έργα στην Ευρωπαϊκή Ένωση

Λαμβάνοντας υπόψη τη στρατηγική σημασία του Industry 4.0 στην οικονομία, η ΕΕ αναλαμβάνει πολυάριθμες πρωτοβουλίες με στόχο την υποστήριξη, την ανάπτυξη και την εφαρμογή αυτής της ιδέας στα κράτη μέλη της. Μεταξύ των κρατών μελών της ΕΕ, η Γερμανία, η Γαλλία και η Μεγάλη Βρετανία έχουν πιο ενεργή συμμετοχή στην ιδέα του Industry 4.0. Πιο συγκεκριμένα, παρακάτω αναφέρονται

---

κάποια ενδεικτικά ερευνητικά έργα:

- Horse [6]

Στόχος τους είναι η πραγματοποίηση ενός άλματος στην κατασκευαστική βιομηχανία μικρομεσαίων επιχειρήσεων προτείνοντας ένα νέο ευέλικτο μοντέλο έξυπνου εργοστασίου. Το μοντέλο αυτό θα περιλαμβάνει τη συνεργασία ανθρώπων, ρομπότ και μηχανημάτων, για την πραγματοποίηση βιομηχανικών εργασιών με αποτελεσματικό τρόπο, μέσω της υιοθέτησης των ρομποτικών τεχνολογιών, των συστημάτων διαχείρισης επιχειρηματικών διεργασιών και των IoT.

- Better Factory [1]

Το Better Factory παρέχει μια μεθοδολογία για το κατασκευαστικό μικρομεσαίων επιχειρήσεων έτσι ώστε να μπορούν να συνεργαστούν με σχεδιαστές έχοντας ως στόχο την ανάπτυξη νέων, εξατομικευμένων προϊόντων και τεχνολογιών καθώς και να εξελιχθούν σε πλήρως συνδεδεμένα κυβερνό-φυσικά συστήματα το οποία θα μετατραπούν σε λεπτές και ευέλικτες εγκαταστάσεις παραγωγής. 500 κατασκευαστικές SME, 320 πάροχοι τεχνολογίας και 100 σχεδιαστές θα κινητοποιηθούν μέσω δύο ανοιχτών ανταγωνιστικών προσκλήσεων με έπαθλο 11 εκατ. ευρώ δημόσιας και ιδιωτικής χρηματοδότησης μεταξύ άλλων υπηρεσιών. Στόχος του έργου είναι να καταδείξει ότι η έρευνα μπορεί να βοηθήσει το κατασκευαστικό μικρομεσαίων επιχειρήσεων και οι μεσαίες κεφαλαιοποιημένες εταιρείες να αποκτήσουν ψηφιακή αριστεία και παγκόσμια ανταγωνιστικότητα μέσω της παραγωγής lean-agile, παρέχοντας στις ευρωπαϊκές κατασκευαστικές μικρομεσαίες επιχειρήσεις πρόσβαση σε εργαλεία όπως το RAMP Internet of Things Infrastructure που θα παρέχει ανάλυση δεδομένων, εργαλεία συγκριτικής αξιολόγησης και άλλες ψηφιακές λύσεις.

- KIT4SME [10]

Το KIT4SME στοχεύει κυρίως σε ευρωπαϊκές κατασκευαστικές μικρομεσαίες επιχειρήσεις και μεσαίες κεφαλαιοποιημένες εταιρείες έχοντας ως στόχο το περιβάλλον κατασκευής. Πιο συγκεκριμένα, εισάγει την τεχνητή νοημοσύνη στα συστήματα παραγωγής. Αξιοποιώντας το δίκτυο των Κέντρων Ψηφιακής

---

Καινοτομίας (Digital Innovation Hubs - DIH), το έργο διασφαλίζει ότι τα αποτελέσματα είναι προσαρμοσμένα στο πεδίο εφαρμογής, έτοιμα για τη βιομηχανία και παραδίδονται σε μια προσαρμόσιμη ψηφιακή πλατφόρμα. Η ομαλή υιοθέτηση του προσαρμοσμένου εξοπλισμού είναι δυνατή χάρη στην υποδομή Powered by FIWARE που συνδυάζει άψογα τα εργοστασιακά συστήματα (όπως τα Σύστημα Εκτέλεσης Κατασκευής (Manufacturing Execution System - MES) και Σχεδιασμού Πόρων Επιχείρησης (Enterprise Resource Planning - ERP)), τους αισθητήρες IoT και τις φορητές συσκευές, τα ρομπότ, τα συνεργατικά ρομπότ και άλλες πηγές δεδομένων του εργοστασίου με λειτουργικές μονάδες που μπορούν να προκαλέσουν τη δημιουργία βάσης δεδομένων. Το KIT4SME προσδιορίζει τρεις κύριους άξονες:

1. Τεχνητή νοημοσύνη για έλεγχο ποιότητας που στοχεύει στον έγκαιρο εντοπισμό σφαλμάτων και στη λήψη καλύτερων αποφάσεων κατά το χρόνο εκτέλεσης
2. Τεχνητή νοημοσύνη για αναδιαμόρφωση ή εξατομίκευση προϊόντων - και στα δύο επίπεδα βελτιστοποίησης της παραγωγής και της βελτιστοποίησης του προγραμματισμού
3. Τεχνητή νοημοσύνη για αλληλεπίδραση ανθρώπου-ρομπότ. Ο τελευταίος αυτός άξονας περιλαμβάνει την παρακολούθηση της ψυχοφυσικής κατάστασης των εργαζομένων σε πραγματικό χρόνο, την ενεργοποίηση της κούρασης και των παρεμβάσεων που ανακουφίζουν από το άγχος, καθώς και τον χαρακτηρισμό και την εξέλιξη της ικανότητας του εργατικού δυναμικού.

- SME 4.0 [12]

Μια μεγάλη πρόκληση για το μέλλον έγκειται στη μεταφορά τεχνογνωσίας και τεχνολογιών του Industry 4.0 σε μικρομεσαίες επιχειρήσεις. Αυτές αποτελούν τον στυλοβάτη της οικονομίας και παίζουν σημαντικό ρόλο στα αναπτυξιακά προγράμματα της Ευρωπαϊκής Ένωσης που έχουν ως στόχο την ενίσχυση της ανταγωνιστικότητας των ευρωπαϊκών επιχειρήσεων. Παρόλες τις προοπτικές του Industry 4.0 στις μικρομεσαίες επιχειρήσεις, ο κύριος περιορισμός έγκειται στην έλλειψη ξεκάθαρων μοντέλων για την υλοποίηση και την εφαρμογή

---

τους σε αυτές. Έτσι, το έργο αυτό στοχεύει στη γεφύρωση αυτού του χάσματος μέσω της δημιουργίας ενός διεθνούς και διεπιστημονικού δικτύου έρευνας. Ο προσδιορισμός των αναγκών και των δυνατοτήτων για μια έξυπνη και ιδιοφυής μικρομεσαία επιχείρηση, η δημιουργία προσαρμοσμένων εννοιών και σχεδιαστικών λύσεων για συστήματα παραγωγής και διοικητικής μέριμνας στις μικρομεσαίες επιχειρήσεις καθώς και η ανάπτυξη κατάλληλων οργανωτικών και επιχειρηματικών μοντέλων θα είναι οι τρεις κύριοι στόχοι του δικτύου έρευνας.

- SERENA [14]

Λόγο της αυξανόμενης πολυπλοκότητας των σύγχρονων μηχανολογικών συστημάτων και των διαδικασιών καθώς και του αυξανόμενου αριθμού των αισθητήρων και συνεπώς των δεδομένων που συλλέγονται στις βιομηχανίες, ανοίγεται ένα όραμα για τη σύνδεση των διαδικασιών παραγωγής. Όλα τα δεδομένα των μηχανών θα είναι προσβάσιμα, επιτρέποντας την ευκολότερη συντήρηση τους σε περίπτωση απροσδόκητων γεγονότων. Η δημιουργία του έργου SERENA βασίζεται σε αυτές τις ανάγκες με σκοπό να πετύχει εξοικονόμηση χρόνου και χρήματος, ελαχιστοποιώντας τις δαπανηρές διακοπές λειτουργίας. Οι προτεινόμενες λύσεις καλύπτουν τις ανάγκες για ευελιξία, δυνατότητα μεταφοράς, απομακρυσμένη παρακολούθηση και έλεγχο από:

1. Μια πλατφόρμα επικοινωνίας βασισμένη στο νέφος plug-and-play για τη διαχείριση δεδομένων και την εξ' αποστάσεως επεξεργασία τους.
2. Προηγμένα IoT συστήματα και έξυπνες συσκευές για τη συλλογή δεδομένων και την παρακολούθηση της κατάστασης των μηχανημάτων.
3. Μεθόδους τεχνητής νοημοσύνης για την προγνωστική συντήρηση και τον προγραμματισμό δραστηριοτήτων συντήρησης και παραγωγής
4. Τεχνολογίες βασισμένες σε AR για την υποστήριξη του ανθρώπινου χειριστή σε δραστηριότητες συντήρησης και παρακολούθησης της κατάστασης των μηχανημάτων παραγωγής.

Το έργο SERENA θα εφαρμοστεί σε διάφορες εφαρμογές. Πιο συγκεκριμένα, θα επικεντρωθεί στην ανάπτυξη του TRL στα επίπεδα TRL5 έως TRL7.

---

Το έργο SERENA θα εφαρμοστεί σε διάφορους βιομηχανικούς τομείς (ηλεκτρικές συσκευές, μετρολογική μηχανική και παραγωγή ανελκυστήρων) και θα διερευνήσει τη δυνατότητα εφαρμογής στη βιομηχανία παραγωγής χαλύβδινων εξαρτημάτων ελέγχοντας τη σύνδεση με άλλες βιομηχανίες (αυτοκινητοβιομηχανία, αεροδιαστημική κ.λπ.) που δείχνουν τον ευέλικτο χαρακτήρα του έργου.

- DIGIMAN4.0 [3]

Το DIGIMAN4.0 ITN θα παρέχει σε παγκόσμιο επίπεδο άριστη ερευνητική κατάρτιση σε 15 ESR (Early Stage Researchers) στον τομέα των ψηφιακών τεχνολογιών κατασκευής παραγωγής του Industry 4.0 προτείνοντας:

1. Καινοτόμες τεχνολογικές λύσεις για υψηλή ποιότητα, υψηλή διεκπεραιωτική ικανότητα και παραγωγή υψηλής ακρίβειας
2. Πρωτοποριακή διεπιστημονική εκπαίδευση σε διάφορους τομείς (Τεχνολογίες Κατασκευής Ακριβείας, Ψηφιακές Τεχνολογίες Κατασκευής, Ολοκληρωμένη Μετρολογία Παραγωγής, Λιτή Παραγωγή, Διαχείριση Παραγωγής)
3. Επικύρωση διαφόρων ψηφιακών τεχνολογιών παραγωγής με την ενσωμάτωση τους σε αλυσίδες διεργασιών με στόχο την παραγωγή προηγμένων εξαρτημάτων σε διάφορους τομείς.

- ORBETEC [9]

Στόχος του προγράμματος ORBETEC, που χρηματοδοτείται από την Ευρωπαϊκή Κοινότητα μέσω του προγράμματος Marie Curie, είναι η δημιουργία ενός μοντέλου διαχείρισης ανθρώπινου δυναμικού που θα εφαρμοστεί στο Industry 4.0 μέσω της μελέτης οργανωτικών, συλλογικών, ατομικών διαδικασιών καινοτομίας σε σχέση με τη νέα λογική της αγοράς. Το καινοτόμο αυτό σχέδιο εντάσσεται σε ένα εξαιρετικά πολύπλοκο τεχνολογικό πλαίσιο, όπως αυτό του UmbraGroup που δραστηριοποιείται στον τομέα της αεροδιαστημικής, με στόχο τη δημιουργία ενός καλύτερου εργασιακού περιβάλλοντος, όπου ολόκληρη η επαγγελματική κοινότητα μπορεί να συμβάλλει συνεχώς στην καινοτομία σκιαγραφώντας νέα μοντέλα για τη διαχείριση και την ενίσχυση των

---

ανθρώπινων πόρων. Υπό αυτήν την έννοια, το ORBETEC διερευνά το ρόλο της εφαρμογής, της επαγγελματικής συμπεριφοράς και της δυναμικής των ανθρώπινων πόρων σε ένα περιβάλλον τύπου 4.0. Το ORBETEC αναλύει επίσης τις ικανότητες του Industry 4.0 περιγράφοντας ένα νέο εμπειρικά επικυρωμένο μοντέλο ικανοτήτων, τις σχετικές στρατηγικές δεξιότητες και τα σχέδια δράσης που απαιτούνται για την προώθηση και την υποστήριξη προληπτικών επαγγελματικών συμπεριφορών σε σχέση με διαδικασίες καινοτομίας σε ατομικό επίπεδο, ομάδας και οργανισμού. Πιο συγκεκριμένα, το ερευνητικό έργο επικεντρώνεται επίσης στη διαδικασία ανάλυσης και στην ανάπτυξη ομάδων για την καινοτομία και υποστήριξη του Industry 4.0. Τέλος, η ORBETEC θα βελτιώσει την απόδοση, την παραγωγικότητα και την ποιότητα ζωής μιας επιχείρησης με οικονομικό αντίκτυπο σε περιφερειακό, εθνικό και ευρωπαϊκό επίπεδο.

### 3.3 Εμπορικές δραστηριότητες

Οι βιομηχανικές εταιρείες σε όλο τον κόσμο λαμβάνουν δράση εν όψει της 4ης βιομηχανικής επανάστασης. Το κύμα των νέων τεχνολογιών ανοίγει ευκαιρίες στις εταιρείες να κάνουν βήματα προς μεγαλύτερη ευελιξία, βιωσιμότητα και παραγωγικότητα. Παρακάτω αναφέρονται μερικές μελέτες περιπτώσεων μεγάλων βιομηχανιών.

Η BJC HealthCare είναι ένας πάροχος υπηρεσιών υγείας. Η εταιρεία χρησιμοποιεί τεχνολογίες αναγνώρισης ραδιοσυχνότητας (Radio Frequency Identification – RFID) για να παρακολουθεί και να διαχειρίζεται χιλιάδες ιατρικές προμήθειες[17]. Μέσω της τεχνολογίας αυτής διαβάζονται και καταγράφονται πληροφορίες που αποθηκεύονται στις ετικέτες των αντικειμένων. Η υλοποίηση της τεχνολογίας αυτής επέφερε σημαντικά οφέλη, η BJC κατάφερε να μειώσει την ποσότητα αποθεμάτων που διατηρούσε σε κάθε εγκατάσταση κατά 23%. Η εταιρεία προβλέπει ότι θα υπάρξουν συνεχείς εξοικονομήσεις περίπου 5 εκατομμυρίων δολαρίων ετησίως.

Η Bosch χρησιμοποιεί το συνδυασμό των IIoT και Big Data για να οδηγήσει το εργοστάσιο Bosch Automotive Diesel System στο ψηφιακό μετασχηματισμό [2]. Η εταιρεία συνδέει τα μηχανήματα της μέσω των ενσωματωμένων αισθητήρων που συλλέγουν δεδομένα σχετικά με τις συνθήκες και την περίοδο κύκλου εργασιών,

---

στον πυρήνα του εργοστασίου έτσι ώστε να μπορεί να παρακολουθεί τη συνολική διαδικασία παραγωγής. Η υιοθέτηση αυτής της προσέγγισης βοηθάει στην πρόβλεψη των αστοχιών του εξοπλισμού, επιτρέποντας στο εργοστάσιο να προγραμματίζει τις λειτουργίες συντήρησης πολύ πριν προκύψουν οποιεσδήποτε αστοχίες. Η εταιρεία δηλώνει ότι αυτός ο τρόπος χρήσης ανάλυσης δεδομένων έχει συμβάλει στην αύξηση της παραγωγής σε ποσοστό μεγαλύτερο του 10% σε ορισμένους τομείς, βελτιώνοντας παράλληλα την παράδοση και την ικανοποίηση των πελατών.

Τα συνδεδεμένα αυτοκίνητα είναι μια τάση στην αυτοκινητοβιομηχανία, η οποία προσφέρει ψηφιακές υπηρεσίες προστιθέμενης αξίας στους πελάτες. Ένας από τους πρώτους κατασκευαστές αυτοκινήτων που έκανε την αρχή είναι η Volkswagen, η οποία ένωσε τις δυνάμεις της με την Microsoft για να αναπτύξουν ένα cloud network, το “Volkswagen Automotive Cloud”. Η τεχνολογία προσφέρει μια σειρά από χαρακτηριστικά όπως η έξυπνη οικιακή συνδεσιμότητα, προσωπικό ψηφιακό βοηθό, υπηρεσία προγνωστικής συντήρησης, ροή πολυμέσων και ενημερώσεις. Η Volkswagen στοχεύει να προσθέσει πάνω από 5 εκατομμύρια προϊόντα ανά έτος στο IoT με τη βοήθεια της υπηρεσίας cloud[35].

Η Fetch Robotics [11] έχει αναπτύξει συνεργατικά αυτόνομα κινητά ρομπότ (Autonomous Mobile Robots – AMR) για τον εντοπισμό, την παρακολούθηση και τη μεταφορά αποθεμάτων σε εγκαταστάσεις αποθήκευσης και εφοδιασμού. Το κέντρο διανομής DHL χρησιμοποιεί AMR Fetch για εργασίες επιλογής και τοποθέτησης. Οι AMR μετακινούνται αυτόνομα στις εγκαταστάσεις μαζί με τους εργαζομένους, μαθαίνοντας και μοιράζοντας αυτόνομα τις πιο αποδοτικές διαδρομές. Η χρήση αυτό-οδηγούμενων ρομπότ μπορεί να βοηθήσει στη μείωση του χρόνου κύκλου παραγγελιών έως 50% και να παρέχει έως και το διπλάσιο κέρδος παραγωγικότητας, σύμφωνα με την εταιρεία.

Στην προσπάθεια να επισπεύσει τη διαδικασία ανάπτυξης του αγωνιστικού αυτοκινήτου, η Team Penske [26] συνεργάστηκε με την Siemen, αποκτώντας πρόσβαση σε προηγμένες λύσεις ψηφιακού σχεδιασμού και προσομοίωσης – συμπεριλαμβανομένων των ψηφιακών διδύμων (digital twins). Μέσω των ψηφιακών διδύμων παρέχονται στους μηχανικούς της Team Penske μια εικονική κλίνη δοκιμών για την έρευνα νέων εξαρτημάτων, βελτιστοποιώντας την απόδοση του αυτοκινήτου. Για την Team Penske αυτό τελικά ισοδυναμεί με μια φθηνότερη, πιο αποδοτική ως προς τους

---

πόρους διαδικασία ελέγχου του προϊόντος και την ανάπτυξη ταχύτερων οχημάτων.

Η General Electric [5] δίνει μια ιδέα για το πως η τεχνολογία AR μπορεί να ενισχύσει την παραγωγή. Η εταιρεία κάνει χρήση των γυαλιών AR στις εγκαταστάσεις παραγωγής κινητήρων τζετ. Πριν από τη χρήση των έξυπνων γυαλιών οι κατασκευαστές μηχανών τζετ έπρεπε συχνά να διακόπτουν την εργασία τους για να ελέγξουν τα εγχειρίδια τους και να διασφαλίσουν ότι οι εργασίες εκτελούνται σωστά. Ωστόσο, με τα γυαλιά AR, μπορούν να λαμβάνουν ψηφιοποιημένες οδηγίες στο οπτικό τους πεδίο καθώς και να έχουν πρόσβαση σε εκπαιδευτικά βίντεο ή να χρησιμοποιούν φωνητικές εντολές για να επικοινωνούν με ειδικούς για άμεση βοήθεια. Κατά τη διάρκεια του προγράμματος, η GE αναφέρει ότι η παραγωγικότητα των εργαζομένων αυξήθηκε έως και κατά 11%, σε σύγκριση με το παρελθόν.



# Κεφάλαιο 4

## Εργαλεία Υλοποίησης

### 4.1 Εισαγωγή

Η ικανότητα εκτέλεσης διαδικασιών ανάλυσης δεδομένων και η δημιουργία προγνωστικών μοντέλων αποτελεί ένα σημαντικό κομμάτι στον σημερινό κόσμο όπου ο ρυθμός παραγωγής δεδομένων ολοένα και αυξάνεται. Για να επιτευχθεί αυτό, οι ερευνητές και οι επαγγελματίες χρησιμοποιούν μια ποικιλία εργαλείων τα οποία έχουν σχεδιαστεί με σκοπό τη διευκόλυνση των εργασιών ανάλυσης δεδομένων και μοντελοποίησης. Σε αυτή την ενότητα θα παρουσιαστούν κάποια από τα πιο γνωστά εργαλεία που χρησιμοποιούνται στην Python για ανάλυση δεδομένων και μηχανική μάθηση συμπεριλαμβανομένων των Jupyter Notebook, Pandas, Scikit-Learn και Matplotlib.

### 4.2 Jupyter Notebook

Το Jupyter Notebook αποτελεί ένα διαδικτυακό, δραστικό, υπολογιστικό περιβάλλον που επιτρέπει σε ερευνητές και επαγγελματίες το δημιουργήσουν και να μοιραστούν αρχεία τα οποία περιέχουν κώδικα, εξισώσεις, οπτικοποιήσεις και αφηγηματικό κείμενο. Πρόκειται για ένα δημοφιλές εργαλείο το οποίο είναι ευρέως χρησιμοποιούμενο στην επιστήμη των δεδομένων για εξόρυξη, ανάλυση και εργασίες μηχανικής μάθησης. Μερικά από τα πλεονεκτήματα του Jupyter Notebook είναι η ικανότητα του να υποστηρίζει πολλαπλές γλώσσες προγραμματισμού, τον ενσωματωμένο γραφικό απεικονισμό με τη χρήση του Matplotlib καθώς και η ικανότητα του να μοιράζεται και να συνεργάζεται σε άλλα σημειωματάρια με άλλους χρήστες.

---

## 4.3 Python

Η Python είναι μια γλώσσα προγραμματισμού υψηλού επιπέδου που χρησιμοποιείται ευρέως στους επιστημονικούς υπολογισμούς, την ανάλυση δεδομένων και στη μηχανική μάθηση. Παρέχει μια μεγάλη και ενεργή κοινότητα προγραμματιστών καθώς και ένα πλούσιο περιβάλλον βιβλιοθηκών και εργαλείων. Η Python είναι γνωστή για την απλότητα, την αναγνωσιμότητα και την ευκολία χρήσης της, καθιστώντας την ιδανική γλώσσα τόσο για αρχάριους όσο και για ειδικούς. Με τις εκτεταμένες βιβλιοθήκες και τα εργαλεία της, η Python επιτρέπει τη γρήγορη δημιουργία πρωτοτύπων και την ανάπτυξη σύνθετων έργων ανάλυσης δεδομένων και μηχανικής μάθησης.

## 4.4 Pandas

Η Pandas είναι μια βιβλιοθήκη της Python που χρησιμοποιείται για την επεξεργασία και την ανάλυση δεδομένων. Παρέχει γρήγορες και ευέλικτες δομές δεδομένων καθώς και εργαλεία κατάλληλα για τις δομές αυτές. Οι δύο κύριες δομές δεδομένων στη Pandas είναι οι Series και DataFrame. Η Pandas περιλαμβάνει επίσης ισχυρές δυνατότητες ευρετηρίασης και επιλογής, επιτρέποντας την αποτελεσματική πρόσβαση και τροποποίηση των δεδομένων. Αυτό την καθιστά ιδιαίτερα χρήσιμη για εργασίες επεξεργασίας δεδομένων, όπως ο καθαρισμός, ο μετασχηματισμός και ο συνδυασμός δεδομένων.

## 4.5 Scikit-Learn

Το Scikit-Learn είναι μια δημοφιλής βιβλιοθήκη μηχανικής μάθησης στη Python και παρέχει ένα ευρύ φάσμα αλγορίθμων για κατηγοριοποίηση, παλινδρόμηση, ομαδοποίηση και μείωση διαστάσεων. Περιλαμβάνει επίσης εργαλεία για την επεξεργασία δεδομένων, την επιλογή χαρακτηριστικών και την αξιολόγηση μοντέλων. Το Scikit-Learn διαθέτει ένα απλό και συνεπές API, καθιστώντας εύκολη τη χρήση και την εναλλαγή μεταξύ των διαφορετικών μοντέλων. Παρέχοντας εύκολη πρόσβαση σε ένα ευρύ φάσμα αλγορίθμων μηχανικής μάθησης, το Scikit-learn διευκολύνει την εξερεύνηση διαφορετικών μοντέλων και τη σύγκριση των επιδόσεων

---

τους.

## 4.6 Matplotlib

Το Matplotlib αποτελεί και αυτό μια βιβλιοθήκη της Python η οποία ειδικεύεται στη δημιουργία στατικών, κινούμενων και διαδραστικών απεικονίσεων στην Python. Παρέχει ένα ευρύ φάσμα διαγραμμάτων και γραφικών παραστάσεων συμπεριλαμβανομένων των γραμμικών διαγραμμάτων, των διαγραμμάτων διασποράς, ιστογράμματα, ραβδογράμματα και πολλά άλλα. Το Matplotlib παρέχει έναν υψηλό βαθμό ευελιξίας και ελέγχου της εμφάνισης των γραφημάτων, επιτρέποντας στους χρήστες να προσαρμόσουν κάθε πτυχή του. Αποτελεί μια εξαιρετικά επεκτάσιμη βιβλιοθήκη, με πολλά εργαλεία και plugins τρίτων κατασκευαστών διαθέσιμα για συγκεκριμένες περιπτώσεις χρήσης.

## Κεφάλαιο 5

# Ανάλυση Εφαρμογών και Προεπεξεργασία Δεδομένων

Στη σημερινή εποχή, ο βιομηχανικός τομέας παράγει τεράστιες ποσότητες δεδομένων που μπορούν να χρησιμοποιηθούν για τη βελτίωση της αποδοτικότητας της εφαρμογής, να μειώσουν το χρόνο της διακοπής λειτουργίας και να βελτιώσουν την ποιότητα παραγωγής. Ωστόσο, για να ξεκλειδωθεί το πλήρες δυναμικό αυτών των δεδομένων, είναι σημαντικό να πραγματοποιηθεί πρώτα η ανάλυση και η προεπεξεργασία τους. Η διαδικασία αυτή περιλαμβάνει τη μετατροπή των ακατέργαστων δεδομένων σε μορφή κατάλληλη για αλγορίθμους μηχανικής μάθησης, οι οποίοι μπορούν στη συνέχεια να χρησιμοποιηθούν για την εξαγωγή πληροφοριών σχετικά με πιθανές βλάβες στο σύστημα.

Μια από τις σημαντικότερες προκλήσεις κατά την εργασία με δεδομένα βιομηχανικών συστημάτων είναι η παρουσία σφαλμάτων ή ανωμαλιών στα δεδομένα. Αυτές μπορεί να οφείλονται σε διάφορους παράγοντες όπως δυσλειτουργίες εξοπλισμού, περιβαλλοντικούς παράγοντες ή ανθρώπινα λάθη. Οι βλάβες αυτές μπορούν να οδηγήσουν σε σημαντικές απώλειες όσον αφορά την αποδοτικότητα της παραγωγής, την ποιότητα και την ασφάλεια. Ως εκ τούτου, είναι σημαντικός ο εντοπισμός και η ανάλυση τέτοιου είδους σφαλμάτων προκειμένου να μετριάζονται οι επιπτώσεις τους. Με την εφαρμογή αλγορίθμων μηχανικής μάθησης σε βιομηχανικά δεδομένα, καθίσταται δυνατός ο εντοπισμός και η διάγνωση αυτών των βλαβών σε πραγματικό χρόνο, επιτρέποντας τη λήψη άμεσων διορθωτικών μέτρων. Στην ενότητα αυτή, θα γίνει μια εισαγωγή στην κάθε εφαρμογή. Έπειτα, θα πραγματοποιηθεί η διαδικασία της ανάλυσης και της προεπεξεργασίας των δεδομένων με σκοπό την ανάλυση τους

---

στο επόμενο κεφάλαιο.

## 5.1 Συστήματα ηλεκτρικής ενέργειας

Σχήμα 5.1: Συστήματα ηλεκτρικής ενέργειας [4]



Σε αυτή την ενότητα εξετάζονται δεδομένα που σχετίζονται με συστήματα ηλεκτρικής ενέργειας (Σχήμα 5.1). Με την πάροδο του χρόνου τα συστήματα αυτά αυξάνονται τόσο σε μέγεθος όσο και σε πολυπλοκότητα και λαμβάνουν χώρα σε διάφορους τομείς όπως η παραγωγή, η μεταφορά, η διανομή και σε συστήματα φόρτωσης. Τύποι σφαλμάτων όπως το βραχυκύκλωμα, σε ένα σύστημα ηλεκτρικής ενέργειας, έχουν ως αποτέλεσμα σοβαρές οικονομικές απώλειες καθώς και τη μείωση της αξιοπιστίας. Το ηλεκτρικό σφάλμα είναι μια ασυνήθιστη κατάσταση που προκαλείται από σφάλματα στον εξοπλισμό, όπως στους μετασχηματιστές και τις περιστρεφόμενες μηχανές, ανθρώπινα λάθη και περιβαλλοντικές συνθήκες. Τα σφάλματα αυτά προκαλούν διακοπές στην ηλεκτρική ροή, βλάβες στον εξοπλισμό και θάνατο σε ζωντανούς οργανισμούς.

Με τον όρο ηλεκτρικό σφάλμα γίνεται αναφορά στην απόκλιση των τάσεων και ρευμάτων από τις ονομαστικές τους τιμές ή καταστάσεις. Στην περίπτωση που προκύψει κάποιο σφάλμα, προκαλείται ροή υπερβολικά υψηλών ρευμάτων και έτσι προκαλούνται βλάβες σε εξοπλισμούς και συσκευές. Ο εντοπισμός των σφαλμάτων και η ανάλυση τους είναι απαραίτητο κομμάτι για την επιλογή ή την κατασκευή κα-

---

τάλληλων διακοπών, ηλεκτρομηχανικών ηλεκτρονόμων και άλλων συσκευών προστασίας. Υπάρχουν κυρίως δύο τύποι σφαλμάτων στα συστήματα ηλεκτρικής ενέργειας, τα Συμμετρικά και τα Μη Συμμετρικά. Στην κατηγορία των Συμμετρικών σφαλμάτων ανήκουν τα πολύ σοβαρά τα οποία εμφανίζονται σπάνια στα συστήματα. Τους αποδίδεται και η ονομασία των εξισορροπημένων σφαλμάτων και είναι δύο τύπων, δηλαδή (L-L-L-G) και (L-L-L) όπου L για γραμμή και G για γείωση. Μόνο το 2-5% των σφαλμάτων είναι συμμετρικά και εάν εμφανιστούν το σύστημα παραμένει ισορροπημένο αλλά οδηγείται σε σοβαρές βλάβες στον εξοπλισμό του. Η ανάλυση τους είναι εύκολη και πραγματοποιείται ανά φάση. Τα Μη Συμμετρικά σφάλματα είναι πολύ συνηθισμένα και λιγότερο σοβαρά σε σχέση με τα προηγούμενα. Υπάρχουν κυρίως τρεις τύποι, (L-G),(L-L) και (LL-G) όπου LL για διπλή γραμμή. Τα σφάλματα γραμμή προς γείωση(L-G) είναι τα πιο κοινά με ποσοστό 65-70% και προκαλούν την επαφή του αγωγού με τη γη. Το 15-20% των σφαλμάτων είναι διπλά σφάλματα γραμμής προς γη (LL-G) και προκαλούν την επαφή των δύο αγωγών με τη γη. Τέλος, τα σφάλματα γραμμής προς γραμμή (L-L) καταλαμβάνουν ένα ποσοστό 5-10% και συμβαίνουν όταν δύο αγωγοί έρχονται σε επαφή μεταξύ τους κατά την ταλάντωση των γραμμών λόγω ανέμων. Τα Μη Συμμετρικά σφάλματα ονομάζονται επίσης και μη ισορροπημένα καθώς η εμφάνισή τους προκαλεί ανισορροπία στο σύστημα, δηλαδή οι τιμές της σύνθετης αντίστασης είναι διαφορετικές σε κάθε φάση με αποτέλεσμα να υπάρχει ασύμμετρη ροή ρεύματος στις φάσεις. Τέτοια σφάλματα είναι πιο δύσκολο να αναλυθούν.

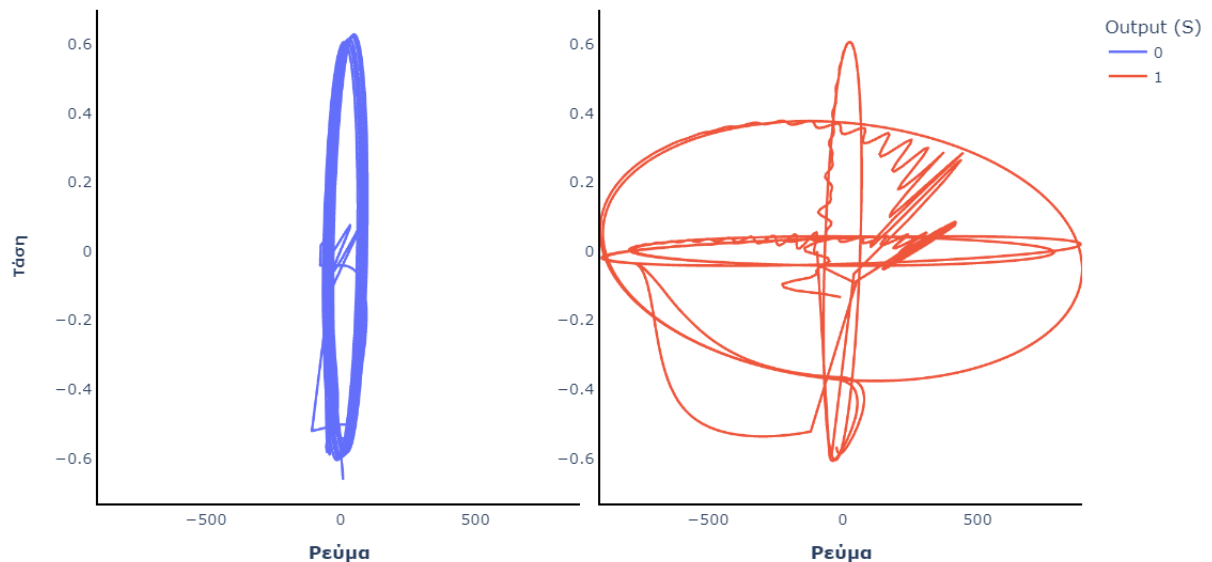
## 5.2 Ανάλυση δεδομένων

### Δυαδική κατηγοριοποίηση

Χρησιμοποιήθηκε το αρχείο detect\_dataset.csv το οποίο περιέχει μετρήσεις που υποδεικνύουν αν στο σύστημα εμφανίζεται κάποιο σφάλμα ή όχι. Αφού έγινε η εισαγωγή των δεδομένων από το αρχείο detect\_dataset.csv διαγράφηκαν οι στήλες “Unnamed:7” και “Unnamed:8” καθώς είχαν μηδενικές τιμές και από τη στιγμή που δεν περιέχουν κάποια πληροφορία μια καλή τεχνική είναι να αφαιρεθούν. Έπειτα, παρατηρήθηκε ότι η έξοδος εμφανίζεται ως 0 ή 1 και ως εκ τούτου το πρόβλημα που αναλύεται εντάσσεται στην κατηγορία των δυαδικών προβλημάτων κατηγοριο-

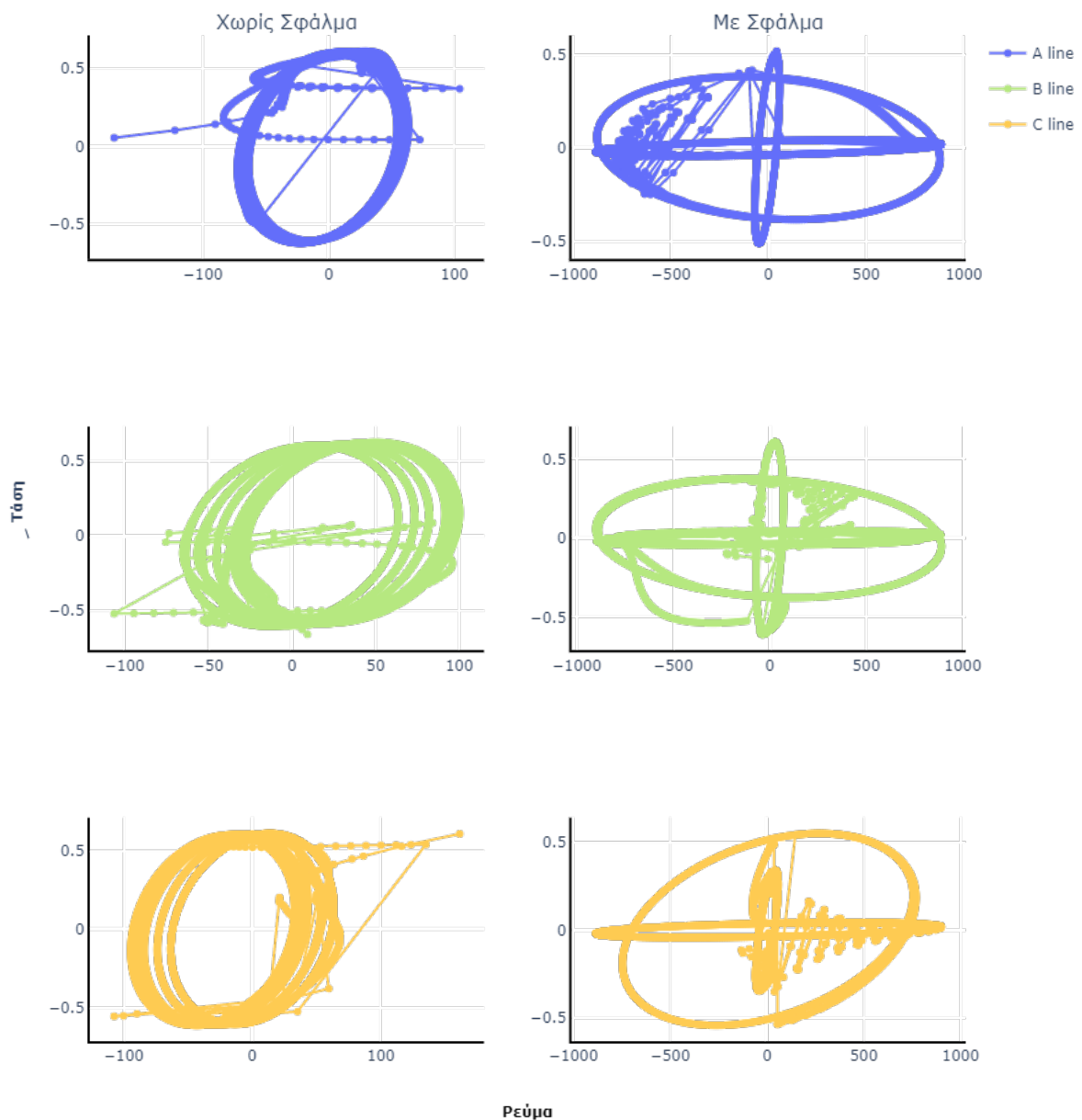
ποίηση. Πριν σχεδιαστεί το μοντέλο γίνεται μια ανάλυση της συμπεριφοράς του ρεύματος και της τάσης σε ελαττωματική και μη κατάσταση. Τα παρακάτω διαγράμματα σχεδιάστηκαν με την χρήση του Plotly.js.

Σχήμα 5.2: Ρεύμα και τάση στη γραμμή Β σε φυσιολογικές και μη φυσιολογικές συνθήκες



Το Σχήμα 5.2 απεικονίζει τη σχέση ανάμεσα στην τάση και το ρεύμα σε φυσιολογικές (0) και μη (1) συνθήκες. Στο πρώτο διάγραμμα υπάρχει μια ομαλή σχέση ανάμεσα στο ρεύμα και την τάση ενώ στο δεύτερο η μεταβολή των δύο τιμών είναι εντελώς ασυνάρτητη. Το διάγραμμα 5.3 απεικονίζει τη σχέση μεταξύ ρεύματος και τάσης στις γραμμές Α, Β και C σε συνθήκες με και χωρίς σφάλματα.

Σχήμα 5.3: Ρεύμα και τάση στις γραμμές σε φυσιολογικές και μη φυσιολογικές συνθήκες



### Πολυταξική κατηγοριοποίηση

Στην περίπτωση αυτή χρησιμοποιήθηκε το αρχείο classData.csv το οποίο περιέχει πληροφορίες σχετικά με τον τύπο του σφάλματος. Με την εισαγωγή των δεδομένων από το αρχείο classData.csv έγινε προσθήκη της στήλης FaultType, του αλφαριθμητικού το οποίο υποδεικνύει τον κωδικό του σφάλματος. Συγκεκριμένα υπάρχουν έξι κατηγορίες κωδικών:



- 
1. Σφάλμα γραμμής – γραμμής : Line – Line (LL)
  2. Σφάλμα γραμμής – γείωσης : Line – Ground(LG)
  3. Σφάλμα γραμμής - γραμμής – γείωσης : Line – Line – Ground (LLG)
  4. Σφάλμα γραμμής – γραμμής – γραμμής – γείωσης: Line – Line – Line –Line- Ground (LLLG)
  5. Σφάλμα γραμμής – γραμμής – γραμμής : Line – Line – Line (LLL)

Τα Ia, Ib, Ic και Va, Vb Vc υποδεικνύουν το ρεύμα και την τάση στις αντίστοιχες γραμμές. Ενώ τα G, C, B, A υποδεικνύουν τη γείωση και τα σφάλματα στις γραμμές C, B και A αντίστοιχα.

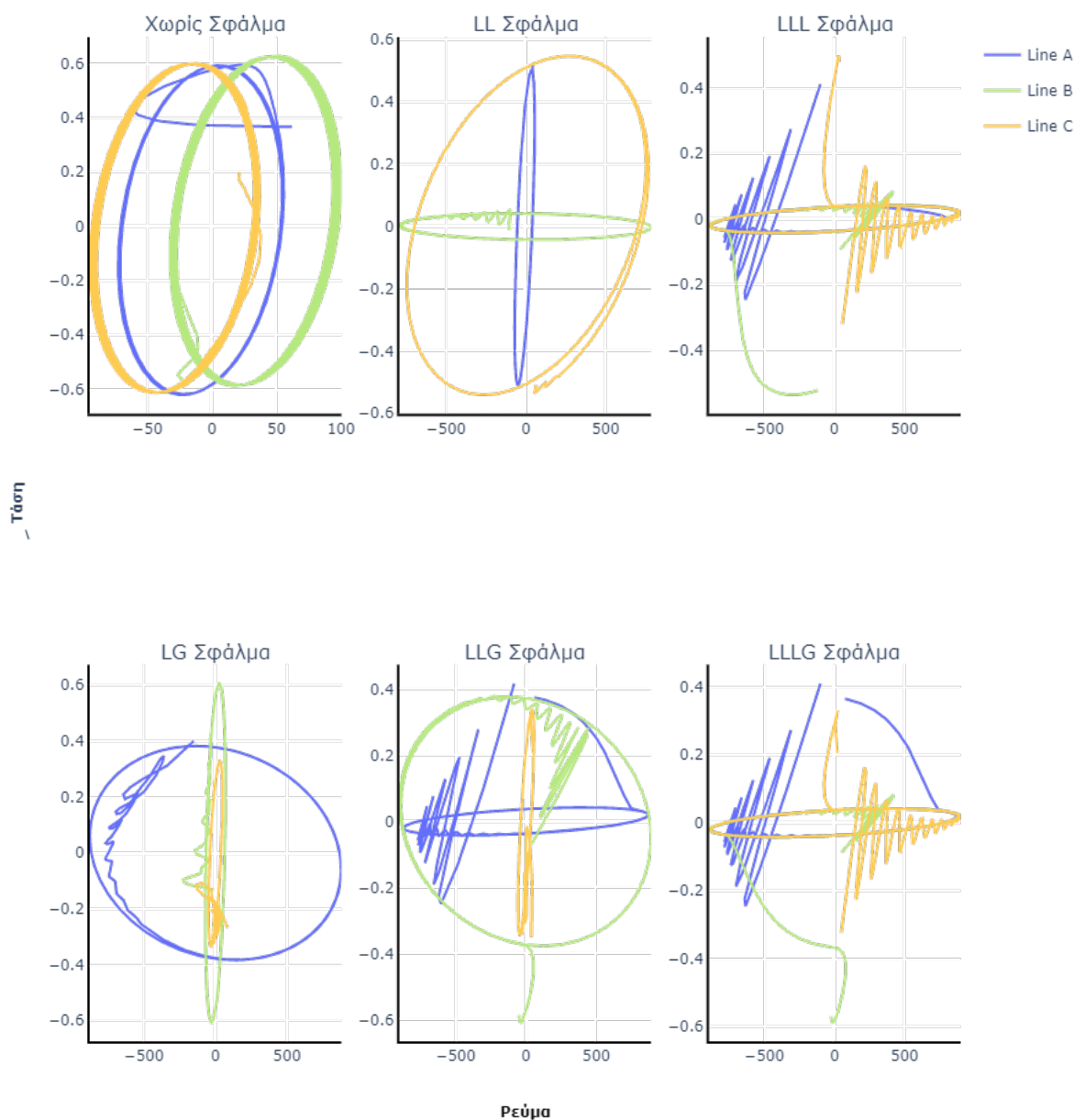
Παρακάτω γίνεται αναφορά των κατηγοριών των εξόδων που υπάρχουν στο πρόβλημα:

1. 0000 – Κανένα σφάλμα
2. 1001 - Σφάλμα LG ( Μεταξύ φάσης A και γείωσης)
3. 0110 – Σφάλμα LL (Μεταξύ φάσης B και C)
4. 1011 – Σφάλμα LLG ( Μεταξύ A, B και γείωσης)
5. 0111 – Σφάλμα LLL (Μεταξύ των τριών φάσεων)
6. 1111 - Σφάλμα LLLG ( Τριφασικό συμμετρικό σφάλμα)

Παρατηρείται ότι υπάρχουν έξι διαφορετικά είδη εξόδων και ως εκ τούτο το πρόβλημα θα χειριστεί ως πολυταξική κατηγοριοποίηση.

Το Σχήμα 5.4 περιγράφει τη σχέση των τάσεων και των ρευμάτων και στις έξι περιπτώσεις των προαναφερθέντων εξόδων. Όπως στη δυαδική κατηγοριοποίηση έτσι και στην πολυταξική παρατηρείται ότι στην περίπτωση που δεν υπάρχει σφάλμα, μεταξύ των τριών γραμμών υπάρχει μια σχέση ανάμεσα στο ρεύμα και στην τάση. Παρόλα αυτά στις περιπτώσεις σφάλματος υπάρχει μια παραμόρφωση της σχέσης. Επίσης στα διαγράμματα σφαλμάτων LLL και LLLG υπάρχει μια ομοιότητα μεταξύ τους.

Σχήμα 5.4: Ρεύμα και τάση στις γραμμές A, B και C σε διάφορες συνθήκες



### 5.3 Πλάνο μελέτης

Το πρώτο βήμα ήταν ο κατάλληλος διαχωρισμός και η προετοιμασία των δεδομένων. Έπειτα, έγινε η επιλογή των αλγορίθμων που χρησιμοποιήθηκαν για το σχεδιασμό του μοντέλου. Το 80% των δεδομένων χρησιμοποιήθηκε για την εκπαίδευση των αλγορίθμων και το 20% για την αξιολόγηση τους. Οι αλγόριθμοι που χρησιμοποιήθηκαν στην κατηγοριοποίηση αναφέρονται παρακάτω:

- 
1. Logistic Regression
  2. KNeighbors Classifier
  3. Random Forest
  4. Gaussian NB
  5. XGB Classifier
  6. D - Tree
  7. MLP
  8. SVC

Σε επόμενο βήμα, εξάγονται μετρήσεις που αφορούν το ποσοστό ακρίβειας και σφάλματος για το κάθε μοντέλο στην κατάσταση δημιουργίας. Έπειτα παρουσιάζεται ένα συγκεντρωτικό διάγραμμα και ο αντίστοιχος πίνακας με τις μετρικές στη κατάσταση της διασταυρούμενης επικύρωσης (cross validation), όπου στόχο έχει την αξιολόγηση του μοντέλου. Τέλος, δίνεται βάση στις χρονικές μεταβλητές και παρατίθενται ένα διάγραμμα και ο αντίστοιχος πίνακας με τις χρονικές μετρικές εκπαίδευσης και αξιολόγησης του κάθε μοντέλου.

## 5.4 Τεχνικές αξιολόγησης

Κατά την ερμηνεία των αποτελεσμάτων στα πειράματα μηχανικής μάθησης έχει σημαντική σημασία να λαμβάνονται υπόψη διάφοροι παράγοντες έτσι ώστε να επιτευχθεί μια ολοκληρωμένη κατανόηση της απόδοσης του μοντέλου. Η αξιολόγηση των μετρικών επιδόσεων που προκύπτουν από το cross-validation, η κατανόηση της τυπικής απόκλισης των μετρικών αυτών και η ανάλυση των σταδίων προσαρμογής και αξιολόγησης έχουν σημαντικό ρόλο στην ερμηνεία και στη σύγκριση αποτελεσμάτων. Το cross-validation βοηθάει στην εκτίμηση της ικανότητας γενίκευσης του μοντέλου, ενώ η τυπική απόκλιση παρέχει πληροφορίες σχετικά με τη συνέπεια και τη σταθερότητα των επιδόσεων του μοντέλου. Τέλος, η εξέταση των χρόνων προσαρμογής και αξιολόγησης επιτρέπει τη βαθύτερη κατανόηση των επιδόσεων του

---

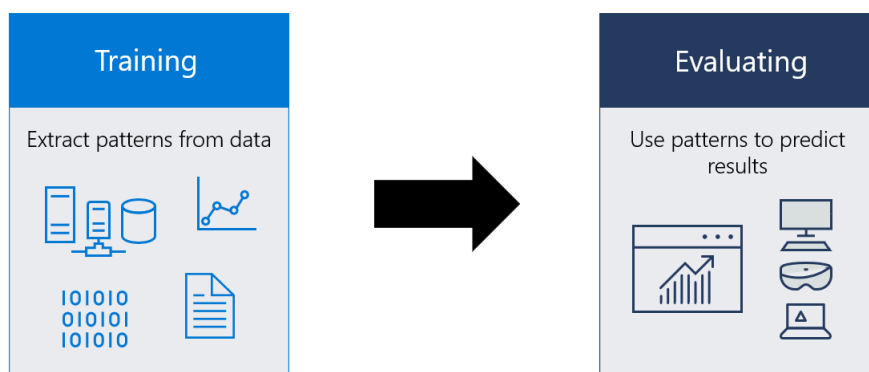
μοντέλου σε άγνωστα δεδομένα. Λαμβάνοντας υπόψη αυτά τα σημεία, μπορεί κανείς να λάβει τεκμηριωμένες αποφάσεις και να βγάλει ουσιαστικά συμπεράσματα από τα αποτελέσματα των πειραμάτων μηχανικής μάθησης.

# Κεφάλαιο 6

## Πειραματικές Εφαρμογές και Αποτελέσματα

Η μηχανική μάθηση έχει εξελιχθεί σε ένα ισχυρό εργαλείο ανάλυσης και μοντελοποίησης συνόλων δεδομένων σε διάφορους τομείς, συμπεριλαμβανομένων και των βιομηχανικών συστημάτων. Με τον ολοένα αυξανόμενο αριθμό δεδομένων που παράγονται από βιομηχανικές διεργασίες, υπάρχει ανάγκη για αποτελεσματικές και ακριβείς μεθόδους ανάλυσης και ερμηνείας των δεδομένων αυτών. Οι τεχνικές μηχανικής μάθησης έχουν αποδείξει ότι είναι μια πολύτιμη λύση σε τέτοιου είδους προβλήματα, καθώς είναι ικανές να μαθαίνουν αυτόματα μοτίβα και σχέσεις από τα δεδομένα, χωρίς την ανάγκη σαφή προγραμματισμού, Σχήμα 6.1.

Σχήμα 6.1: Διαδικασία μελέτης [13]



Στο πλαίσιο αυτό, η σύγκριση διαφορετικών μεθόδων μηχανικής μάθησης είναι σημαντική για την επιλογή της καταλληλότερης προσέγγισης για μια συγκεκριμένη εργασία. Η αποτελεσματικότητα μιας μεθόδου μηχανικής μάθησης μετριέται συνήθως με την ακρίβεια, την ερμηνευσιμότητα, την υπολογιστική αποτελεσματικότητα

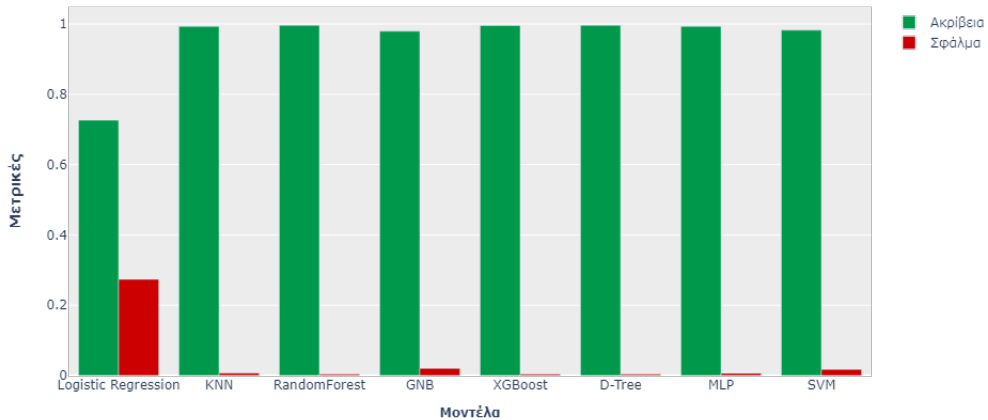
και την ανθεκτικότητα του στο θόρυβο και στις ακραίες τιμές. Συγκρίνοντας διαφορετικές μεθόδους μπορούν να εξαχθούν γνώσεις σχετικά με τα δυνατά και αδύνατα σημεία κάθε μεθόδου και να επιλεχθεί εκείνη που ταιριάζει καλύτερα στις ανάγκες του προβλήματος. Στη συνέχεια της ενότητας θα διερευνηθούν οι εφαρμογές διάφορων μεθόδων μηχανικής μάθησης και θα συγκριθούν οι αποδόσεις τους με βάση ένα σύνολο μετρήσεων.

## 6.1 Σχεδιασμός μοντέλου δυαδικής κατηγοριοποίησης

### 6.1.1 Ανάλυση

Όπως αναφέρθηκε σε προηγούμενο κεφάλαιο σχετικά με το πλάνο της εργασίας η ανάλυση ξεκινάει με την εξαγωγή του διαγράμματος Ακρίβειας – Σφάλματος, Σχήμα 6.2 και του αντίστοιχου Πίνακα 6.1.

Σχήμα 6.2: Ακρίβεια και σφάλμα στη δυαδική κατηγοριοποίηση



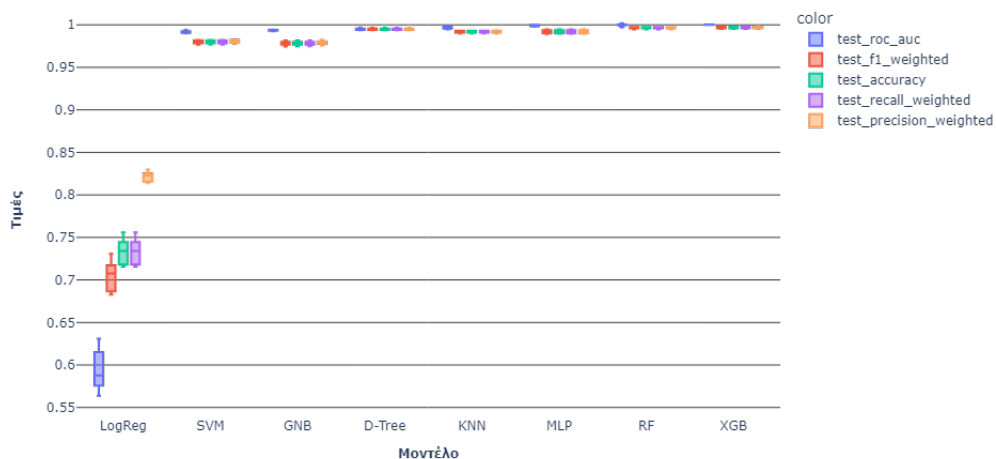
Αφού ολοκληρώθηκε η εκμάθηση των αλγορίθμων υπολογίστηκε για το κάθε ένα το ποσοστό ακρίβειας και σφάλματος. Όπως παρατηρείται στο διάγραμμα του Σχήματος 6.2 και στον αντίστοιχο Πίνακα 6.1 το χειρότερο ποσοστό κατέχει ο αλγόριθμος Logistic Regression, ενώ τα καλύτερα ποσοστά οι αλγόριθμοι Random Forest και Decision Tree, έχοντας μάλιστα και τις ίδιες τιμές. Την αμέσως επόμενη καλύτερη θέση κατέχει ο XGBoost με ελάχιστη διαφορά από τους Random Forest και Decision – Tree.

Μέθοδος	Ακρίβεια	Σφάλμα
LogReg	0.726299	0.273701
KNN	0.993257	0.006743
Random Forest	0.996429	0.003570
GNB	0.979769	0.020230
XGBoost	0.996033	0.003967
D-Tree	0.996429	0.003570
MLP	0.993653	0.006347
SVM	0.982547	0.017453

Πίνακας 6.1: Πίνακας τιμών ακρίβειας - σφάλματος στη δυαδική κατηγοριοποίηση

Έκτος από την ακρίβεια και το σφάλμα έγινε ανάλυση του κάθε μοντέλου μέσω του cross validation για να αναλυθούν οι επιδόσεις τους. Επομένως, παρατίθενται και άλλες μετρικές στο διάγραμμα του Σχήματος 6.3 και στον Πίνακα 6.2 ώστε να υπάρξει μια πλήρης εικόνα για την επιλογή του.

Σχήμα 6.3: Μετρικές κατάστασης cross validation δυαδικής κατηγοριοποίησης



Μοντέλο	Ακρίβεια		Roc_Auc		Ανάκληση		F1		Ακρίβεια Θετικών Προβλέψεων	
	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος
LogReg	0.016286	0.733122	0.026244	0.594376	0.016286	0.733122	0.019410	0.704351	0.006313	0.821269
SVM	0.002129	0.979747	0.001461	0.991455	0.002129	0.979747	0.002141	0.979702	0.001964	0.980467
GNB	0.002557	0.978376	0.000709	0.993383	0.002557	0.978376	0.002570	0.978324	0.002353	0.979197
D-Tree	0.001214	0.995042	0.001087	0.994973	0.001214	0.995042	0.001214	0.995042	0.001214	0.995046
KNN	0.001415	0.991983	0.001417	0.996482	0.001415	0.991983	0.001418	0.991978	0.001378	0.992048
MLP	0.002110	0.992089	0.000732	0.999081	0.002110	0.992089	0.002113	0.992084	0.002082	0.992134
RF	0.001094	0.996519	0.000240	0.999850	0.001094	0.996519	0.001093	0.996519	0.001093	0.996525
XGB	0.001203	0.996730	0.000011	0.999970	0.001203	0.996730	0.001203	0.996730	0.001203	0.996737

Πίνακας 6.2: Πίνακας μετρικών κατάστασης cross validation στη δυαδική κατηγοριοποίηση

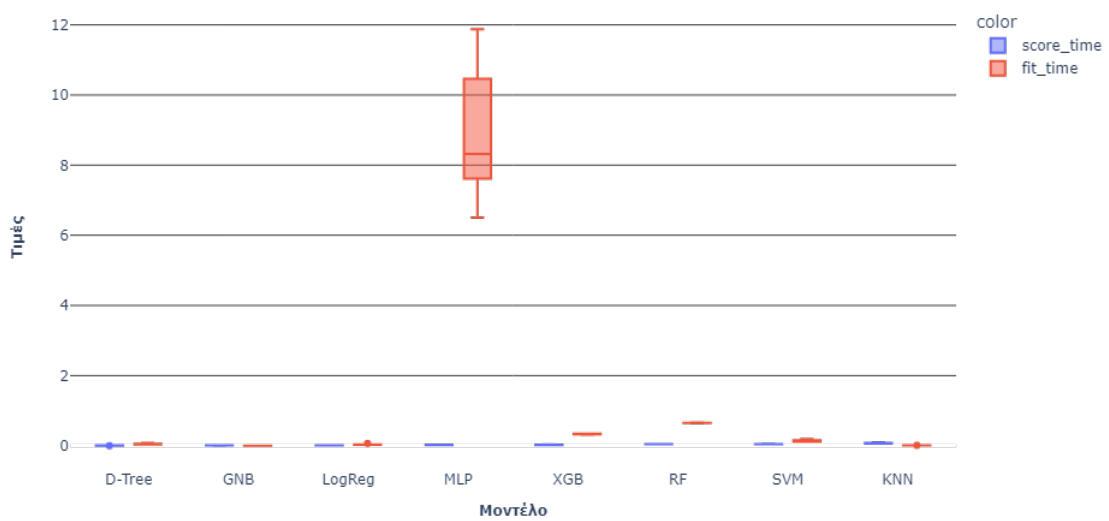
Παρατηρείται ότι το μοντέλο του Logistic Regression έχει τις χειρότερες αποδόσεις, ενώ οι Random Forest και XGBoost έχουν τις καλύτερες, με τον τελευταίο να

έχει την πρώτη θέση.

Προκειμένου να γίνει πιο κατανοητή η δομή των δεδομένων χρησιμοποιήθηκε επίσης η τυπική απόκλιση και ο μέσος των μετρικών. Παρατηρείται επίσης ότι οι τιμές της τυπικής απόκλισης σε κάθε μετρική είναι πιο κοντά στο μέσο, πράγμα που επιβεβαιώνει την αξιοπιστία των δεδομένων.

Σε επόμενο βήμα γίνεται μια χρονική ανάλυση των μοντέλων με τη χρήση των μετρικών fit και score time:

Σχήμα 6.4: Σύγκριση των μοντέλων μέσω του χρόνου fit και score δυαδικής κατηγοριοποίησης



Μοντέλο	Χρόνος Εκπαίδευσης		Χρόνος Αξιολόγησης	
	Απόκλιση	Μέσος	Απόκλιση	Μέσος
D-Tree	0.018011	0.044890	0.002603	0.004331
GNB	0.000646	0.003398	0.003092	0.005779
LogReg	0.016287	0.034266	0.001805	0.008785
MLP	2.058908	8.934588	0.003726	0.023914
XGB	0.008778	0.330690	0.004206	0.024080
RF	0.008145	0.650559	0.002584	0.042374
SVM	0.031784	0.136744	0.002619	0.044368
KNN	0.002081	0.005990	0.015549	0.069994

Πίνακας 6.3: Πίνακας χρονικών τιμών fit - score δυαδικής κατηγοριοποίησης



---

Σύμφωνα με το διάγραμμα του Σχήματος 6.4 και τον αντίστοιχο Πίνακα 6.3, μικρότερο `fit_time` και ως εκ τούτου πιο γρήγορος στη δημιουργία είναι ο GNB ενώ το χειρότερο χρόνο κατέχει ο MLP. Παρατηρείται ότι, τα μοντέλα XGB, RF και D-Tree που σύμφωνα με τις μετρικές θεωρούνται τα πιο αξιόπιστα να έχουν τις χειρότερες τιμές όσον αφορά το `fit_time`. Σε σχέση με το `score_time` μικρότερο χρόνο κατέχει ο D-Tree και χειρότερο ο KNN. Οι XGB και RF κατέχουν μια μεσαία τιμή στη χρονική κατάταξη. Επομένως, το αντίτιμο της αξιοπιστίας αποτελεί η χρονική διάρκεια των διεργασιών μέχρι την εξαγωγή των αποτελεσμάτων.

### 6.1.2 Σύνοψη

Στην υποενότητα αυτή εξετάστηκε το σύνολο των δεδομένων της Δυαδικής Κατηγοριοποίησης. Από τις μετρικές του `cross validation` όπως αναφέρθηκε παραπάνω, ο αλγόριθμος XGBoost είχε τις καλύτερες μέσες τιμές σε όλες τις μετρικές με ελάχιστη διαφορά όμως από τον RF, με τον τελευταίο να έχει καλύτερες τιμές στην τυπική απόκλιση από τον XGBoost. Στη χειρότερη θέση βρίσκεται ο αλγόριθμος Logistic Regression.

Όμως από τις μετρικές της ακρίβειας και του σφάλματος, μετά τη σχεδίαση του μοντέλου, οι αλγόριθμοι RF και D-Tree είχαν τις καλύτερες επιδόσεις, εκτοπίζοντας στην τρίτη θέση τον XGBoost. Με βάση αυτές τις παρατηρήσεις, ο RF και ο D-Tree φαίνεται να είναι καλύτερες επιλογές για το πρόβλημα της δυαδικής κατηγοριοποίησης.

Οι λόγοι για τους οποίους ο RF και ο D-Tree απέδωσαν καλύτερα από τον XGBoost αποδίδονται στους παρακάτω παράγοντες:

1. Μπορούν να χειριστούν αποτελεσματικά τόσο αριθμητικά όσο και κατηγορικά χαρακτηριστικά. Στην περίπτωση που το σύνολο δεδομένων περιέχει ένα μείγμα αυτών των τύπων χαρακτηριστικών, οι αλγόριθμοι αυτοί μπορούν να αποτυπώσουν τις σχέσεις και τα μοτίβα με μεγαλύτερη ακρίβεια. Από την άλλη πλευρά ο XGBoost έχει σχεδιαστεί κυρίως για αριθμητικά χαρακτηριστικά, και παρόλο που μπορεί να χειριστεί κατηγορικά με κατάλληλη κωδικοποίηση, ενδέχεται να μην έχει τόσο καλές επιδόσεις όσο οι RF και D-Tree.
2. Είναι ικανά να καταγράφουν μη γραμμικές σχέσεις μεταξύ των χαρακτηριστικών και της μεταβλητής στόχου. Επιπρόσθετα, μπορούν να χειριστούν σύν-

---

θετα όρια αποφάσεων και αλληλεπιδράσεις μεταξύ μεταβλητών. Ο XGBoost από την άλλη βασίζεται στην ενίσχυση της κλίσης και μοντελοποιεί κυρίως γραμμικούς συνδυασμούς αδύναμων μαθητών.

3. Είναι πιο ερμηνεύσιμοι σε σύγκριση με τον XGBoost. Τα δέντρα αποφάσεων παρέχουν διαφανείς κανόνες απόφαση, καθιστώντας ευκολότερη την κατανόηση και την ερμηνεία των προβλέψεων του μοντέλου. Ο RF συνδυάζει πολλαπλά δέντρα απόφασης και διατηρεί και αυτό με τη σειρά του το επίπεδο ερμηνευσιμότητας, καθώς και την αξιολόγηση της σημασίας των χαρακτηριστικών. Ο XGBoost από την άλλη είναι μια μέθοδος συνόλου και τείνει να είναι πιο πολύπλοκος και λιγότερο ερμηνεύσιμος.
4. Είναι πιο ανθεκτικά στο θόρυβο ή σε ακραία σημεία δεδομένων. Μπορούν να απομονώσουν και να αγνοήσουν τις ακραίες τιμές καθώς λαμβάνουν αποφάσεις με βάση την πλειοψηφία ή το μέσο όρο. Ενώ, ο XGBoost μπορεί να είναι πιο ευαίσθητος σε θορυβώδη σημεία δεδομένων, καθώς προσπαθεί να βελτιστοποιήσει τη συνάρτηση ελαχιστοποίησης των σφαλμάτων.
5. Η απόδοση του XGBoost εξαρτάται σε μεγάλο βαθμό από τον κατάλληλο συντονισμό των υπερπαραμέτρων. Επομένως, αν οι υπερπαραμέτροι του XGBoost δεν είχαν ρυθμιστεί με βέλτιστο τρόπο στα πειράματα, αυτό μπορεί να έχει ως αποτέλεσμα ελαφρώς κατώτερες επιδόσεις σε σύγκριση με τους RF και D-Tree. Οι τελευταίες μέθοδοι γενικά είναι λιγότερο ευαίσθητες στις ρυθμίσεις των υπερπαραμέτρων και μπορούν να αποδώσουν καλύτερα με τις προεπιλεγμένες ή τις βασικές ρυθμίσεις.

Σε δεύτερη φάση γίνεται μια σύγκριση μεταξύ του D-Tree και του RF βασισμένη στις αποδόσεις και στις μετρικές. Από τα αποτελέσματα του cross validation ο RF, κατά μέσο όρο, πέτυχε καλύτερα αποτελέσματα κατηγοριοποίησης σε σύγκριση με τον D-Tree. Ο RF αποτελεί εξέλιξη του D-Tree, συνδυάζει πολλαπλά δέντρα απόφασης, πράγμα που του επιτρέπει να καταγράφει πιο σύνθετες σχέσεις και να βελτιώνει τη συνολική απόδοση σε σύγκριση με ένα μεμονωμένο δέντρο απόφασης που χρησιμοποιεί ο D-Tree.

Τελικά, η επιλογή μεταξύ του RF και του D-Tree εξαρτάται από τις συγκεκριμένες απαιτήσεις και τους περιορισμούς του εκάστοτε προβλήματος. Πρέπει να

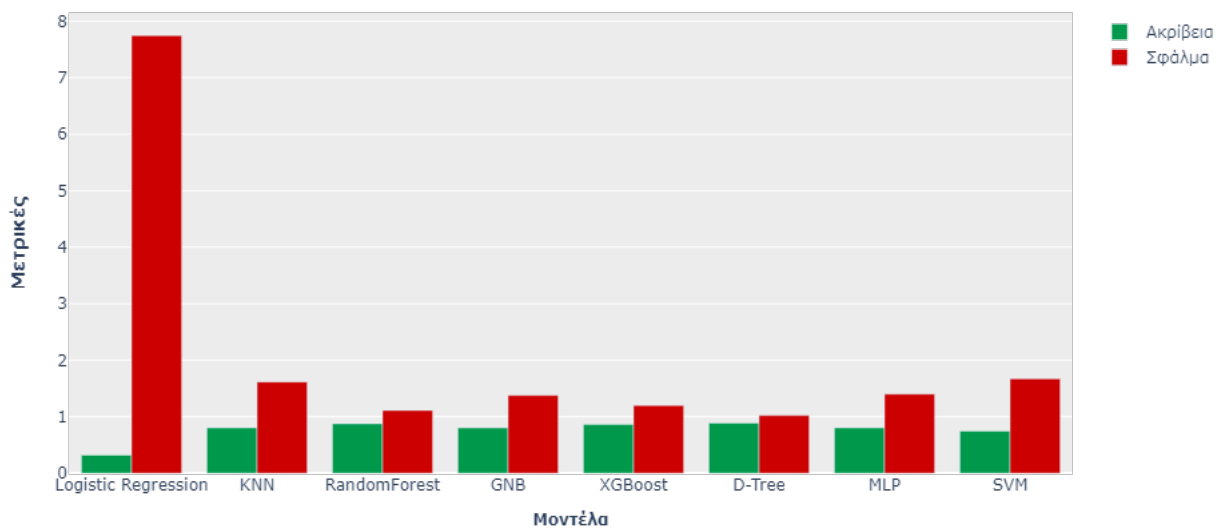
λαμβάνονται υπ' όψιν παράγοντες όπως το μέγεθος του συνόλου δεδομένων, οι διαθέσιμοι υπολογιστικοί πόροι, η σημασία της ερμηνευσιμότητας και η ανάγκη για προβλέψεις σε πραγματικό χρόνο κατά τη λήψη μιας απόφασης.

## 6.2 Σχεδιασμός μοντέλου πολυταξικής κατηγοριοποίησης

### 6.2.1 Ανάλυση

Όπως στη δυαδική, έτσι και στην πολυταξική κατηγοριοποίηση η ανάλυση ξεκινάει με την εξαγωγή του διαγράμματος ακρίβειας – σφάλματος.

Σχήμα 6.5: Ακρίβεια και σφάλμα στην πολυταξική κατηγοριοποίηση



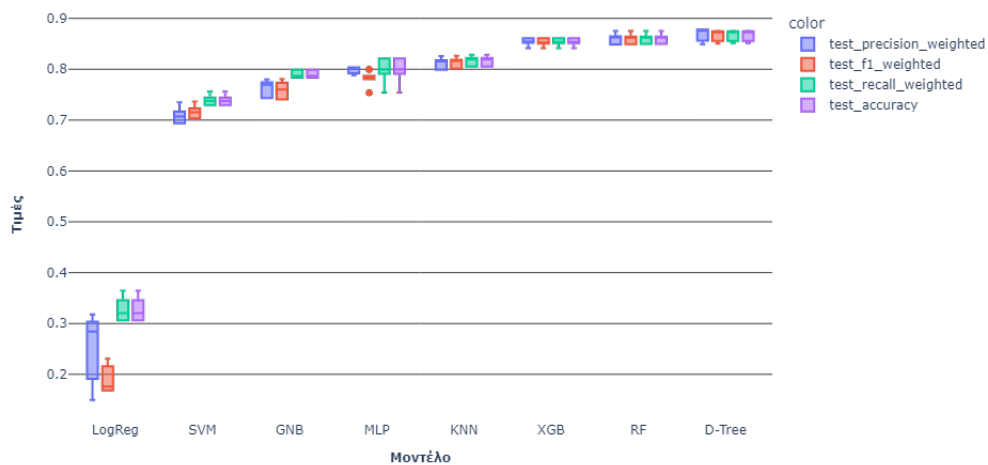
Μέθοδος	Ακρίβεια	Σφάλμα
LogReg	0.319135	7.745709
KNN	0.802924	1.615385
Random Forest	0.875397	1.111252
GNB	0.801653	1.376351
XGBoost	0.864589	1.198347
D-Tree	0.884933	1.020343
MLP	0.802924	1.399873
SVM	0.743166	1.670693

Πίνακας 6.4: Πίνακας τιμών ακρίβειας - σφάλματος στην πολυταξική κατηγοριοποίηση

Σύμφωνα με το διάγραμμα του Σχήματος 6.5 και τον Πίνακα 6.4 ο αλγόριθμος D-Tree έχει τις καλύτερες επιδόσεις σε σχέση με τους υπόλοιπους αφού έχει μεγαλύτερη ακρίβεια και μικρότερο σφάλμα. Αντιθέτως, ο Logistic Regression παρουσιάζει τα χειρότερα ποσοστά.

Σε επόμενο βήμα πραγματοποιείται η σύγκριση των μοντέλων μέσω των μετρικών κατηγοριοποίησης. Για την καλύτερη κατανόηση παρατίθεται η γραφική τους απεικόνιση, Σχήμα 6.6 αλλά και αριθμητική απεικόνιση, Πίνακας 6.5.

Σχήμα 6.6: Μετρικές κατάστασης cross validation στην πολυταξική κατηγοριοποίηση



Μοντέλο	Ανάκληση		Ακρίβεια Θετικών Προβλέψεων		F1		Ακρίβεια	
	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος
LogReg	0.020381	0.326754	0.062986	0.243479	0.023527	0.187428	0.020381	0.326754
SVM	0.010790	0.739805	0.015710	0.709566	0.012851	0.716725	0.010790	0.739805
GNB	0.006389	0.790424	0.014667	0.763393	0.014742	0.758616	0.006389	0.790424
D-Tree	0.009631	0.867171	0.011542	0.868509	0.010029	0.867083	0.009631	0.867171
KNN	0.008931	0.814731	0.010273	0.809584	0.009662	0.811470	0.008931	0.814731
MLP	0.020245	0.798463	0.005627	0.798088	0.012296	0.783864	0.020245	0.798463
RF	0.009929	0.860804	0.010604	0.860891	0.010174	0.860454	0.009929	0.860804
XGB	0.006843	0.853906	0.006952	0.853983	0.006977	0.853643	0.006843	0.853906

Πίνακας 6.5: Πίνακας μετρικών κατάστασης cross validation στην πολυταξική κατηγοριοποίηση

Παρατηρείται ότι από τη διαδικασία του cross validation καλύτερες επιδόσεις είχε ο αλγόριθμος D-Tree ενώ ο Logistic Regression τις χειρότερες. Από τον παραπάνω πίνακα παρατηρείτε επίσης ότι ενώ ο D-Tree έχει τις καλύτερες μέσες τιμές σε όλες τις μετρικές η τυπική του απόκλιση δεν συμφωνεί. Πιο συγκεκριμένα, σε μερικές περιπτώσεις καλύτερη τυπική απόκλιση έχουν οι αλγόριθμοι GNB, MLP και XGB. Επομένως, μπορεί ο D-Tree να έχει καλύτερες αποδόσεις αλλά οι αλγό-

ριθμοι που αναφέρθηκαν έχουν μια πιο σταθερή απόδοση στα διάφορα k-folds που εξετάστηκαν.

Για να αξιολογηθούν καλύτερα οι επιδόσεις των μοντέλων συχνά γίνεται χρήση του Πίνακα Σύγκυσης (Confusion Matrix). Η ανάλυση αυτού του πίνακα προσφέρει πληροφορίες για τη συμπεριφορά των μοντέλων και βοηθάει στον εντοπισμό περιοχών που επιφέρουν βελτίωση. Στο Σχήμα 6.7 παρουσιάζονται οι πίνακες για το κάθε μοντέλο με σκοπό να αποδοθεί οπτικά το μοντέλο με την καλύτερη απόδοση και να εξαχθούν πληροφορίες σχετικά με τις σχέσεις των τιμών.

Σχήμα 6.7: Confusion matrices

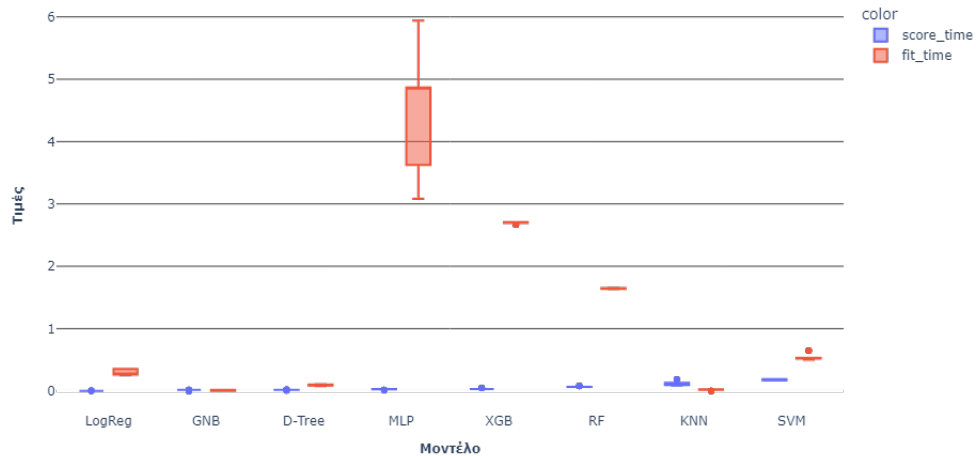


Αθροίζοντας τα στοιχεία της κύριας διαγωνίου παρατηρείται ότι ο αλγόριθμος D-Tree έκανε τις καλύτερες προβλέψεις, πιο συγκεκριμένα 1.392 και στη δεύτερη θέση είναι ο RF με 1.377 σωστές προβλέψεις.

Τέλος, πραγματοποιείται η χρονική ανάλυση των μοντέλων και παρατηρείται από το διάγραμμα του Σχήματος 6.8 και τον Πίνακα 6.6 ότι ο D-Tree έχει ένα

σχετικά καλό fit\_time, δηλαδή η διαδικασία προετοιμασίας του μοντέλου ήταν σχετικά αργή σε σύγκριση με άλλους αλγορίθμους όπως ο GNB που είχε το καλύτερο fit\_time. Παρόλα αυτά, ο D-Tree είχε το τρίτο καλύτερο score\_time και ως εκ τούτο ο χρόνος αναμονής αποτελεσμάτων είναι μικρός.

Σχήμα 6.8: Σύγκριση μοντέλων πολυταξικής κατηγοριοποίησης μέσω του χρόνου fit και score



Μοντέλο	Χρόνος Εκπαίδευσης		Χρόνος Αξιολόγησης	
	Απόκλιση	Αξιολόγηση	Απόκλιση	Αξιολόγηση
D-Tree	0.010889	0.096930	0.003035	0.015948
GNB	0.007026	0.004166	0.005962	0.013103
KNN	0.008712	0.015684	0.031930	0.114801
LogReg	0.041207	0.289255	0.001520	0.000668
MLP	0.998779	4.433978	0.006106	0.026530
RF	0.007654	1.641927	0.006461	0.066420
SVM	0.049209	0.540620	0.008455	0.180630
XGB	0.016447	2.695712	0.005920	0.033849

Πίνακας 6.6: Πίνακας χρονικών τιμών fit - score

### 6.2.2 Σύνοψη

Το παραπάνω πρόβλημα που αναλύθηκε ανήκει στην κατηγορία της Πολυταξικής Κατηγοριοποίησης. Τα βήματα που ακολουθήθηκαν ήταν τα ίδια με της Δυαδικής Κατηγοριοποίησης. Από τις μετρικές παρατηρήθηκε ότι ο αλγόριθμος D-Tree είχε καλύτερες αποδόσεις στη διαδικασία του cross validation αλλά και στην εξέταση των προβλέψεων μετά την εκπαίδευση του. Μερικοί λόγοι για τους οποίους ο αλγό-

---

ριθμος D-Tree αποδίδει καλύτερα αναλύθηκαν στην προηγούμενη σύνοψη της δυαδικής κατηγοριοποίησης. Συγκεκριμένα όμως για την πολυταξική κατηγοριοποίηση αναφέρονται παρακάτω οι εξής λόγοι:

1. Μπορούν να χειριστούν φυσικά προβλήματα κατηγοριοποίησης πολλαπλών κατηγοριών. Μπορούν να διαμερίσουν το χώρο των χαρακτηριστικών και να δημιουργήσουν κανόνες απόφασης για κάθε κλάση, καθιστώντας τα κατάλληλα για προβλήματα με περισσότερες από δύο κλάσεις
2. Τα δέντρα αποφάσεων είναι λιγότερο ευαίσθητα στην κλιμάκωση δεδομένων και δεν απαιτούν εκτεταμένη προεπεξεργασία των δεδομένων σε σύγκριση με άλλους αλγορίθμους

Επιπλέον, αξίζει να αναφερθεί ότι οι μέθοδοι, όπως ο Random Forest και το XGBoost, συχνά βασίζονται στις αρχές των δέντρων απόφασης και στοχεύουν στην περαιτέρω βελτίωση της απόδοσης συνδυάζοντας πολλαπλά δέντρα απόφασης. Αυτό θα μπορούσε να εξηγήσει γιατί τα RF και XGBoost είχαν επίσης καλή απόδοση στην ανάλυση.

## **6.3 Χειρισμός πολυταξικής κατηγοριοποίησης ως δυαδική**

### **6.3.1 Ανάλυση**

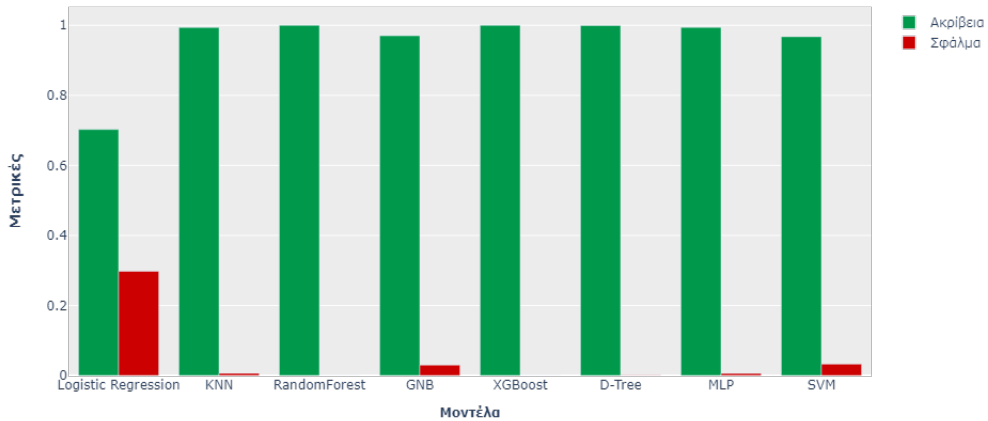
Η διαχείριση μιας πολυταξικής κατηγοριοποίησης ως δυαδική μπορεί να είναι μια χρήσιμη προσέγγιση σε πολλές περιπτώσεις. Πρώτον, προσφέρει μια πιο εύκολη σύγκριση και ερμηνεία των αποτελεσμάτων, ειδικά όταν πρόκειται για μεγάλο αριθμό κλάσεων. Η δυαδική κατηγοριοποίηση επικεντρώνεται στη σχέση μεταξύ δύο κλάσεων, απλοποιώντας την ανάλυση και διευκολύνοντας την εξαγωγή συμπερασμάτων. Επιπλέον, πολλοί αλγόριθμοι μηχανικής μάθησης αποδίδουν καλύτερα στη δυαδική κατηγοριοποίηση, καθώς έχουν σχεδιαστεί να δουλεύουν με δυαδικά δεδομένα. Ως εκ τούτου, στην υποενότητα αυτή το πρόβλημα ανάγεται σε δυαδική κατηγοριοποίηση με τη χρήση της τεχνικής “one – vs – all” ή “one – vs – rest”. Η προσέγγιση αυτή περιλαμβάνει την εκπαίδευση πολλαπλών δυαδικών κατηγοριοποιητών για την κάθε μια κλάση ξεχωριστά.

Συνολικά, υλοποιήθηκαν έξι δυαδικές αναλύσεις, όσες δηλαδή και οι κλάσεις του προβλήματος. Διαχωρίστηκαν τα δεδομένα σε δύο κατηγορίες. Σε κάθε επανάληψη θεωρήθηκε μία μόνο κλάση θετική και οι υπόλοιπες εντάχθηκαν στην κατηγορία της αρνητικής. Με την ίδια λογική που ακολουθήθηκε στις προηγούμενες αναλύσεις εξάχθηκε αρχικά το διάγραμμα ακρίβειας και σφάλματος. Έπειτα τα διαγράμματα των μετρικών από τα cross validation και τις χρονικές μετρήσεις.

### 1. Θετική Κλάση Χωρίς Σφάλμα - 0000

#### Ακρίβεια Σφάλμα

Σχήμα 6.9: Ακρίβεια - σφάλμα σε 0000



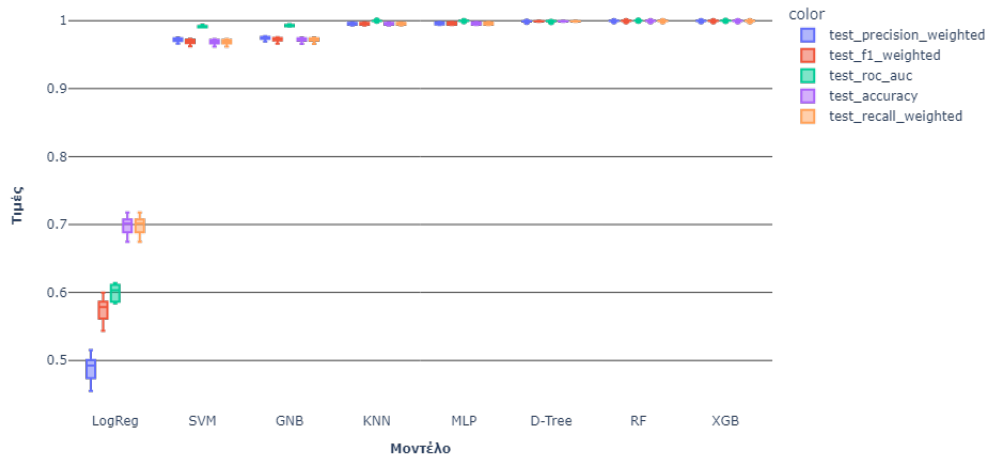
Μέθοδος	Ακρίβεια	Σφάλμα
LogReg	0.702479	0.297521
KNN	0.993643	0.297521
Random Forest	1.0	0.0
GNB	0.970121	0.297521
XGBoost	1.0	0.0
D-Tree	0.999364	0.000636
MLP	0.994278	0.005722
SVM	0.967578	0.032422

Πίνακας 6.7: Πίνακας τιμών ακρίβειας - σφάλματος σε 0000



## Cross Validation

Σχήμα 6.10: Σύγκριση μοντέλων μέσω μετρικών κατηγοριοποίησης 0000

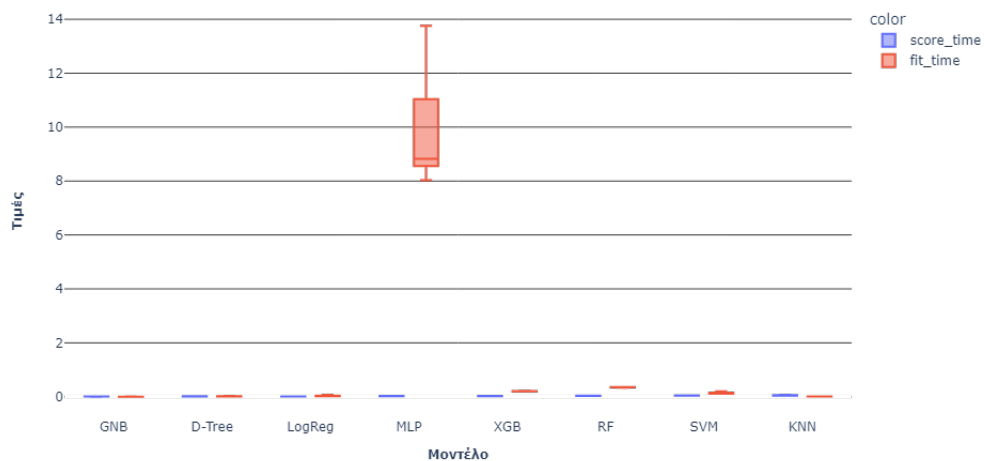


Μοντέλο	Ακρίβεια Θετικών Προβλέψεων		Ανάληψη		Roc_Auc		F1		Ακρίβεια	
	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος
LogReg	0.022156	0.487841	0.015928	0.698311	0.013542	0.599858	0.020772	0.574346	0.015928	0.698311
SVM	0.003451	0.971663	0.004275	0.968670	0.001619	0.991368	0.004154	0.969092	0.004275	0.968670
GNB	0.003024	0.974420	0.003670	0.972009	0.001306	0.992643	0.003576	0.972349	0.003670	0.972009
D-Tree	0.000563	0.999206	0.000563	0.999205	0.000704	0.999278	0.000563	0.999205	0.000563	0.999205
KNN	0.001843	0.995465	0.001900	0.995387	0.000257	0.999541	0.001892	0.995398	0.001900	0.995387
MLP	0.002107	0.996386	0.002149	0.996342	0.000302	0.999852	0.002142	0.996349	0.002149	0.996342
RF	0.000355	0.999841	0.000356	0.999841	0.000001	0.999999	0.000356	0.999841	0.000356	0.999841
XGB	0.000355	0.999841	0.000356	0.999841	0.000138	0.999938	0.000356	0.999841	0.000356	0.999841

Πίνακας 6.8: Πίνακας μετρικών κατάστασης cross validation σε 0000

## Χρόνοι Εκπαίδευσης και Αξιολόγησης

Σχήμα 6.11: Χρόνοι εκπαίδευσης και αξιολόγησης σε 0000



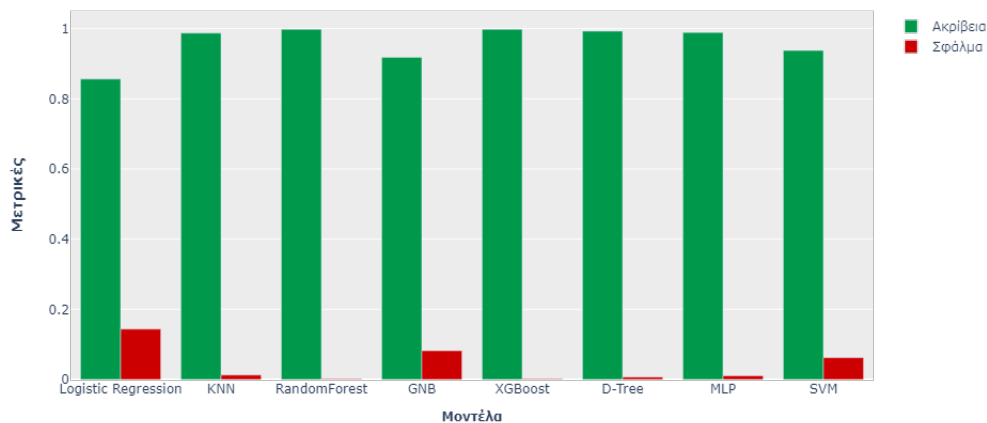
Μοντέλο	Χρόνος Εκπαίδευσης		Χρόνος Αξιολόγησης	
	Απόκλιση	Μέσος	Απόκλιση	Μέσος
D-Tree	0.008210	0.015606	0.007183	0.011381
GNB	0.000891	0.004197	0.004251	0.007303
KNN	0.002165	0.006595	0.011486	0.051036
LogReg	0.019692	0.035976	0.001978	0.010083
MLP	2.288266	9.894145	0.005355	0.026437
RF	0.008231	0.349716	0.001728	0.034567
SVM	0.033362	0.135638	0.005645	0.045058
XGB	0.007279	0.205191	0.002916	0.025779

Πίνακας 6.9: Πίνακας χρονικών τιμών εκπαίδευσης - αξιολόγησης σε 0000

## 2. Θετική Κλάση Σφάλμα 1011

### Ακρίβεια Σφάλμα

Σχήμα 6.12: Ακρίβεια - σφάλμα σε σφάλμα 1011

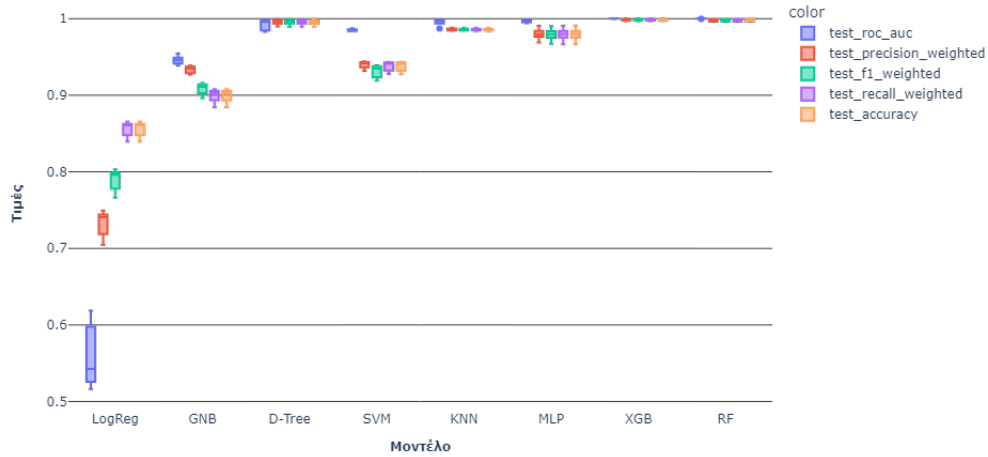


Μέθοδος	Ακρίβεια	Σφάλμα
LogReg	0.856325	0.143675
KNN	0.987285	0.012715
Random Forest	0.998093	0.001907
GNB	0.917991	0.082009
XGBoost	0.998093	0.001907
D-Tree	0.993007	0.006993
MLP	0.989193	0.010807
SVM	0.937699	0.062301

Πίνακας 6.10: Πίνακας τιμών ακρίβειας - σφάλματος σε σφάλμα 1011

## Cross Validation

Σχήμα 6.13: Σύγκριση μοντέλων μέσω μετρικών κατηγοριοποίησης σε σφάλμα 1011

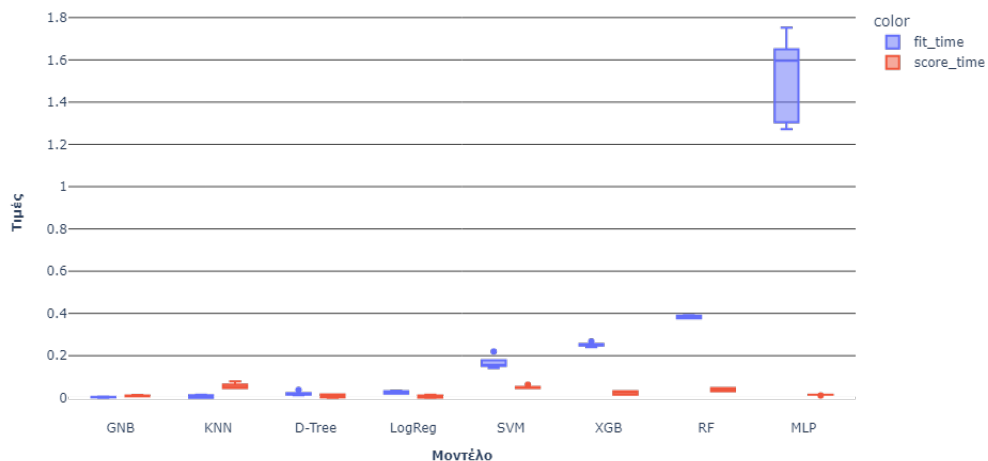


Μοντέλο	Ανάκλιση		Ακρίβεια Θετικών Προβλέψεων		F1		Roc_Auc		Ακρίβεια	
	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος
LogReg	0.010605	0.855600	0.018077	0.732142	0.015028	0.789047	0.043382	0.559443	0.010605	0.855600
SVM	0.006810	0.937819	0.005055	0.939834	0.008611	0.931179	0.001441	0.985121	0.006810	0.937819
GNB	0.009271	0.899490	0.004876	0.933274	0.007919	0.908355	0.005702	0.945815	0.009271	0.899490
D-Tree	0.003669	0.995706	0.003640	0.995737	0.003657	0.995712	0.007341	0.992021	0.003669	0.995706
KNN	0.001531	0.985846	0.001500	0.985881	0.001486	0.985836	0.004732	0.995550	0.001531	0.985846
MLP	0.008549	0.979485	0.007628	0.980068	0.008224	0.979349	0.001713	0.996485	0.008549	0.979485
RF	0.001550	0.998092	0.001542	0.998098	0.001559	0.998084	0.000122	0.999939	0.001550	0.998092
XGB	0.001180	0.998568	0.001176	0.998571	0.001186	0.998564	0.000139	0.999907	0.001180	0.998568

Πίνακας 6.11: Πίνακας μετρικών κατάστασης cross validation σε σφάλμα 1011

## Χρόνοι Εκπαίδευσης και Αξιολόγησης

Σχήμα 6.14: Χρόνοι εκπαίδευσης και αξιολόγησης σε σφάλμα 1011



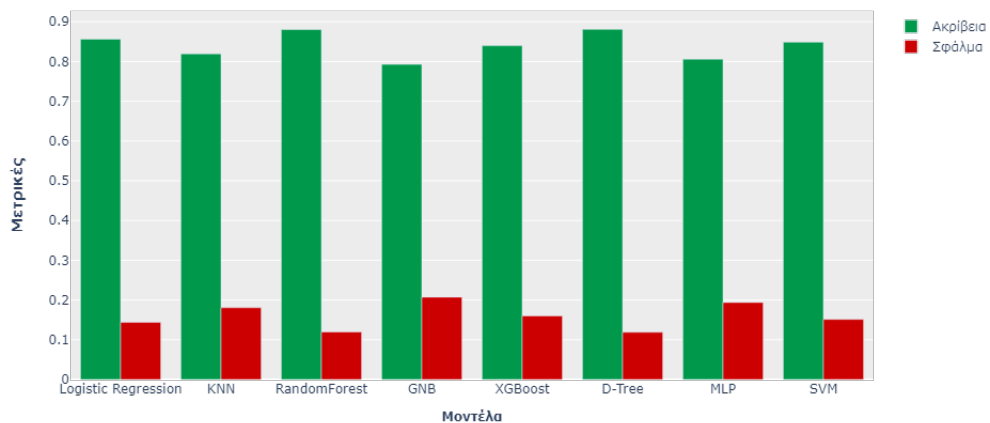
Μοντέλο	Χρόνος Εκπαίδευσης		Χρόνος Αξιολόγησης	
	Απόκλιση	Μέσος	Απόκλιση	Μέσος
D-Tree	0.010291	0.019722	0.007947	0.010351
GNB	0.002143	0.003179	0.002307	0.009823
KNN	0.007674	0.005523	0.013718	0.055565
LogReg	0.006670	0.027695	0.006543	0.007022
MLP	0.207394	1.510242	0.001855	0.014792
RF	0.006957	0.382589	0.008635	0.037699
SVM	0.030333	0.167550	0.006696	0.050584
XGB	0.009590	0.252104	0.008544	0.024980

Πίνακας 6.12: Πίνακας χρονικών τιμών εκπαίδευσης - αξιολόγησης σε 1011

### 3. Θετική Κλάση 1111

#### Ακρίβεια Σφάλμα

Σχήμα 6.15: Ακρίβεια - σφάλμα σε 1111

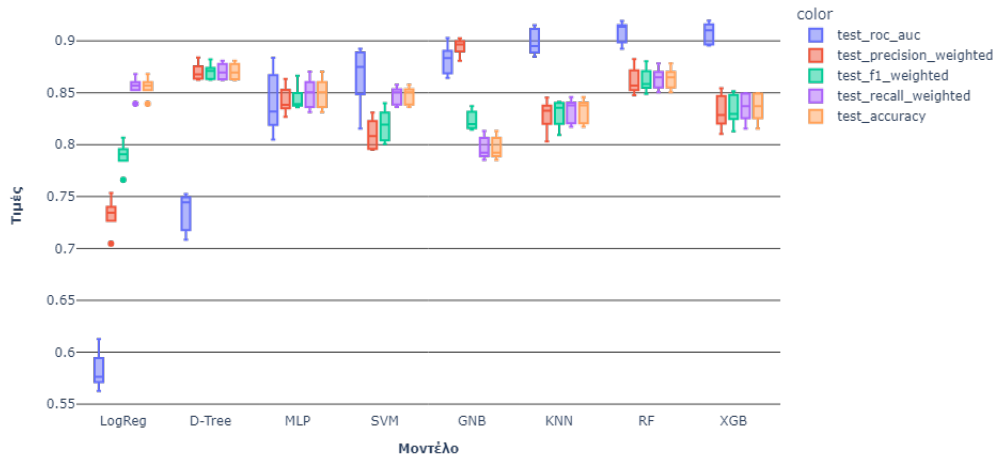


Μέθοδος	Ακρίβεια	Σφάλμα
LogReg	0.856325	0.143675
KNN	0.819453	0.180547
Random Forest	0.880483	0.119517
GNB	0.792753	0.207247
XGBoost	0.839797	0.160203
D-Tree	0.881119	0.118881
MLP	0.806110	0.193897
SVM	0.848697	0.151303

Πίνακας 6.13: Πίνακας τιμών ακρίβειας - σφάλματος σε σφάλμα 1111

#### Cross Validation

Σχήμα 6.16: Σύγκριση μοντέλων μέσω μετρικών κατηγοριοποίησης σε σφάλμα 1111

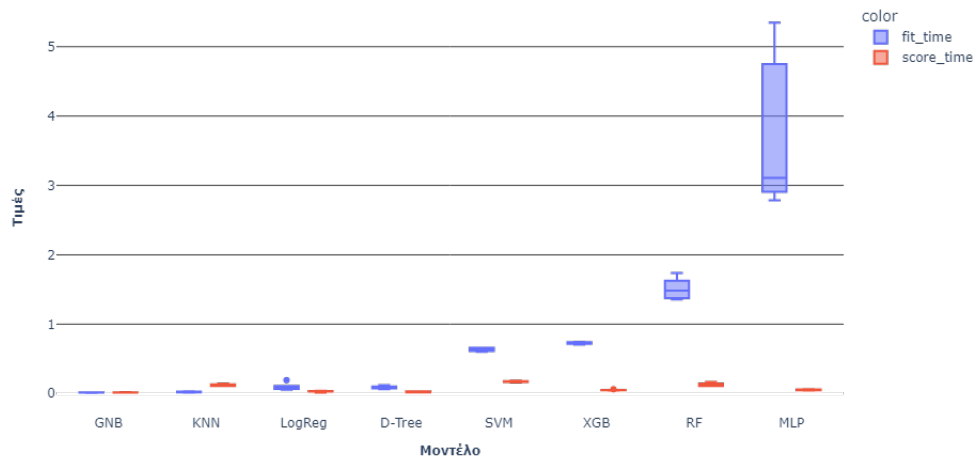


Μοντέλο	Ανάκλιση		Ακρίβεια		Roc_Auc		Ακρίβεια Θετικών Προβλέψεων		F1	
	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος
LogReg	0.010289	0.855757	0.010289	0.855757	0.018982	0.582792	0.017550	0.732406	0.014583	0.789269
SVM	0.009047	0.847169	0.009047	0.847169	0.030902	0.866049	0.015420	0.810034	0.016184	0.818634
GNB	0.011393	0.797074	0.011393	0.797074	0.015019	0.881429	0.008299	0.894311	0.009613	0.823613
D-Tree	0.008125	0.870546	0.008125	0.870546	0.019340	0.734901	0.008603	0.870129	0.007895	0.870089
KNN	0.012147	0.832219	0.012147	0.832219	0.012983	0.899059	0.015738	0.828467	0.013445	0.829957
MLP	0.015476	0.849394	0.015476	0.849394	0.031315	0.841183	0.013763	0.843143	0.012423	0.844790
RF	0.010594	0.863549	0.010594	0.863549	0.011004	0.907968	0.013712	0.861837	0.011976	0.862438
XGB	0.014334	0.836038	0.014334	0.836038	0.010687	0.907397	0.017324	0.832298	0.015544	0.833923

Πίνακας 6.14: Πίνακας μετρικών κατάστασης cross validation σε σφάλμα 1111

### Χρόνοι Εκπαίδευσης και Αξιολόγησης

Σχήμα 6.17: Χρόνοι εκπαίδευσης και αξιολόγησης σε σφάλμα 1111



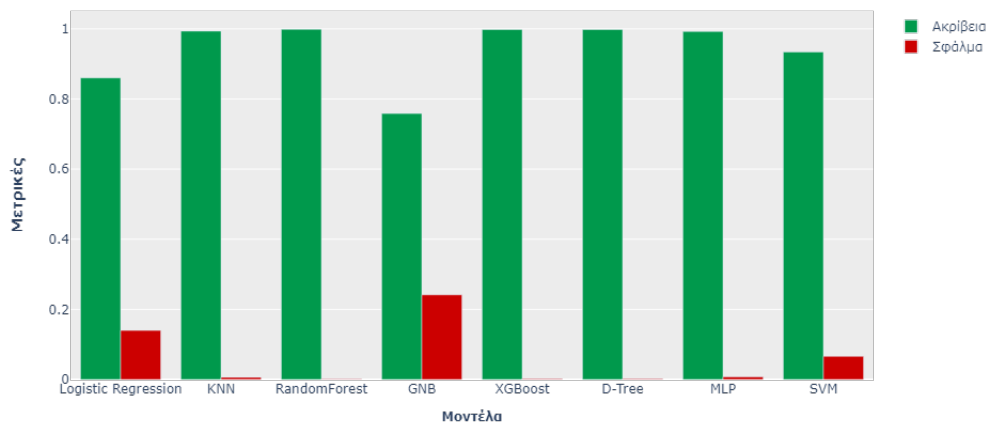
Μοντέλο	Χρόνος Αξιολόγησης		Χρόνος Εκπαίδευσης	
	Απόκλιση	Μέσος	Απόκλιση	Μέσος
D-Tree	0.008700	0.025157	0.020465	0.088169
GNB	0.003060	0.015736	0.003896	0.014985
KNN	0.013901	0.123506	0.008178	0.020096
LogReg	0.010432	0.029743	0.055637	0.093754
MLP	0.007206	0.052298	1.137143	3.748585
RF	0.024294	0.125310	0.157327	1.510297
SVM	0.010070	0.172096	0.027523	0.629727
XGB	0.006442	0.050224	0.014242	0.728225

Πίνακας 6.15: Πίνακας χρονικών τιμών εκπαίδευσης - αξιολόγησης

#### 4. Θετική Κλάση 1001

##### Ακρίβεια Σφάλμα

Σχήμα 6.18: Ακρίβεια - σφάλμα σε σφάλμα 1001

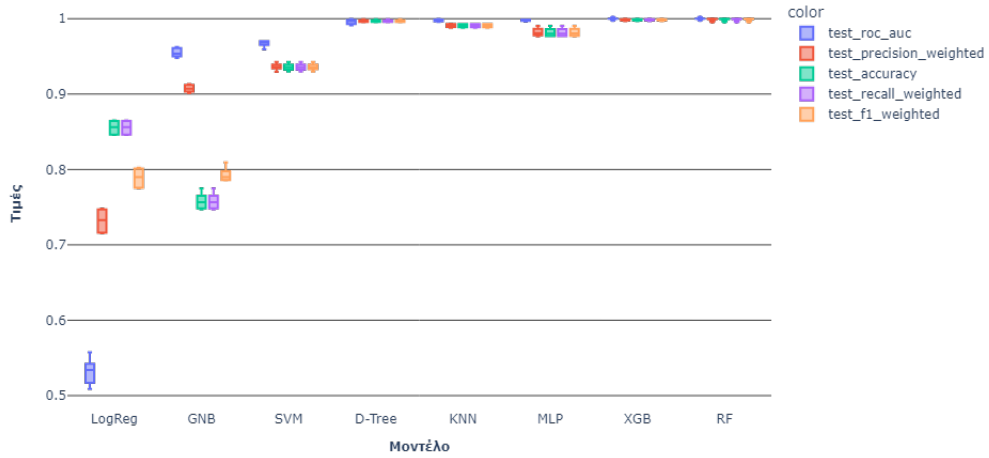


Μέθοδος	Ακρίβεια	Σφάλμα
LogReg	0.860140	0.139860
KNN	0.993643	0.006357
Random Forest	0.998729	0.001271
GNB	0.758423	0.241577
XGBoost	0.998093	0.001907
D-Tree	0.998093	0.001907
MLP	0.992371	0.007629
SVM	0.933884	0.066116

Πίνακας 6.16: Πίνακας τιμών ακρίβειας - σφάλματος σε σφάλμα 1001

##### Cross Validation

Σχήμα 6.19: Σύγκριση μοντέλων μέσω μετρικών κατηγοριοποίησης σε σφάλμα 1001

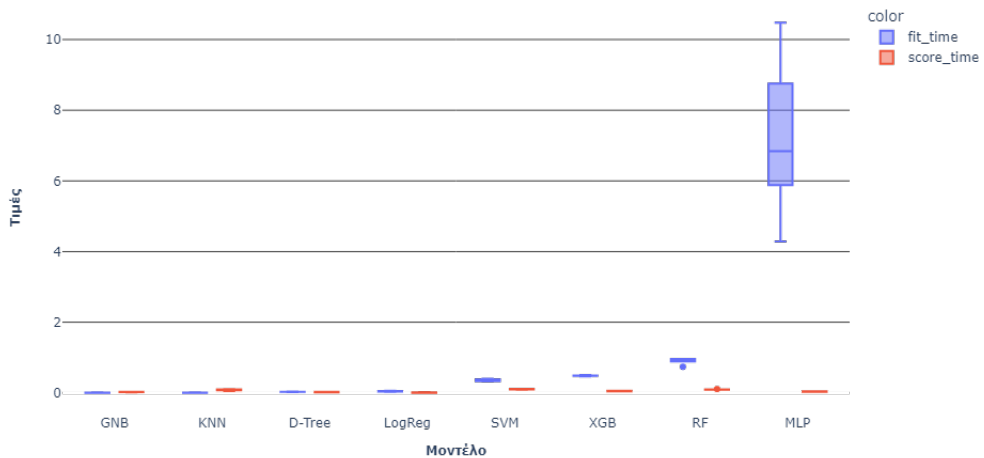


Μοντέλο	Roc_Auc		Ακρίβεια		Ακρίβεια Θετικών Προβλέψεων		Ανάκληση		F1	
	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος
LogReg	0.018365	0.531401	0.009162	0.855439	0.015670	0.731843	0.009162	0.855439	0.012999	0.788811
SVM	0.004476	0.967029	0.004823	0.935910	0.004708	0.936254	0.004823	0.935910	0.004749	0.936040
GNB	0.006191	0.954961	0.011289	0.757791	0.005060	0.907459	0.011289	0.757791	0.009642	0.792961
D-Tree	0.003413	0.995806	0.001591	0.996819	0.001539	0.996877	0.001591	0.996819	0.001578	0.996830
KNN	0.001555	0.997626	0.002288	0.990936	0.002277	0.991206	0.002288	0.990936	0.002283	0.991002
MLP	0.001220	0.997848	0.005803	0.981869	0.005634	0.982343	0.005803	0.981869	0.005803	0.981967
RF	0.000027	0.999988	0.001067	0.999523	0.001064	0.999524	0.001067	0.999523	0.001066	0.999523
XGB	0.000052	0.999977	0.001258	0.998409	0.001244	0.998426	0.001258	0.998409	0.001254	0.998413

Πίνακας 6.17: Πίνακας μετρικών κατάστασης cross validation σε σφάλμα 1001

### Χρόνοι Εκπαίδευση Αξιολόγησης

Σχήμα 6.20: Χρόνοι εκπαίδευσης αξιολόγησης σε σφάλμα 1001



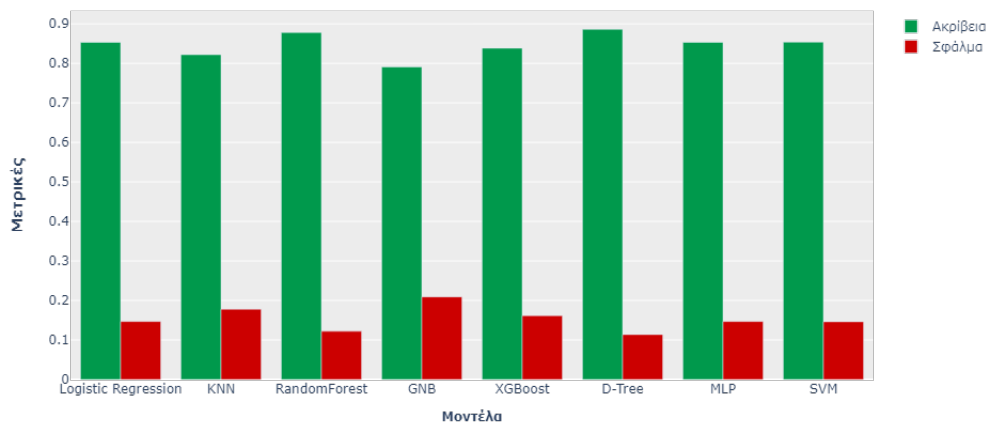
Μοντέλο	Χρόνος Εκπαίδευσης		Χρόνος Αξιολόγησης	
	Απόκλιση	Μέσος	Απόκλιση	Μέσος
D-Tree	0.005955	0.030301	0.005605	0.022466
GNB	0.006303	0.004481	0.005933	0.024249
KNN	0.005776	0.005823	0.019837	0.083540
LogReg	0.013637	0.045955	0.009247	0.008855
MLP	2.286890	7.239083	0.007409	0.040530
RF	0.096848	0.913192	0.007252	0.097488
SVM	0.032593	0.355497	0.007799	0.105465
XGB	0.011302	0.486752	0.008679	0.056980

Πίνακας 6.18: Πίνακας χρονικών τιμών εκπαίδευσης - αξιολόγησης

## 5. Θετική Κλάση 0111

### Ακρίβεια Σφάλμα

Σχήμα 6.21: Ακρίβεια - σφάλμα σε σφάλμα 0111



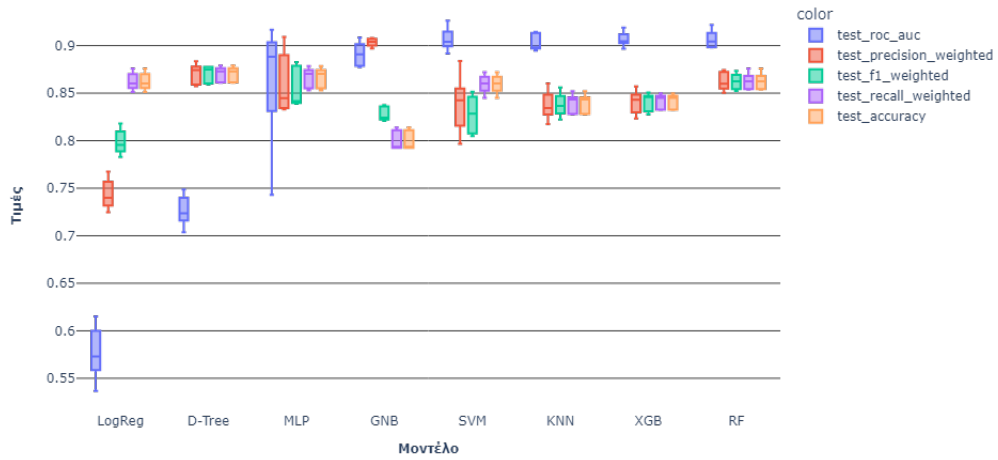
Μέθοδος	Ακρίβεια	Σφάλμα
LogReg	0.853147	0.146853
KNN	0.821996	0.178004
Random Forest	0.877940	0.122060
GNB	0.790845	0.209154
XGBoost	0.838525	0.161475
D-Tree	0.886205	0.113795
MLP	0.853147	0.146853
SVM	0.853783	0.146217

Πίνακας 6.19: Πίνακας τιμών ακρίβειας - σφάλματος σε σφάλμα 0111

### Cross Validation



Σχήμα 6.22: Σύγκριση μοντέλων μέσω μετρικών κατηγοριοποίησης σε σφάλμα 0111

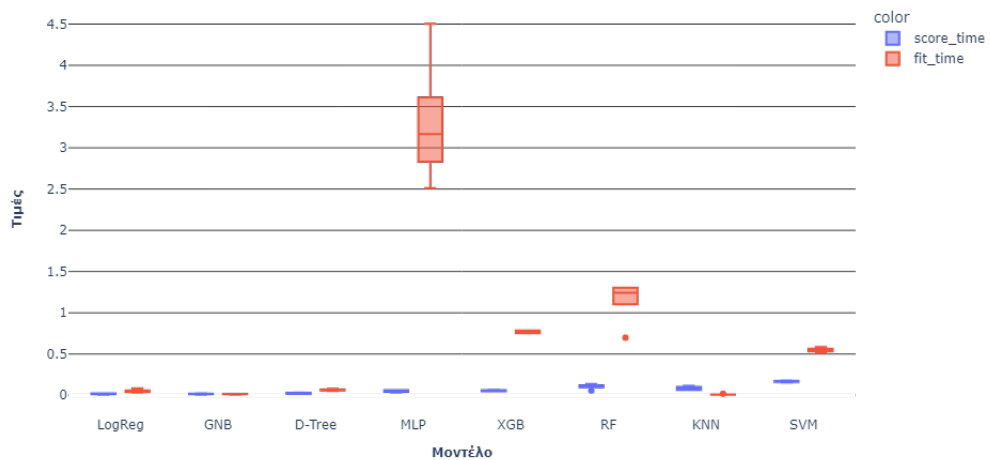


Μοντέλο	Ακρίβεια		Ανάκλιση		F1		Ακρίβεια Θετικών Προβλέψεων		Roc_Auc	
	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος
LogReg	0.009690	0.862435	0.009690	0.862435	0.013801	0.798757	0.016741	0.743870	0.029752	0.577104
SVM	0.010295	0.859573	0.010295	0.859573	0.020914	0.827551	0.032124	0.838102	0.012797	0.906931
GNB	0.010601	0.800414	0.010601	0.800414	0.007748	0.828163	0.004460	0.903407	0.013239	0.891045
D-Tree	0.008235	0.869911	0.008235	0.869911	0.009251	0.869886	0.011205	0.870211	0.017124	0.726587
KNN	0.010725	0.838898	0.010725	0.838898	0.012816	0.837860	0.015957	0.837500	0.009287	0.903625
MLP	0.011054	0.866093	0.011054	0.866093	0.021983	0.856043	0.033818	0.861050	0.069320	0.861578
RF	0.009132	0.862436	0.009132	0.862436	0.008782	0.862188	0.010140	0.862452	0.009833	0.906680
XGB	0.008113	0.841284	0.008113	0.841284	0.009954	0.840464	0.012989	0.840137	0.008184	0.906822

Πίνακας 6.20: Πίνακας μετρικών κατάστασης cross validation σε σφάλμα 0111

### Χρόνοι Εκπαίδευσης και Αξιολόγησης

Σχήμα 6.23: Χρόνοι εκπαίδευσης και αξιολόγησης σε σφάλμα 0111



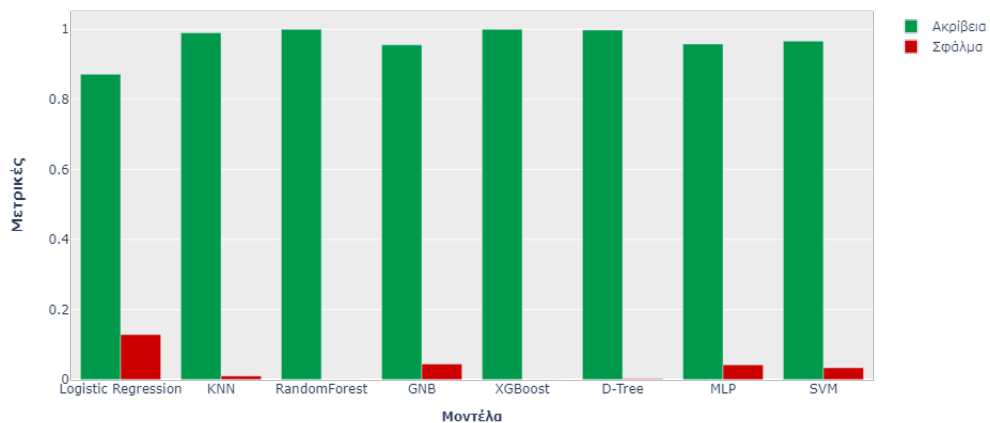
Μοντελο	Χρόνος Εκπαίδευσης		Χρόνος Αξιολόγησης	
	Απόκλιση	Μέσος	Απόκλιση	Μέσος
D-Tree	0.009759	0.063437	0.007552	0.023316
GNB	0.003277	0.013670	0.005530	0.016111
KNN	0.004652	0.009877	0.021200	0.082314
LogReg	0.021055	0.049613	0.007384	0.016915
MLP	0.744379	3.286591	0.013170	0.052337
RF	0.258527	1.157698	0.029836	0.106594
SVM	0.025431	0.547992	0.008792	0.168057
XGB	0.016484	0.770109	0.008382	0.054357

Πίνακας 6.21: Πίνακας χρονικών τιμών εκπαίδευσης - αξιολόγησης

## 6. Θετική Κλάση 0110

### Ακρίβεια Σφάλμα

Σχήμα 6.24: Ακρίβεια - σφάλμα σε σφάλμα 0110

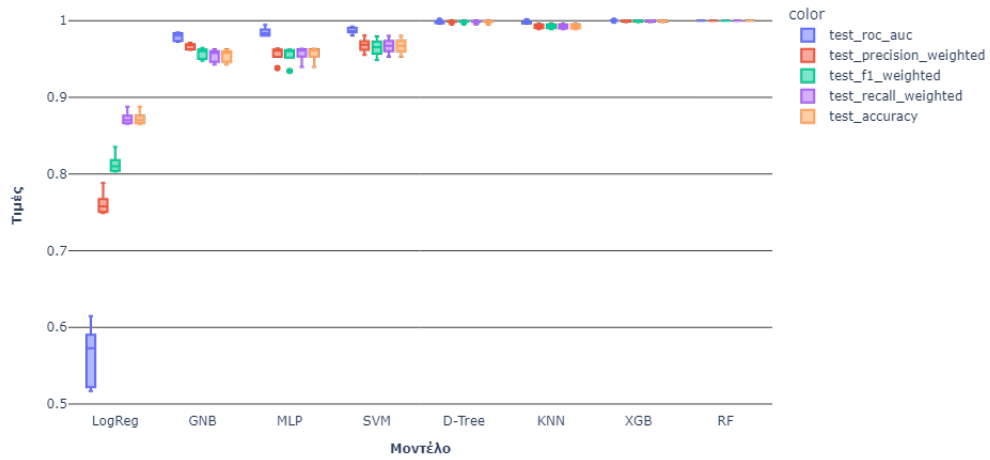


Μέθοδος	Ακρίβεια	Σφάλμα
LogReg	0.871583	0.128417
KNN	0.989828	0.010172
Random Forest	1.0	0.0
GNB	0.955499	0.0445010
XGBoost	1.0	0.0
D-Tree	0.998093	0.001907
MLP	0.958042	0.041958
SVM	0.966306	0.033694

Πίνακας 6.22: Πίνακας τιμών ακρίβειας - σφάλματος σε σφάλμα 0111

### Cross Validation

Σχήμα 6.25: Σύγκριση μοντέλων μέσω μετρικών κατηγοριοποίησης σε σφάλμα 0110

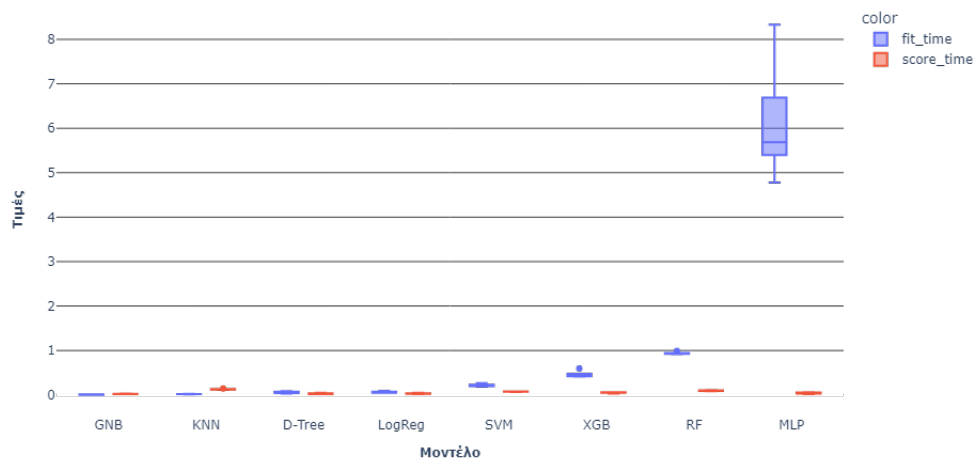


Μοντέλο	F1		Roc_Auc		Ακρίβεια		Ανάκληση		Ακρίβεια Θετικών Προβλέψεων	
	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος	Απόκλιση	Μέσος
LogReg	0.012897	0.813049	0.041237	0.561970	0.009004	0.872457	0.009004	0.872457	0.015796	0.761246
SVM	0.011447	0.964732	0.004331	0.987208	0.010059	0.966920	0.010059	0.966920	0.009193	0.967940
GNB	0.007269	0.957035	0.005498	0.978974	0.008486	0.954038	0.008486	0.954038	0.003910	0.966536
D-Tree	0.000870	0.999046	0.001501	0.997395	0.000872	0.999046	0.000872	0.999046	0.000868	0.999048
KNN	0.002890	0.992350	0.001205	0.997184	0.002901	0.992366	0.002901	0.992366	0.002890	0.992375
MLP	0.011914	0.955407	0.005744	0.985090	0.009911	0.956904	0.009911	0.956904	0.010520	0.956632
RF	0.000000	1.000000	0.000000	1.000000	0.000000	1.000000	0.000000	1.000000	0.000000	1.000000
XGB	0.000714	0.999522	0.000003	0.999999	0.000711	0.999523	0.000711	0.999523	0.000710	0.999524

Πίνακας 6.23: Πίνακας μετρικών κατάστασης cross validation σε σφάλμα 0110

### Χρόνοι Εκπαίδευσης και Αξιολόγησης

Σχήμα 6.26: Χρόνοι εκπαίδευσης και αξιολόγησης σε σφάλμα 0110



Model	Χρόνος Αξιολόγησης		Χρόνος Εκπαίδευσης	
	std	mean	std	mean
D-Tree	0.012970	0.026161	0.020629	0.055274
GNB	0.005427	0.021400	0.005343	0.006182
KNN	0.010710	0.127384	0.003898	0.017642
LogReg	0.011180	0.029105	0.017223	0.060069
MLP	0.017154	0.038195	1.335429	6.108652
RF	0.007547	0.096525	0.028628	0.939011
SVM	0.005414	0.076358	0.030278	0.215405
XGB	0.011757	0.051242	0.076990	0.457759

Πίνακας 6.24: Πίνακας Χρονικών Τιμών Εκπαίδευσης - Αξιολόγησης

Από τις παραπάνω περιπτώσεις γίνεται αντιληπτό ότι υπάρχει μια γενική αύξηση τόσο στην ακρίβεια και στο σφάλμα όσο και στις μετρικές απόδοσης και χρόνων όλων των αλγορίθμων. Αυτό μπορεί να ισχύει για τους παρακάτω λόγους.

1. Σε πολυταξικά προβλήματα, στην περίπτωση που υπάρχει ανισορροπία κλάσεων μπορεί να επηρεαστεί η ικανότητα του μοντέλου να εκπαιδευτεί και να γενικεύει καλά. Με τη μετατροπή του σε δυαδικό πρόβλημα μπορεί να επέλθει εξισορρόπηση των κλάσεων γεγονός που μπορεί να βελτιώσει την απόδοση των αλγορίθμων. Αυτή η ανισορροπία ορισμένες φορές οδηγεί και σε μεροληπτική πρόβλεψη του μοντέλου.
2. Στα δυαδικά προβλήματα απλοποιείται η πολυπλοκότητα των αλγορίθμων. Επίσης, μπορεί να είναι πιο εύκολο για τους αλγορίθμους να εκπαιδευτούν και να μάθουν τα υποκείμενα πρότυπα έτσι ώστε να κάνουν ακριβείς προβλέψεις.
3. Σε ορισμένες περιπτώσεις, τα σύνολα δεδομένων πολλαπλών κατηγοριών μπορεί να περιέχουν θόρυβο ετικετών δηλαδή δεδομένα με λανθασμένη ή διαφορετική ετικέτα. Με τη μετατροπή σε δυαδικό πρόβλημα μπορεί να μειωθεί το αποτέλεσμα του θορύβου αφού η κατηγοριοποίηση έγινε πιο απλή και αυτό μπορεί να οδηγήσει σε βελτίωση απόδοσης.
4. Στα δυαδικά προβλήματα κατηγοριοποίησης υπάρχει βελτιωμένη ερμηνευσιμότητα με αποτέλεσμα να υπάρχει συχνά σαφέστερη διάκριση μεταξύ των δύο κλάσεων. Αυτό μπορεί να είναι ιδιαίτερα σημαντικό για αλγορίθμους που βασίζονται σε δέντρα αποφάσεων, όπως ο D-Tree.

### 6.3.2 Σύνοψη

Στο προηγούμενο πρόβλημα η πολυταξική κατηγοριοποίηση μετατράπηκε σε δυαδική. Παρακάτω παρουσιάζονται κάποιοι συγκριτικοί πίνακες των μετρικών της πολυταξικής κατηγοριοποίησης και των μέσων τιμών της δυαδικής.

- Ακρίβεια Σφάλμα

Μέθοδος	Ακρίβεια Πολυ- ταξικής	Ακρίβεια Δυα- δικής
LogReg	0.319135	0.833333
KNN	0.802924	0.934308
Random Forest	0.875397	0.959207
GNB	0.801653	0.864272
XGBoost	0.864589	0.945751
D-Tree	0.884933	0.959313
MLP	0.802924	0.932190
SVM	0.743166	0.917991

Πίνακας 6.25: Συγκριτικός πίνακας τιμών ακρίβειας σε πολυταξική και δυαδική κατηγοριοποίηση

Συγκρίνοντας τους πίνακες ακρίβειας και σφάλματος γίνεται αντιληπτό ότι υπάρχει μια σημαντική διαφορά μεταξύ των τιμών. Πιο συγκεκριμένα, στο πρόβλημα της δυαδικής τις καλύτερες αποδόσεις συνεχίζει να έχει ο D-Tree αλλά με ποσοστό ακρίβειας 95% σε σχέση με πριν που είχε 88%. Έπειτα, ο RF και ο XGBoost συνεχίζουν να έχουν τις επόμενες καλύτερες τιμές. Σημαντική αύξηση είχε και ο αλγόριθμος Logistic Regression που από 32% έφτασε το 83% αλλά παρόλα αυτά συνεχίζει να έχει τη χειρότερη απόδοση σε σχέση με τους υπόλοιπους, Πίνακας 6.25.

Μέθοδος	Σφάλμα Πολυ- ταξικής	Σφάλμα Δυαδι- κής
LogReg	7.745709	0.166667
KNN	1.615385	0.114219
Random Forest	1.111252	0.040793
GNB	1.376351	0.180335
XGBoost	1.198347	0.054249
D-Tree	1.020343	0.040687
MLP	1.399873	0.067811
SVM	1.670693	0.082009

Πίνακας 6.26: Συγκριτικός πίνακας τιμών σφάλματος σε πολυταξική και δυαδική κατηγοριοποίηση

Όσον αφορά τα ποσοστά σφάλματος πάλι παρατηρείται μια σημαντική μείωση σε όλους τους αλγορίθμους. Πάλι ο D-Tree έχει το μικρότερο σφάλμα με ποσοστό 0.4% σε σχέση με πριν που ήταν 1.02%, Πίνακας 6.26.

- Cross Validation

Μοντέλο	Ανάκληση		Ακρίβεια Θετικών Προβλέψεων		F1		Ακρίβεια	
	Μέσος_Πολ.	Μέσος_Δυαδ.	Μέσος_Πολ.	Μέσος_Δυαδ.	Μέσος_Πολ.	Μέσος_Δυαδ.	Μέσος_Πολ.	Μέσος_Δυαδ.
LogReg	0.326754	0.833333	0.243479	0.698225	0.187428	0.758880	0.326754	0.825508
SVM	0.739805	0.919343	0.709566	0.910638	0.716725	0.907871	0.739805	0.9193435
GNB	0.790424	0.863469	0.763393	0.929901	0.758616	0.880413	0.790424	0.863469
D-Tree	0.867171	0.955206	0.868509	0.955201	0.867083	0.955128	0.867171	0.955206
KNN	0.814731	0.939275	0.809584	0.938482	0.811470	0.938734	0.814731	0.939275
MLP	0.798463	0.938348	0.798088	0.936604	0.783864	0.935651	0.798463	0.938348
RF	0.860804	0.953907	0.860891	0.953625	0.860454	0.953679	0.860804	0.953907
XGB	0.853906	0.945611	0.853983	0.944799	0.853643	0.945121	0.853906	0.945611

Πίνακας 6.27: Πίνακας μετρικών κατάστασης cross validation σε δυαδική και πολυταξική κατ

Και στη φάση του cross validation παρατηρείται μια γενική αύξηση των ποσοστών σε όλες τις μετρικές για όλα τα μοντέλα. Η κατάταξη τους είναι η ίδια, δηλαδή στην πρώτη θέση βρίσκεται ο D-Tree μετά ο RF και ο XGBoost και στην χειρότερη ο Logistic Regression. Υπάρχει όμως μια διαφορά μεταξύ του GNB και του SVM. Ενώ στην πολυταξική ο GNB έχει καλύτερα ποσοστά από τον SVM στις μετρικές ανάκλησης, f1 σκόρ και ακρίβειας, στη δυαδική ο SVM κατέχει καλύτερη θέση, Πίνακας 6.27.

Στην πολυταξική κατηγοριοποίηση ο GNB δίνει καλύτερα ποσοστά καθώς υποθέτει ότι υπάρχει μια ανεξαρτησία μεταξύ των χαρακτηριστικών. Από την άλλη ο SVM δεν κάνει καμία υπόθεση και μπορεί να χειριστεί πιο πολύπλοκες σχέσεις μεταξύ των χαρακτηριστικών. Ωστόσο, αυτό μπορεί να οδηγήσει σε μειωμένη απόδοση όταν τα δεδομένα παρουσιάζουν ανεξαρτησία στα χαρακτηριστικά.

Στη δυαδική κατηγοριοποίηση ο SVM έχει καλύτερες επιδόσεις καθώς στοχεύει στην εύρεση του βέλτιστου ορίου απόφασης που διαχωρίζει σε μέγιστο βαθμό τις δύο κλάσεις. Επίσης, ο SVM μπορεί να χειριστεί καλύτερα τις ανισοροπίες μεταξύ των κλάσεων μετριάζοντας το όριο απόφασης για να μετριάσει τη μεροληψία προς την κλάση με τα περισσότερα δεδομένα.

- Χρονικές Μετρικές

Από τις χρονικές μετρικές παρατηρείται ότι υπάρχουν αυξομειώσεις στα μοντέλα. Πιο συγκεκριμένα στο fit\_time οι αλγόριθμοι GNB και MLP παρουσία-

Μοντέλο	Χρόνος Εκπαίδευσης		Χρόνος Αξιολόγησης	
	Μέσος_Πολ.	Μέσος_Δυαδ.	Μέσος_Πολ.	Μέσος_Δυαδ.
D-Tree	0.096930	0.045418	0.015948	0.019805
GNB	0.004166	0.007782	0.013103	0.015770
KNN	0.015684	0.010926	0.114801	0.087224
LogReg	0.289255	0.052177	0.000668	0.016954
MLP	4.433978	5.297883	0.026530	0.037432
RF	1.641927	0.875417	0.066420	0.0830301
SVM	0.540620	0.341968	0.180630	0.102936
XGB	2.695712	0.483357	0.033849	0.043927

Πίνακας 6.28: Πίνακας χρονικών τιμών fit - score σε πολυταξική και δυαδική κατηγοριοποίηση

σαν μια αύξηση στη δυαδική κατηγοριοποίηση, Πίνακας 6.28. Ο GNB χρειάζεται να υπολογίζει τις υπό συνθήκη πιθανότητες για κάθε πρόβλημα δυαδικής κατηγοριοποίησης ξεχωριστά. Από την άλλη στον MLP αυξάνεται ο αριθμός των προβλημάτων δυαδικής κατηγοριοποίησης και συνεπώς αυξάνεται και η πολυπλοκότητα της εκπαίδευσης του μοντέλου.

Όσον αφορά το score\_time παρατηρείται αύξηση στους D-Tree, GNB, LogReg, MLP, RF και XGBoost αλγόριθμους. Αυτοί οι αλγόριθμοι απαιτούν συνήθως περισσότερους υπολογισμούς κατά τη διάρκεια της πρόβλεψης σε σύγκριση με το χρόνο προσαρμογής, καθώς μπορεί να χρειαστεί να διασχίσουν το δέντρο απόφασης, να υπολογίσουν πιθανότητες ή να εκτελέσουν πολύπλοκους υπολογισμούς για κάθε περίπτωση εισόδου. Ο αυξημένος χρόνος αποτελεσμάτων για αυτούς τους αλγόριθμους στη δυαδική κατηγοριοποίηση θα μπορούσε να οφείλεται στους λόγους που αναφέρθηκαν.

Το συμπέρασμα που εξάγετε από τα παραπάνω πειράματα είναι ότι η αντιμετώπιση ενός προβλήματος κατηγοριοποίησης πολλαπλών κατηγοριών ως δυαδικό πρόβλημα μπορεί να αποδώσει καλύτερα αποτελέσματα σε ορισμένες περιπτώσεις. Συγκεκριμένα, οι αποδώσεις του D-Tree υποδηλώνουν ότι είναι ιδιαίτερα κατάλληλος για εργασίες κατηγοριοποίησης όπου τα όρια απόφασης είναι σχετικά απλά.

Ένας πιθανός λόγος είναι η μείωση της πολυπλοκότητας του προβλήματος κατηγοριοποίησης, με την αντιμετώπιση του προβλήματος ως δυαδικό τα μοντέλα θα πρέπει να διακρίνουν μόνο μεταξύ δύο κλάσεων κάθε φορά και όχι πολλών ταυτόχρονα. Η απλούστευση αυτή μπορεί να διευκολύνει τα μοντέλα να μάθουν τα υποκείμενα μοτίβα στα δεδομένα και να παράγουν ακριβέστερες προβλέψεις. Επι-

---

πλέον, οι καλύτερες μετρικές του D-Tree οφείλονται στο γεγονός ότι το μοντέλο μπορεί να παράξει απλούστερα όρια αποφάσεων όταν αντιμετωπίζει δυαδικά προβλήματα κατηγοριοποίησης. Τα απλούστερα αυτά όρια μπορεί να μειώσουν το ρίσκο της υπερπροσαρμογής στα δεδομένα εκπαίδευσης, με αποτέλεσμα να οδηγηθούν σε καλύτερες αποδόσεις γενίκευσης.

Επομένως, η προσέγγιση αυτή μπορεί να θεωρηθεί ως εναλλακτική λύση στη χρήση πολύ πιο σύνθετων μοντέλων απευθείας για την κατηγοριοποίηση πολλαπλών κατηγοριών, ειδικά όταν η πολυπλοκότητα του προβλήματος κατηγοριοποίησης είναι αρκετά υψηλή.



# Κεφάλαιο 7

## Βελτιστοποίηση Παραμέτρων

Τα μοντέλα μηχανικής μάθησης συχνά χρειάζονται μια διαδικασία λεπτομερούς ρύθμισης προκειμένου να επιτύχουν βέλτιστες αποδόσεις σε μια δεδομένη εργασία. Αυτή η διαδικασία γνωστή και ως συντονισμός μοντέλου ή ρύθμιση υπερπαραμέτρων, περιλαμβάνει τη ρύθμιση των παραμέτρων του μοντέλου για την εύρεση του καλύτερου συνδυασμού που μεγιστοποιεί τις αποδόσεις του.

Οι υπερπαραμέτροι αρχικοποιούνται πριν από την έναρξη της διαδικασίας εκπαίδευσης και καθορίζουν τη συμπεριφορά του μοντέλου. Μερικά παραδείγματα υπερπαραμέτρων είναι ο ρυθμός μάθησης, ο αριθμός των κρυφών στρωμάτων σε ένα νευρωνικό δίκτυο, η ισχύς κανονικοποίησης και ο αριθμός των δέντρων σε ένα τυχαίο δάσος. Δεδομένου ότι αυτές οι υπερπαραμέτροι δεν μαθαίνονται από τα δεδομένα εκπαίδευσης, θα πρέπει να οριστούν με βάση προηγούμενης γνώσης ή μέσω μιας σειράς πειραμάτων.

Η διαδικασία ρύθμισης των υπερπαραμέτρων μπορεί να είναι χρονοβόρα και να απαιτεί πολλούς πόρους, καθώς συχνά περιλαμβάνει την εκπαίδευση και την αξιολόγηση πολλαπλών παραλλαγών του ίδιου μοντέλου. Ωστόσο, τα οφέλη της ρύθμισης μπορεί να είναι σημαντικά.

Επιλέγοντας τις βέλτιστες υπερπαραμέτρους, μπορεί να βελτιωθεί η ακρίβεια του μοντέλου, να μειωθεί η υπερπροσαρμογή και να βελτιωθεί η απόδοση γενίκευσης. Υπάρχουν διάφορες τεχνικές για τον συντονισμό των υπερπαραμέτρων, που κυμαίνονται από το Grid Search έως πιο προηγμένες μεθόδους όπως η Bayesian βελτιστοποίηση και οι εξελικτικοί αλγόριθμοι. Το Grid Search ή αλλιώς αναζήτηση πλέγματος περιλαμβάνει τη συστηματική δοκιμή όλων των πιθανών συνδυασμών υπερπαραμέτρων εντός ενός προκαθορισμένου εύρους. Ενώ η Bayesian βελτιστο-

ποίηση και οι εξελικτικοί αλγόριθμοι χρησιμοποιούν στατιστικές τεχνικές για την επιλογή των θεωρητικά κατάλληλων παραμέτρων για να αξιολογηθούν έπειτα.

Συνοψίζοντας, ο συντονισμός του μοντέλου ή ο συντονισμός των υπερπαραμέτρων είναι ένα κρίσιμο βήμα στη ροή εργασίας της μηχανικής μάθησης που βοηθά στη βελτιστοποίηση της απόδοσης των μοντέλων. Υπάρχουν διάφορες τεχνικές διαθέσιμες για τον συντονισμό των υπερπαραμέτρων και η επιλογή της μεθόδου εξαρτάται από το συγκεκριμένο πρόβλημα και τους διαθέσιμους πόρους.

## 7.1 Παραμετροποίηση αλγορίθμων δυαδικής κατηγοριοποίησης

### 7.1.1 Grid search CV

Σύμφωνα με την υποενότητα 6.1 όπου αναλύθηκε το πρόβλημα της δυαδικής κατηγοριοποίησης, οι αλγόριθμοι με τις καλύτερες επιδόσεις όπως απεικονίζονται στον Πίνακα 7.1 ήταν οι:

Μοντέλα	Ακρίβεια	Ακρίβεια Θετικών Προβλέψεων	Ανάκληση	ROC AUC	F1
D-Tree	0.995042	0.995046	0.995042	0.994973	0.995042
RF	0.996519	0.996525	0.996519	0.999850	0.996519
XGBoost	0.996730	0.996737	0.996730	0.999970	0.996730

Πίνακας 7.1: Αποτελέσματα μετρικών στη δυαδική κατηγοριοποίηση

Χρησιμοποιώντας, την τεχνική του Grid Search CV έγινε η ανάλυση για τον κάθε αλγόριθμο και βρέθηκαν οι κατάλληλοι παράμετροι για τα παρακάτω σενάρια:

- D-Tree

Παράμετροι	Ακρίβεια	Ακρίβεια Θετ. Προβλ.	Ανάκληση	F1	ROC AUC	Συνδυασμός
class_weight	None	None	balanced	None	None	None
criterion	entropy	gini	entropy	entropy	entropy	entropy
max_depth	16	3	None	16	16	16
min_impurity_decrease	0.0	0.0	0.0	0.0	0.0	0.0
min_samples_leaf	1	1	1	1	1	1
min_samples_split	6	6	6	6	6	6
presort	False	True	True	False	False	False
splitter	best	random	random	best	best	best

Πίνακας 7.2: Τιμές παραμέτρων D-Tree στη δυαδική κατηγοριοποίηση

Μετρικές	Ακρίβεια	Ακρίβεια Θετ. Προβλ.	Ανάκληση	F1	ROC AUC	Συνδυασμός
Ακρίβεια	0.996941	0.880380	0.996624	0.996941	0.996941	0.996941
Ακρίβεια Θετ.Προβλ.	0.996079	1.0	0.994026	0.996079	0.996079	0.996079
Ανάκληση	0.997229	0.738206	0.998615	0.997229	0.997229	0.997229
F1	0.996654	0.847013	0.996315	0.996654	0.996654	0.996654
Roc_Auc	0.996964	0.869103	0.996782	0.996964	0.996964	0.996964
Συνδυασμός	0.996773	0.866940	0.996473	0.996773	0.996773	0.996773

Πίνακας 7.3: Αποτελέσματα μετρικών D-Tree στη δυαδική κατηγοριοποίηση

- Random Forest

Παράμετροι	Ακρίβεια	Ακρίβεια Θετ. Προβλ.	Ανάκληση	F1	ROC AUC	Συνδυασμός
bootstrap	False	False	True	False	False	False
max_depth	None	5	None	None	None	None
max_features	sqrt	sqrt	sqrt	log2	sqrt	log2
min_samples_leaf	1	1	2	1	1	1
min_samples_split	5	10	5	2	5	2
n_estimators	50	100	100	100	50	100

Πίνακας 7.4: Τιμές παραμέτρων Random Forest στη δυαδική κατηγοριοποίηση

Μετρικές	Ακρίβεια	Ακρίβεια Θετ. Προβλ.	Ανάκληση	F1	ROC AUC	Συνδυασμός
Ακρίβεια	0.997468	0.990506	0.996835	0.997468	0.997468	0.997468
Ακρίβεια Θετ.Προβλ.	0.997001	0.998591	0.995168	0.997001	0.997001	0.997001
Ανάκληση	0.997461	0.980607	0.997922	0.997461	0.997461	0.997461
F1	0.997229	0.989499	0.996540	0.997229	0.997229	0.997229
Roc_Auc	0.997468	0.989721	0.996921	0.997468	0.997468	0.997468
Συνδυασμός	0.997325	0.989785	0.996677	0.997325	0.997325	0.997325

Πίνακας 7.5: Αποτελέσματα μετρικών Random Forest στη δυαδική κατηγοριοποίηση

- XGBoost

Παράμετροι	Ακρίβεια	Ακρίβεια Θετ. Προβλ.	Ανάκληση	F1	ROC AUC	Συνδυασμός
colsample_bytree	0.8	0.8	0.8	0.8	0.8	0.8
max_depth	5	None	3	5	5	5
n_estimators	100	200	100	100	100	100
scale_pos_weight	5	1	5	5	5	5
subsample	1.0	0.9	0.9	1.0	1.0	1.0

Πίνακας 7.6: Τιμές παραμέτρων XGBoost στη δυαδική κατηγοριοποίηση

Μετρικές	Ακρίβεια	Ακρίβεια Θετ. Προβλ.	Ανάκληση	F1	ROC AUC	Συνδυασμός
Ακρίβεια	0.997785	0.997574	0.997257	0.997785	0.997785	0.997785
Ακρίβεια Θετ. Προβλ.	0.996773	0.997230	0.995627	0.996773	0.996773	0.996773
Ανάκληση	0.998384	0.997461	0.998384	0.998384	0.998384	0.998384
F1	0.997577	0.997344	0.997002	0.997577	0.997577	0.997577
Roc_Auc	0.997832	0.997565	0.997346	0.997832	0.997832	0.997832
Συνδυασμός	0.997670	0.997435	0.997123	0.997832	0.997670	0.997670

Πίνακας 7.7: Αποτελέσματα μετρικών XGBoost στη δυαδική κατηγοριοποίηση

### 7.1.2 Randomized search CV

- D-Tree

Παράμετροι	Ακρίβεια	Ακρίβεια Θετ. Προβλ.	Ανάκληση	F1	ROC AUC	Συνδυασμός
splitter	best	random	best	best	best	best
presort	True	False	True	True	True	True
min_samples_split	8	10	8	8	8	8
min_samples_leaf	4	4	4	4	4	4
min_impurity_decrease	0.0	0.0	0.0	0.0	0.0	0.0
max_depth	20	5	20	20	20	20
criterion	entropy	gini	entropy	entropy	entropy	entropy
class_weight	balanced	None	balanced	balanced	balanced	balanced

Πίνακας 7.8: Τιμές παραμέτρων D-Tree στη δυαδική κατηγοριοποίηση

Μετρικές	Ακρίβεια	Ακρίβεια Θετ. Προβλ.	Ανάκληση	F1	ROC AUC	Συνδυασμός
Ακρίβεια	0.996097	0.944515	0.996097	0.996097	0.996097	0.996097
Ακρίβεια Θετ.Προβλ.	0.995158	0.999498	0.995158	0.995158	0.995158	0.995158
Ανάκληση	0.996306	0.879023	0.996306	0.996306	0.996306	0.996306
F1	0.995730	0.934746	0.995730	0.995730	0.995730	0.995730
Roc_Auc	0.996114	0.939317	0.996114	0.996114	0.996114	0.996114
Συνδυασμός	0.995881	0.939420	0.995881	0.995881	0.995881	0.995881

Πίνακας 7.9: Αποτελέσματα μετρικών D-Tree στη δυαδική κατηγοριοποίηση

- Random Forest

Παράμετροι	Ακρίβεια	Ακρίβεια Θετ. Προβλ.	Ανάκληση	F1	ROC AUC	Συνδυασμός
n_estimators	300	100	200	300	300	300
min_samples_split	10	10	10	10	10	10
min_samples_leaf	1	1	2	1	1	1
max_depth	None	3	None	None	None	None
criterion	gini	entropy	gini	gini	gini	gini
class_weight	None	None	None	None	None	None

Πίνακας 7.10: Τιμές παραμέτρων Random Forest στη δυαδική κατηγοριοποίηση

Μετρικές	Ακρίβεια	Ακρίβεια Θετ. Προβλ.	Ανάκληση	F1	ROC AUC	Συνδυασμός
Ακρίβεια	0.996941	0.985338	0.996730	0.996941	0.996941	0.996941
Ακρίβεια Θετ.Προβλ.	0.995854	0.998339	0.995395	0.995854	0.995854	0.995854
Ανάκληση	0.997461	0.969528	0.997461	0.997461	0.997461	0.997461
F1	0.996653	0.983704	0.996424	0.996653	0.996653	0.996653
Roc_Auc	0.996982	0.984084	0.996788	0.996982	0.996982	0.996982
Συνδυασμός	0.996778	0.984199	0.996559	0.996778	0.996778	0.996778

Πίνακας 7.11: Αποτελέσματα μετρικών Random Forest στη δυαδική κατηγοριοποίηση

- XGBoost

Παράμετροι	Ακρίβεια	Ακρίβεια Θετ. Προβλ.	Ανάκληση	F1	ROC AUC	Συνδυασμός
subsample	0.8	1.0	0.8	0.8	0.8	0.8
reg_lambda	0.01	0.1	0.01	0.01	0.01	0.01
reg_alpha	0.01	0.1	0.01	0.01	0.01	0.01
n_estimators	100	300	100	100	100	100
min_child_weight	5	1	5	5	5	5
max_depth	None	3	None	None	None	None
learning_rate	0.1	0.1	0.1	0.1	0.1	0.01
gamma	0	5	0	0	0	0
colsample_bytree	1.0	0.9	1.0	1.0	1.0	1.0

Πίνακας 7.12: Τιμές παραμέτρων XGBoost στη δυαδική κατηγοριοποίηση

Μετρικές	Ακρίβεια	Ακρίβεια Θετ. Προβλ.	Ανάκληση	F1	ROC AUC	Συνδυασμός
Ακρίβεια	0.996817	0.994409	0.996817	0.996817	0.996817	0.996817
Ακρίβεια Θετ. Προβλ.	0.995856	0.997903	0.995856	0.995856	0.995856	0.995856
Ανάκληση	0.997111	0.989842	0.997111	0.997111	0.997111	0.997111
F1	0.996923	0.993850	0.996923	0.996923	0.996923	0.996923
Roc_Auc	0.996405	0.994047	0.996405	0.996405	0.996405	0.996405
Συνδυασμός	0.996211	0.994010	0.996211	0.996211	0.996211	0.996211

Πίνακας 7.13: Αποτελέσματα μετρικών XGBoost στη δυαδική κατηγοριοποίηση

### 7.1.3 Σύνοψη

Στην υποενότητα αυτή χρησιμοποιήθηκαν οι τεχνικές του Grid Search CV και του Randomized Search CV για την παραμετροποίηση των αλγορίθμων D-Tree, Random Forest και XGBoost στο κομμάτι της δυαδικής κατηγοριοποίησης.

Από την ανάλυση παρατηρήθηκε ότι ο D – Tree κατάφερε με τους επιλεγμένους συνδυασμούς να βελτιώσει τα ποσοστά όλων των μετρικών σε σχέση με τον RF και τον XGBoost που δεν μπόρεσαν να βελτιώσουν τη μετρική ROC AUC. Όπως προαναφέρθηκε σε προηγούμενη ενότητα, η μετρική ROC AUC μετρά την ικανότητα του μοντέλου να διακρίνει μεταξύ των θετικών και αρνητικών δειγμάτων σε διαφορετικά όρια πιθανότητας. Όταν μετρικές όπως η ακρίβεια, η ακρίβεια των θετικών προβλέψεων, η ανάκληση και το f1 σκόρ βελτιώνονται αυτό σημαίνει ότι η συνολική απόδοση του μοντέλου όσον αφορά τη σωστή κατηγοριοποίηση των θετικών και

αρνητικών δειγμάτων έχει βελτιωθεί. Ωστόσο, η ROC AUC εστιάζει στην ικανότητα του μοντέλου να κατατάσσει τα δείγματα και όχι να αποδίδει τις συγκεκριμένες ετικέτες κλάσης.

Επομένως είναι πιθανό βελτιστοποιήσεις σε μετρικές κατηγοριοποίησης να μην συνδέονται απαραίτητα σε βελτιώσεις στην ROC AUC, ειδικά σε περιπτώσεις που υπάρχει ανισότητα δεδομένων, μη βέλτιστης επιλογής του threshold, τυχαιότητα στις μεθόδους ενσωμάτωσης ή επίτευξη του ανώτατου ορίου απόδοσης.

## 7.2 Παραμετροποίηση αλγορίθμων πολυταξικής κατηγοριοποίησης

Μοντέλο	Ανάκληση	Ακρίβεια Θετ. Προβλ.	F1	Ακρίβεια
D-Tree	0.867171	0.868509	0.867083	0.867171
RF	0.860804	0.860891	0.860454	0.860804
XGB	0.853906	0.853983	0.853643	0.853906

Πίνακας 7.14: Πίνακας μετρικών κατάστασης Cross Validation στην πολυταξική κατηγοριοποίηση

### 7.2.1 Grid search SV

- D-Tree

Παράμετροι	Ακρίβεια
class_weight	balanced
criterion	gini
max_depth	None
min_impurity_decrease	0.0
min_samples_leaf	1
min_samples_split	2
presort	True
splitter	random

Πίνακας 7.15: Τιμές παραμέτρων D-Tree στην πολυταξική κατηγοριοποίηση

Μετρικές	Ακρίβεια
Ακρίβεια	0.867206
Ακρίβεια Θετ. Προβλ.	0.842781
Ανάκληση	0.843024
F1	0.842802
Συνδυασμός	0.842869

Πίνακας 7.16: Αποτελέσματα μετρικών D-Tree στην πολυταξική κατηγοριοποίηση

- Random Forest

Παράμετροι	Ακρίβεια
class_weight	None
criterion	entropy
max_depth	None
min_samples_leaf	1
min_samples_split	2
n_estimators	200

Πίνακας 7.17: Τιμές παραμέτρων Random Forest στην πολυταξική κατηγοριοποίηση

Μετρικές	Ακρίβεια
Ακρίβεια	0.862436
Ακρίβεια Θετ. Προβλ.	0.837531
Ανάκληση	0.837389
F1	0.837270
Συνδυασμός	0.837397

Πίνακας 7.18: Αποτελέσματα μετρικών Random Forest στην πολυταξική κατηγοριοποίηση

- XGBoost

Παράμετροι	Ακρίβεια
colsample_bytree	0.5
gamma	0
learning_rate	0.01
max_depth	10
min_child_weight	1
n_estimators	300
subsample	1.0

Πίνακας 7.19: Τιμές παραμέτρων XGBoost στην πολυταξική κατηγοριοποίηση



Μετρικές	Ακρίβεια
Ακρίβεια	0.934011
Ακρίβεια Θετ.Προβλ.	0.937356
Ανάκληση	0.932256
F1	0.933577
Συνδυασμός	0.933757

Πίνακας 7.20: Αποτελέσματα μετρικών XGBoost στην πολυταξική κατηγοριοποίηση

### 7.2.2 Randomized search CV

- D-Tree

Παράμετροι	Ακρίβεια
splitter	best
presort	False
min_samples_split	10
min_samples_leaf	2
min_impurity_decrease	0.0
max_depth	20
criterion	entropy
class_weight	balanced

Πίνακας 7.21: Τιμές παραμέτρων D-Tree στην πολυταξική κατηγοριοποίηση

Μετρικές	Ακρίβεια
Ακρίβεια	0.857029
Ακρίβεια Θετ.Προβλ.	0.832563
Ανάκληση	0.831989
F1	0.826920
Συνδυασμός	0.830491

Πίνακας 7.22: Αποτελέσματα μετρικών D-Tree στην πολυταξική κατηγοριοποίηση

- Random Forest

Παράμετροι	Ακρίβεια
n_estimators	300
min_samples_split	10
min_samples_leaf	1
max_features	log2
max_depth	None
criterion	gini

Πίνακας 7.23: Τιμές παραμέτρων Random Forest στην πολυταξική κατηγοριοποίηση

Μετρικές	Ακρίβεια
Ακρίβεια	0.838263
Ακρίβεια Θετ.Προβλ.	0.808400
Ανάκληση	0.808912
F1	0.808300
Συνδυασμός	0.808537

Πίνακας 7.24: Αποτελέσματα μετρικών Random Forest στην πολυταξική κατηγοριοποίηση

- XGBoost

Παράμετροι	Ακρίβεια
subsample	1.0
n_estimators	300
min_child_weight	1
max_depth	10
learning_rate	0.01
gamma	0
colsample_bytree	0.5

Πίνακας 7.25: Τιμές παραμέτρων XGBoost στην πολυταξική κατηγοριοποίηση

Μετρικές	Ακρίβεια
Ακρίβεια	0.934345
Ακρίβεια Θετ.Προβλ.	0.939764
Ανάκληση	0.934429
F1	0.934267
Συνδυασμός	0.936153

Πίνακας 7.26: Αποτελέσματα μετρικών XGoost στην πολυταξική κατηγοριοποίηση

### 7.2.3 Σύνοψη

Στην υποενότητα αυτή παραμετροποιήθηκαν οι αλγόριθμοι D-Tree, RF και XGBoost οι οποίοι είχαν τις καλύτερες επιδόσεις στη πολυταξική κατηγοριοποίηση. Σε σύγκριση με την παραμετροποίηση που αναλύθηκε στην προηγούμενη ενότητα, εδώ παρατηρείται ότι υπάρχει μόνο ένας συνδυασμός μετρικών.

Επίσης, παρατηρείται ότι δεν αυξάνονται όλες οι μετρικές. Υπάρχει ένας συμβιβασμός μεταξύ των μετρικών κατά το συντονισμό των παραμέτρων. Πιο συγκεκριμένα η τιμή της ακρίβειας αυξήθηκε ενώ οι τιμές των άλλων μετρικών μειώθηκαν. Η βελτίωση και η μείωση των τιμών μπορεί να αποδοθεί στις αντισταθμίσεις μεταξύ των διάφορων μετρικών αξιολόγησης. Κάθε μετρική δίνει έμφαση σε διαφορετικές

---

πτυχές της απόδοσης του μοντέλου και η βελτιστοποίηση μιας μετρικής μπορεί να αποβεί εις βάρος μιας άλλης.

Όσον αφορά τον XGBoost, εμφανίζει σημαντική βελτιστοποίηση σε όλες τις μετρικές κατά τη διάρκεια του Randomized Search CV. Αυτό υποδηλώνει ότι η τυχαία δειγματοληψία συνδυασμών υπερπαραμέτρων στο Randomized Search CV εξερεύνησε ένα ευρύτερο εύρος τιμών υπερπαραμέτρων, οδηγώντας στην εύρεση ενός βέλτιστου συνδυασμού για το XGBoost .

# Κεφάλαιο 8

## Συμπεράσματα

Στην παρούσα διπλωματική εργασία, διερευνήθηκε η εφαρμογή αλγορίθμων μηχανικής μάθησης με στόχο την πρόβλεψη ανωμαλιών σε ένα ηλεκτρικό σύστημα. Τα σύνολα δεδομένων που αναλύθηκαν παρείχαν πολύτιμες πληροφορίες σχετικά με την εμφάνιση και τους τύπους σφαλμάτων εντός του συστήματος. Πραγματοποιήθηκαν δυαδικές και πολυταξικές κατηγοριοποιήσεις με στόχο την ανάπτυξη ακριβών μοντέλων ανίχνευσης ανωμαλιών που θα μπορούσαν να βοηθήσουν στην προληπτική συντήρηση και τον εντοπισμό σφαλμάτων.

Τα αποτελέσματα έδειξαν ότι οι αλγόριθμοι D-Tree, RF και XGBoost επέδειξαν σταθερά καλύτερες επιδόσεις τόσο στις δυαδικές όσο και στις πολυταξικές κατηγοριοποιήσεις. Οι αλγόριθμοι αυτοί πέτυχαν σταθερά υψηλά επίπεδα ακρίβειας, με τα καλύτερα μοντέλα να επιτυγχάνουν ακρίβεια περίπου 99% στη δυαδική και 85-86% στην πολυταξική κατηγοριοποίηση. Αυτό αποδεικνύει την αποτελεσματικότητα αυτών των αλγορίθμων στον εντοπισμό ανωμαλιών και στην κατηγοριοποίηση τους σε συγκεκριμένους τύπους σφαλμάτων.

Επιπλέον, όταν αντιμετωπίστηκε η κατηγοριοποίηση πολλαπλών κατηγοριών ως δυαδικό πρόβλημα, παρατηρήθηκε αξιοσημείωτη βελτίωση της συνολικής απόδοσης με ακρίβεια που έφτασε το 94-95%. Αυτό υποδηλώνει ότι η απλούστευση του έργου κατηγοριοποίησης με το συνδυασμό των τύπων σφάλματος σε μία ενιαία κλάση “ανωμαλίας” μπορεί να ενισχύσει την προγνωστική ισχύ των μοντέλων.

Διερευνήθηκαν επίσης, οι δυνατότητες συντονισμού των υπερπαραμέτρων χρησιμοποιώντας τεχνικές όπως ο Grid Search CV και ο Randomized Search CV. Μέσω του συντονισμού ήταν δυνατή η βελτιστοποίηση της απόδοσης των αλγορίθμων. Στην περίπτωση της δυαδικής κατηγοριοποίησης παρατηρήθηκε μικρή βελτίωση της

---

ακρίβειας μετά τον συντονισμό των υπερπαραμέτρων. Ωστόσο, στην κατηγοριοποίηση των πολλαπλών κατηγοριών, ο XGBoost παρουσίασε σημαντική βελτίωση των επιδόσεων μετά τη ρύθμιση, επισημαίνοντας τη σημασία της σωστής ρύθμισης των υπερπαραμέτρων για την επίτευξη καλύτερων αποτελεσμάτων.

Ενώ η παρούσα μελέτη παρέχει πολύτιμες πληροφορίες για την ανίχνευση ανωμαλιών στο ηλεκτρικό σύστημα με τη χρήση αλγορίθμων μηχανικής μάθησης, υπάρχουν αρκετοί δρόμοι για μελλοντική ανάλυση που μπορούν να διερευνηθούν.

1. Μηχανική χαρακτηριστικών: Διερεύνηση πρόσθετων χαρακτηριστικών ή τεχνικών μηχανικής που μπορούν να καταγράψουν διαφοροποιημένα μοτίβα στα δεδομένα του ηλεκτρικού συστήματος. Αυτό μπορεί να προϋποθέτει την ενσωμάτωση γνώσεων του τομέα ή την εξαγωγή χρονικών ή χωρικών χαρακτηριστικών.
2. Μέθοδοι ενσωμάτωσης: Διερεύνηση των δυνατοτήτων μεθόδων ενσωμάτωσης όπως το bagging, το boosting ή το stacking για περαιτέρω βελτίωση της προβλεπτικής απόδοσης των μοντέλων. Οι μέθοδοι ενσωμάτωσης συνδυάζουν πολλαπλά μοντέλα για να αξιοποιήσουν τα επιμέρους πλεονεκτήματα και να βελτιώσουν τη συνολική ακρίβεια και ανθεκτικότητα.
3. Τοπικοποίηση ανωμαλιών: Η ανάλυση επεκτείνεται έτσι ώστε να συμπεριλαμβάνει τεχνικές για την τοπικοποίηση των ανωμαλιών, οι οποίες περιλαμβάνουν τον εντοπισμό των συγκεκριμένων εξαρτημάτων ή υποσυστημάτων εντός του ηλεκτρικού συστήματος που συμβάλλουν στις ανωμαλίες που ανιχνεύονται. Οι πληροφορίες αυτές μπορεί να παρέχουν πολύτιμες πληροφορίες με στόχο τη στοχευμένη συντήρηση και αντιμετώπιση προβλημάτων.
4. Online ανίχνευση ανωμαλιών: Διερεύνηση της διαδικασίας ανάπτυξης online μοντέλων ανίχνευσης ανωμαλιών που μπορούν να παρακολουθούν συνεχώς το ηλεκτρικό σύστημα σε πραγματικό χρόνο και να παρέχουν άμεσες ειδοποιήσεις όταν εντοπίζονται ανωμαλίες. Αυτό θα επιτρέψει την προληπτική συντήρηση και την ταχεία αντίδραση σε πιθανές βλάβες ή σφάλματα.

# Βιβλιογραφία

- [1] Better factory project.
- [2] Bosch iot insights: your end-to-end service solution for iot data management.
- [3] Digiman4.0.
- [4] Electricity system.
- [5] The future is our point.
- [6] Horse project.
- [7] K-nearest neighbors algorithm.
- [8] Maintenance strategy - the choices available.
- [9] Organisational behaviour with new technologies: a human resources management model for industry 4.0.
- [10] Platform-enabled kits of artificial intelligence for an easy uptake by smes.
- [11] Robotics automation for warehousing, 3pls, distribution, manufacturing.
- [12] Sme 4.0.
- [13] Training and evaluating process.
- [14] Versatile plug-and-play platform enabling remote predictive maintenance.
- [15] Mlp from scratch, 2020.
- [16] Ethem Alpaydin. *Introduction to machine learning*. MIT press, 2020.
- [17] Tom Alumni. Healthcare digitization: saving lives one data point at a time at bjc healthcare. 2017.
- [18] Angelos Angelopoulos, Emmanouel T Michailidis, Nikolaos Nomikos, Panagiotis Trakadas, Antonis Hatziefremidis, Stamatis Voliotis, and Theodore Zahariadis. Tackling faults in the industry 4.0 era—a survey of machine-learning solutions and key aspects. *Sensors*, 20(1):109, 2019.
- [19] Ebru Turanoglu Bekar, Per Nyqvist, and Anders Skoogh. An intelligent approach for data pre-processing and analysis in predictive maintenance with an industrial case study. *Advances in Mechanical Engineering*, 12(5):1687814020919207, 2020.

- 
- [20] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, volume 4. Springer, 2006.
- [21] Brad Boehmke and Brandon M Greenwell. *Hands-on machine learning with R*. CRC press, 2019.
- [22] Matteo Calabrese, Martin Cimmino, Francesca Fiume, Martina Manfrin, Luca Romeo, Silvia Ceccacci, Marina Paolanti, Giuseppe Toscano, Giovanni Ciandrini, Alberto Carrotta, et al. Sophia: An event-based iot and machine learning architecture for predictive maintenance in industry 4.0. *Information*, 11(4):202, 2020.
- [23] Thyago P Carvalho, Fabrízio AAMN Soares, Roberto Vita, Roberto da P Francisco, João P Basto, and Symone GS Alcalá. A systematic literature review of machine learning methods applied to predictive maintenance, 2019.
- [24] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. pages 785–794, 2016.
- [25] D Richard Cutler, Thomas C Edwards Jr, Karen H Beard, Adele Cutler, Kyle T Hess, Jacob Gibson, and Joshua J Lawler. Random forests for classification in ecology. *Ecology*, 88(11):2783–2792, 2007.
- [26] Katie Dudek. Team penske and the digital twin make a speedy pair. 2019.
- [27] James Gareth, Witten Daniela, Hastie Trevor, and Tibshirani Robert. *An introduction to statistical learning: with applications in R*. Spinger, 2013.
- [28] Hardik A Gohel, Himanshu Upadhyay, Leonel Lagos, Kevin Cooper, and Andrew Sanzetenea. Predictive maintenance architecture development for nuclear infrastructure using machine learning. *Nuclear Engineering and Technology*, 52(7):1436–1442, 2020.
- [29] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- [30] Trevor Hastie, Robert Tibshirani, Jerome H Friedman, and Jerome H Friedman. *The elements of statistical learning: data mining, inference, and prediction*, volume 2. Springer, 2009.
- [31] Mario Hermann, Tobias Pentek, Boris Otto, et al. Design principles for industrie 4.0 scenarios: a literature review. *Technische Universität Dortmund, Dortmund*, 45, 2015.
- [32] Paul Iusztin. Decision trees: from 0 to xgboost & lightgbm, 2022.
- [33] Ameeth Kanawaday and Aditya Sane. Machine learning for predictive maintenance of industrial machines using iot sensor data. pages 87–90, 2017.
- [34] DonHee Lee and Seong No Yoon. Application of artificial intelligence-based technologies in the healthcare industry: Opportunities and challenges, 2021.
- [35] Christopher Ludwig. Building volkswagen’s industrial cloud. 2019.
- [36] Vimala Mathew, Tom Toby, Vikram Singh, B Maheswara Rao, and M Goutham Kumar. Prediction of remaining useful lifetime (rul) of turbofan engine using machine learning. pages 306–311, 2017.

- 
- [37] Marc Morisson, Arnault and Pattinson. Industry 4.0 - interreg europe policy brief, 2019.
- [38] Marina Paolanti, Luca Romeo, Andrea Felicetti, Adriano Mancini, Emanuele Frontoni, and Jelena Loncarski. Machine learning approach for predictive maintenance in industry 4.0. pages 1–6, 2018.
- [39] J Ross Quinlan. *C4. 5: programs for machine learning*. Elsevier, 2014.
- [40] Khushwant Rai. The math behind logistic regression, 2020.
- [41] Rajeev DS Raizada and Yune-Sang Lee. Smoothness without smoothing: why gaussian naive bayes is not naive for multi-subject searchlight studies, 2013.
- [42] Rahul Rastogi. Support vector regression and it’s mathematical implementation, 2020.
- [43] Irina Rish et al. An empirical study of the naive bayes classifier. 3(22):41–46, 2001.
- [44] Smriti Saini. A gentle introduction to xgboost for applied machine learning, 2021.
- [45] Maicon Saturno, V Moura Pertel, Fernando Deschamps, and EDFR Loures. Proposal of an automation solutions architecture for industry 4.0, 2017.
- [46] Gustavo Scalabrini Sampaio, Arnaldo Rabello de Aguiar Vallim Filho, Leilton Santos da Silva, and Leandro Augusto da Silva. Prediction of motor failure time using an artificial neural network. *Sensors*, 19(19):4342, 2019.
- [47] Prashanth Subramaniam and Maninder Jeet Kaur. Review of security in mobile edge computing with deep learning, 2019.
- [48] Muhammad Syafrudin, Ganjar Alfian, Norma Latif Fitriyani, and Jongtae Rhee. Performance analysis of iot-based sensor, big data processing, and machine learning model for real-time monitoring system in automotive manufacturing. *Sensors*, 18(9):2946, 2018.