



ΠΑΝΕΠΙΣΤΗΜΙΟ ΠΕΙΡΑΙΩΣ
ΣΧΟΛΗ ΤΕΧΝΟΛΟΓΙΩΝ ΠΛΗΡΟΦΟΡΙΚΗΣ ΚΑΙ ΕΠΙΚΟΙΝΩΝΙΩΝ
ΤΜΗΜΑ ΨΗΦΙΑΚΩΝ ΣΥΣΤΗΜΑΤΩΝ

ΠΡΟΓΡΑΜΜΑ ΜΕΤΑΠΤΥΧΙΑΚΩΝ ΣΠΟΥΔΩΝ
“ΠΛΗΡΟΦΟΡΙΑΚΑ ΣΥΣΤΗΜΑΤΑ & ΥΠΗΡΕΣΙΕΣ”

«Εφαρμογή Τεχνικών Εξόρυξης Γνώσης σε Οικονομικά Δεδομένα –
Πλεονεκτήματα και Μειονεκτήματα σε μια Τράπεζα και στις Πιστωτικές με Χρήση
Python»

Ιλιάννα Ελένη Σδράλια

Υποβάλλεται για την εκπλήρωση των προϋποθέσεων λήψης Μεταπτυχιακού
Διπλώματος στην ειδίκευση «ΜΔΑ/ΠΠΣ/ΠΔ» του ΠΜΣ “Πληροφοριακά
Συστήματα & Υπηρεσίες” στο ΠΑΝΕΠΙΣΤΗΜΙΟ ΠΕΙΡΑΙΩΣ

2025

Επιβλέπων: Φιλιππάκης Μιχαήλ

Ακαδημαϊκή Θέση: Καθηγητής

ΣΕΛΙΔΑ ΕΓΚΥΡΟΤΗΤΑΣ

Όνοματεπώνυμο Φοιτήτριας: Σδράλια Ιλιάνα Ελένη

Τίτλος Μεταπτυχιακής Διπλωματικής Εργασίας: Εφαρμογή τεχνικών εξόρυξης γνώσης σε οικονομικά δεδομένα – πλεονεκτήματα και μειονεκτήματα σε μια τράπεζα και στις πιστωτικές με χρήση Python. Η παρούσα Μεταπτυχιακή Διπλωματική Εργασία υποβάλλεται ως μερική εκπλήρωση των απαιτήσεων του Προγράμματος Μεταπτυχιακών Σπουδών “Πληροφοριακά Συστήματα & Υπηρεσίες” του Τμήματος Ψηφιακών Συστημάτων του Πανεπιστημίου Πειραιώς και εγκρίθηκε στις 02/09/2025 από τα μέλη της Εξεταστικής Επιτροπής.

Εξεταστική Επιτροπή

Επιβλέπων (Τμήμα Ψηφιακών Συστημάτων, Πανεπιστήμιο Πειραιώς): Μιχαήλ Φιλιππάκης, Καθηγητής
Μέλος Εξεταστικής Επιτροπής: Μαρία Χαλκίδη, Καθηγήτρια

Μέλος Εξεταστικής Επιτροπής: Δημοσθένης Κυριαζής, Καθηγητής


ΥΠΕΥΘΥΝΗ ΔΗΛΩΣΗ ΑΥΘΕΝΤΙΚΟΤΗΤΑΣ

Η Ιλιάνα Ελένη Σδράλια, γνωρίζοντας τις συνέπειες της λογοκλοπής, δηλώνω υπεύθυνα ότι η παρούσα εργασία με τίτλο «Εφαρμογή τεχνικών εξόρυξης γνώσης σε οικονομικά δεδομένα – πλεονεκτήματα και μειονεκτήματα σε μια τράπεζα και στις πιστωτικές με χρήση Python», αποτελεί προϊόν αυστηρά προσωπικής εργασίας και όλες οι πηγές που έχω χρησιμοποιήσει, έχουν δηλωθεί κατάλληλα στις βιβλιογραφικές παραπομπές και αναφορές. Τα σημεία όπου έχω χρησιμοποιήσει ιδέες, κείμενο 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. Επιπλέον δηλώνω υπεύθυνα ότι η συγκεκριμένη Μεταπτυχιακή Διπλωματική Εργασία έχει συγγραφεί από εμένα προσωπικά και δεν έχει υποβληθεί ούτε έχει αξιολογηθεί στο πλαίσιο κάποιου άλλου μεταπτυχιακού ή προπτυχιακού τίτλου σπουδών, στην Ελλάδα ή στο εξωτερικό. Παράβαση της ανωτέρω ακαδημαϊκής μου ευθύνης αποτελεί ουσιώδη λόγο για την ανάκληση του πτυχίου μου. Σε κάθε περίπτωση, αναληθούς ή ανακριβούς δηλώσεως, υπόκειμαι στις συνέπειες που προβλέπονται τις διατάξεις που προβλέπει η Ελληνική και Κοινοτική Νομοθεσία περί πνευματικής ιδιοκτησίας.

Η ΔΗΛΟΥΣΑ

Όνοματεπώνυμο: Ιλιάνα Ελένη Σδράλια

Αριθμός Μητρώου: ME2352

Υπογραφή: 

ΕΥΧΑΡΙΣΤΙΕΣ

Θα ήθελα να εκφράσω τις θερμές μου ευχαριστίες στον επιβλέποντα καθηγητή μου κ. Φιλιππάκη Μιχαήλ για την πολύτιμη καθοδήγηση, στήριξη και εποικοδομητική συμβολή του καθ' όλη τη διάρκεια της εκπόνησης της παρούσας διπλωματικής εργασίας. Επίσης, ευχαριστώ την οικογένειά μου για την αμέριστη στήριξη και υπομονή που επέδειξε σε όλη τη διάρκεια των σπουδών μου.

Περιεχόμενα

Περίληψη	9
1. Μεγάλα Δεδομένα (Big Data)	10
1.1. Ορισμός και Ιδιότητες των Big Data.....	10
1.2. Τεχνολογίες και Υποδομές Big Data.....	12
1.3. Τομείς Εφαρμογής Μεγάλων Δεδομένων.....	13
2. Τα Μεγάλα Δεδομένα στον Τραπεζικό Τομέα	15
2.1. Ανίχνευση Απάτης (Fraud Detection)	17
2.2. Προσωποποιημένη Προώθηση Προϊόντων (Personalized Marketing).....	19
2.3. Διαχείριση Ρίσκου (Risk Management)	20
2.4. Αξία του Χρόνου Ζωής του Πελάτη (Customer Lifetime Value – CLV)	22
2.5. Τμηματοποίηση Πελατών (Customer Segmentation).....	23
2.6. Συστήματα Συστάσεων (Recommendation Systems)	24
3. Εξόρυξη Δεδομένων (Data Mining).....	26
3.1. Στάδια της Διαδικασίας KDD.....	28
3.2. Τεχνικές Μάθησης στην Εξόρυξη Δεδομένων	30
3.3. Μοντέλα Μηχανικής Μάθησης	32
3.3.1. Random Forest Classifier.....	33
3.3.2. Logistic Regression	34
3.3.3. Gradient Boost Classifier	36
3.4. Μετρικές Αξιολόγησης.....	37
4. Ηθικά Ζητήματα και Κανονιστικές Προκλήσεις στη Χρήση Τεχνητής Νοημοσύνης στον Τραπεζικό Τομέα	39
4.1. Ζητήματα Μεροληψίας και Διακρίσεων (Algorithmic Bias)	40
4.2. Διαφάνεια και Ερμηνευσιμότητα Μοντέλων (Explainability).....	43
4.3. Προστασία Προσωπικών Δεδομένων και Κανονιστική Συμμόρφωση	45

4.4. Δίκαιη Πρόσβαση σε Χρηματοοικονομικές Υπηρεσίες.....	47
4.5. Προτάσεις Ορθής Διαχείρισης Ηθικών Κινδύνων.....	49
5. Βαθιά Μηχανική Μάθηση (Deer Learning) σε Οικονομικά και Τραπεζικά Δεδομένα.....	51
5.1. Νευρωνικά Δίκτυα και Εφαρμογές σε Τραπεζικά και Οικονομικά Δεδομένα	52
5.2. Δομικά Στοιχεία Νευρωνικών Δικτύων σε Τραπεζικά Συστήματα.....	53
5.3. Βασικοί Αλγόριθμοι Βαθιάς Μάθησης σε Οικονομικά και Τραπεζικά Δεδομένα	55
6. Πειραματικό Μέρος / Υλοποίηση	58
6.1. Περιγραφή Συνόλων Δεδομένων.....	58
6.2. Προ-επεξεργασία Δεδομένων.....	60
6.3. Αποτελέσματα – Σύγκριση Αποτελεσμάτων.....	61
6.3.1. Bank Marketing Dataset.....	62
6.3.2. Credit Card Clients Dataset	77
7. Συμπεράσματα.....	97
Βιβλιογραφία	100

Πίνακας Εικόνων

Εικόνα 1: Οι 5 διαστάσεις (5 Vs) των Big Data, Πηγή: BotPenguin. (χ.χ.). Five V of Big Data	11
Εικόνα 2: Διαγραμματική απεικόνιση της διαδικασίας εξόρυξης γνώσης (KDD), από το αρχικό σύνολο δεδομένων έως την παραγωγή πληροφορίας. Πηγή: Morgan Kaufmann (Han et al., Data Mining: Concepts and Techniques).	27
Εικόνα 3: Κατηγοριοποίηση της Μηχανικής Μάθησης με βάση το είδος της μάθησης και τους τύπους εφαρμογών, Πηγή: ResearchGate.....	32
Εικόνα 4: Κατανομή της μεταβλητής στόχου Subscription (Y) του συνόλου δεδομένων Bank Marketing.63	
Εικόνα 5: Πίνακας συσχέτισης μεταξύ των αριθμητικών μεταβλητών του συνόλου δεδομένων Bank Marketing.....	65
Εικόνα 6: Σύγκριση ROC καμπυλών για τα μοντέλα ταξινόμησης υπό διαφορετικές στρατηγικές διαχείρισης ανισορροπίας δεδομένων: Baseline, SMOTE και Undersampling – Bank Marketing.....	74
Εικόνα 7: Κατανομή συχνοτήτων της μεταβλητής στόχου Default (Y) του συνόλου δεδομένων Credit Card Clients.....	78
Εικόνα 8: Boxplot της μεταβλητής X1 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y. .80	
Εικόνα 9: Boxplot της μεταβλητής X2 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y. .80	
Εικόνα 10: Boxplot της μεταβλητής X3 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.81	
Εικόνα 11: Boxplot της μεταβλητής X4 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.	81
Εικόνα 12: Boxplot της μεταβλητής X5 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.	82
Εικόνα 13: Boxplot της μεταβλητής X6 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.83	
Εικόνα 14: Boxplot της μεταβλητής X7 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.83	
Εικόνα 15: Boxplot της μεταβλητής X8 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.83	
Εικόνα 16: Boxplot της μεταβλητής X9 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.83	
Εικόνα 17: Boxplot της μεταβλητής X10 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.	84
Εικόνα 18: Boxplot της μεταβλητής X11 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.	84
Εικόνα 19: Boxplot της μεταβλητής X12 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.	84

Εικόνα 20: Βοχplot της μεταβλητής X13 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.	84
Εικόνα 21: Βοχplot της μεταβλητής X14 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.	85
Εικόνα 22: Βοχplot της μεταβλητής X15 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.	85
Εικόνα 23: Βοχplot της μεταβλητής X16 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.	85
Εικόνα 24: Βοχplot της μεταβλητής X17 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.	85
Εικόνα 25: Βοχplot της μεταβλητής X18 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.	86
Εικόνα 26: Βοχplot της μεταβλητής X19 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.	86
Εικόνα 27: Βοχplot της μεταβλητής X20 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.	86
Εικόνα 28: Βοχplot της μεταβλητής X21 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.	86
Εικόνα 29: Βοχplot της μεταβλητής X22 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.	87
Εικόνα 30: Βοχplot της μεταβλητής X23 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.	87
Εικόνα 31: Πίνακας συσχέτισης των αριθμητικών μεταβλητών του συνόλου δεδομένων Credit Card Clients.....	88
Εικόνα 32: Σύγκριση ROC καμπυλών για διαφορετικούς ταξινομητές και στρατηγικές διαχείρισης ανισορροπίας: Random Forest (Baseline και SMOTE), Logistic Regression, Gradient Boosting και Undersampling RF.	95

Λίστα Πινάκων

Table 1: Αποτελέσματα ταξινόμησης: Precision, Recall και F1-Score ανά κατηγορία για το baseline μοντέλο της Random Forest, καθώς και συνολικές επιδόσεις (Accuracy, Macro και Weighted Averages) στο σύνολο δεδομένων Bank Marketing.....	67
Table 2: Αποτελέσματα ταξινόμησης: Precision, Recall και F1-Score ανά κατηγορία για το μοντέλο της Random Forest, μετά την εφαρμογή της τεχνικής SMOTE, καθώς και συνολικές επιδόσεις (Accuracy, Macro και Weighted Averages) στο σύνολο δεδομένων Bank Marketing.....	69
Table 3: Αποτελέσματα ταξινόμησης: Precision, Recall και F1-Score ανά κατηγορία για το μοντέλο της Random Forest, μετά την εφαρμογή της τεχνικής Undersampling, καθώς και συνολικές επιδόσεις (Accuracy, Macro και Weighted Averages) στο σύνολο δεδομένων Bank Marketing.	71
Table 4: Αποτελέσματα ταξινόμησης (Precision, Recall, F1-Score) για τον βελτιστοποιημένο Random Forest με Grid Search και εφαρμογή SMOTE στο σύνολο δεδομένων Bank Marketing.	76
Table 5: Αποτελέσματα ταξινόμησης του baseline μοντέλου Random Forest χωρίς εξισορρόπηση δεδομένων για το σύνολο δεδομένων Credit Card Clients.	90
Table 6: Αποτελέσματα ταξινόμησης του μοντέλου Random Forest, μετά την εφαρμογή της τεχνικής SMOTE, για το σύνολο δεδομένων Credit Card Clients.	92
Table 7: Αποτελέσματα ταξινόμησης του μοντέλου Random Forest, μετά την εφαρμογή της μεθόδου Undersampling, για το σύνολο δεδομένων Credit Card Clients.	93
Table 8: Τιμές ROC-AUC για τα εξεταζόμενα μοντέλα (Random Forest, Logistic Regression, Gradient Boosting) με και χωρίς εφαρμογή τεχνικών εξισορρόπησης δεδομένων για το σύνολο δεδομένων Credit Card Clients.	94

Περίληψη

Η παρούσα διπλωματική εργασία εξετάζει την εφαρμογή τεχνικών εξόρυξης γνώσης και μηχανικής μάθησης σε οικονομικά δεδομένα, επικεντρωμένη ειδικά στον τραπεζικό τομέα και στη διαχείριση πιστωτικού κινδύνου. Βασικός στόχος είναι η αξιοποίηση μεγάλων συνόλων δεδομένων για την εξαγωγή χρήσιμης γνώσης που μπορεί να υποστηρίξει καλύτερες στρατηγικές λήψης αποφάσεων, βελτιώνοντας παράλληλα τη διαχείριση ρίσκου και την απόδοση τραπεζικών προϊόντων.

Για την υλοποίηση, χρησιμοποιήθηκαν δύο δημόσια διαθέσιμα σύνολα δεδομένων: το Bank Marketing Dataset, που αφορά τη στόχευση πελατών σε τραπεζικά προϊόντα μέσω καμπανιών μάρκετινγκ, και το Default of Credit Card Clients Dataset, το οποίο σχετίζεται με την πρόβλεψη αθέτησης πληρωμών από πελάτες πιστωτικών καρτών. Τα δεδομένα προεπεξεργάστηκαν με τεχνικές μετασχηματισμού, όπως το One-Hot Encoding για τις κατηγορικές μεταβλητές και κανονικοποίηση για τις αριθμητικές.

Η ανισορροπία δεδομένων αντιμετωπίστηκε με τις μεθόδους SMOTE (Synthetic Minority Oversampling Technique) και Random Undersampling, προκειμένου να ενισχυθεί η ακρίβεια των ταξινομητών στις υποεκπροσωπούμενες κατηγορίες. Για τη δημιουργία μοντέλων χρησιμοποιήθηκαν αλγόριθμοι επιβλεπόμενης μάθησης, όπως Random Forest, Logistic Regression και Gradient Boosting. Τα αποτελέσματα αξιολογήθηκαν με βάση μετρικές όπως Accuracy, Precision, Recall, F1-Score και ROC-AUC.

Στο σύνολο δεδομένων Bank Marketing, ο Random Forest επέδειξε ROC-AUC 0.92 με υψηλή σταθερότητα, ενώ ο Gradient Boosting προσέγγισε το 0.91. Στο σύνολο δεδομένων πιστωτικών καρτών, οι βέλτιστες αποδόσεις επιτεύχθηκαν και πάλι με Random Forest και Gradient Boosting, αναδεικνύοντας την ευελιξία και αποτελεσματικότητα αυτών των αλγορίθμων στη διαχείριση πολύπλοκων οικονομικών δεδομένων.

Η εργασία καταλήγει στο συμπέρασμα ότι η συνδυαστική χρήση τεχνικών προεπεξεργασίας και προηγμένων μοντέλων μηχανικής μάθησης προσφέρει σημαντικές δυνατότητες βελτίωσης της ακρίβειας πρόβλεψης σε πραγματικές τραπεζικές εφαρμογές. Επιπλέον, παρέχονται προτάσεις για περαιτέρω έρευνα με χρήση πραγματικών τραπεζικών δεδομένων και μελέτη πιο εξελιγμένων μεθόδων βαθιάς μάθησης.

1. Μεγάλα Δεδομένα (Big Data)

Η ραγδαία τεχνολογική εξέλιξη των τελευταίων δεκαετιών έχει επιφέρει βαθιές αλλαγές στον τρόπο με τον οποίο οι άνθρωποι επικοινωνούν, εργάζονται, καταναλώνουν και ζουν. Η ενσωμάτωση των τεχνολογιών πληροφορικής και επικοινωνιών στην καθημερινότητα δημιούργησε έναν τεράστιο όγκο δεδομένων, ο οποίος αυξάνεται με εκθετικούς ρυθμούς. Τα δεδομένα αυτά δημιουργούνται από ποικίλες πηγές: έξυπνα κινητά, αισθητήρες, συστήματα συναλλαγών, μέσα κοινωνικής δικτύωσης, μηχανές αναζήτησης, συσκευές IoT, εφαρμογές υγείας, και άλλες.

Σε αυτό το περιβάλλον, αναδύθηκε ο όρος «Μεγάλα Δεδομένα» (Big Data), ο οποίος περιγράφει το σύνολο αυτών των τεράστιων, ποικίλων και ταχύρρυθμα παραγόμενων πληροφοριών, που καθίστανται δύσκολα διαχειρίσιμες με τα συμβατικά εργαλεία βάσεων δεδομένων και ανάλυσης [1]. Η ανάγκη για την ορθή αξιοποίηση αυτών των δεδομένων είναι ζωτικής σημασίας για την ανταγωνιστικότητα και την επιβίωση των σύγχρονων επιχειρήσεων.

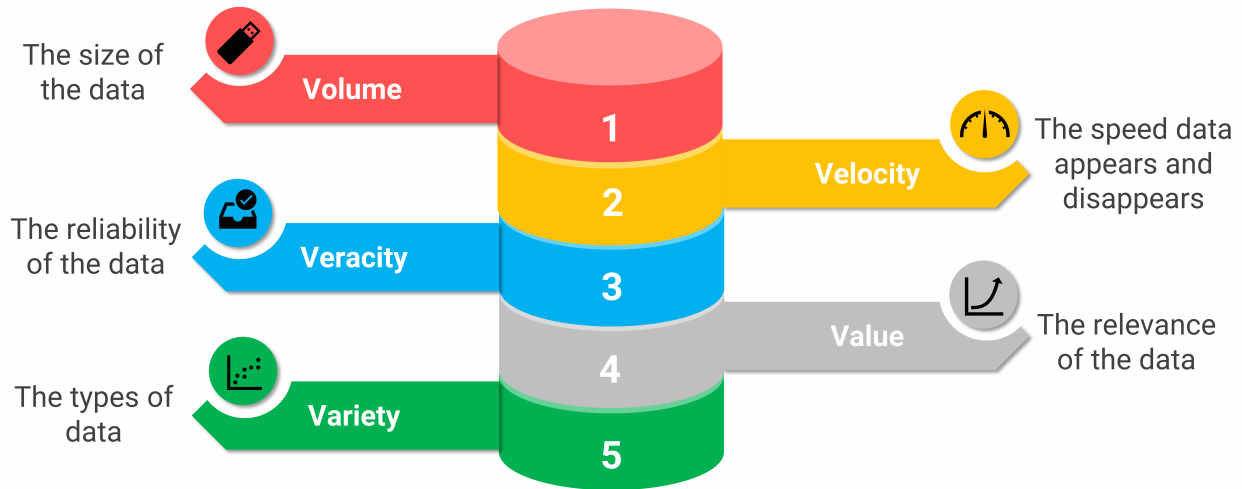
Η σημασία των Big Data αντικατοπτρίζεται ξεκάθαρα και σε οικονομικό επίπεδο. Σύμφωνα με έκθεση του Forbes, η παγκόσμια αγορά για Big Data software και υπηρεσίες αυξάνεται με ρυθμό πάνω από 10% ετησίως, από 42 δισ. δολάρια το 2018 σε πάνω από 100 δισ. το 2026. Αυτή η ανάπτυξη δεν περιορίζεται σε έναν μόνο τομέα – σχεδόν κάθε κλάδος της οικονομίας, από την υγεία και τις τηλεπικοινωνίες έως τις μεταφορές και τη ναυτιλία, βασίζεται πλέον στην ανάλυση δεδομένων για την επίτευξη στρατηγικών στόχων [1], [2].

Η μετάβαση αυτή απαιτεί όχι μόνο τεχνολογικές υποδομές και συστήματα για την αποθήκευση και επεξεργασία τεράστιων όγκων πληροφορίας, αλλά και εξειδικευμένο ανθρώπινο δυναμικό: data scientists, analysts, engineers, αρχιτέκτονες δεδομένων. Η διαχείριση, ανάλυση και προστασία των Big Data έχει αναδειχθεί σε ανεξάρτητο επιστημονικό και επαγγελματικό πεδίο. Στη συνέχεια, εξετάζονται οι βασικές τεχνολογίες, τα εργαλεία και οι υπολογιστικές αρχιτεκτονικές που καθιστούν δυνατή τη συστηματική αξιοποίηση αυτών των δεδομένων.

1.1. Ορισμός και Ιδιότητες των Big Data

Τα Μεγάλα Δεδομένα δεν ορίζονται μόνο από το μέγεθος τους. Η επιστημονική κοινότητα, για να αποδώσει πλήρως τα χαρακτηριστικά τους, έχει θεσπίσει ένα πλαίσιο που συνοψίζεται στα **5V's** [2] (Volume, Velocity, Variety, Veracity, Value), το οποίο αναλύεται παρακάτω:

The 5 Vs of Big Data



Εικόνα 1: Οι 5 διαστάσεις (5 Vs) των Big Data, Πηγή: BotPenguin. (χ.χ.). Five V of Big Data

- Όγκος (Volume): Ο καθαρός όγκος των δεδομένων είναι τεράστιος. Εταιρείες και οργανισμοί αποθηκεύουν καθημερινά gigabytes, terabytes, ακόμη και petabytes πληροφοριών. Μόνο το 2018 δημιουργήθηκαν 33 zettabytes δεδομένων παγκοσμίως, και η πρόβλεψη είναι ότι έως το 2035 θα ξεπεράσουν τα 2.000 zettabytes. Η ανάγκη για αποδοτικές λύσεις αποθήκευσης (π.χ. cloud storage, κατανεμημένες βάσεις δεδομένων) είναι επιτακτική.
- Ταχύτητα (Velocity): Η ταχύτητα με την οποία δημιουργούνται, επεξεργάζονται και μεταδίδονται τα δεδομένα είναι χωρίς προηγούμενο. Η ανάγκη για real-time ανάλυση καθιστά τις συμβατικές μεθόδους επεξεργασίας παρωχημένες. Η επεξεργασία σε πραγματικό χρόνο επιτρέπει, για παράδειγμα, την άμεση ανίχνευση απάτης σε τραπεζικές συναλλαγές ή την παροχή προσωποποιημένων προσφορών στους καταναλωτές.
- Ποικιλία (Variety): Τα δεδομένα σήμερα δεν είναι μόνο δομημένα (structured), όπως π.χ. αριθμητικές βάσεις δεδομένων. Το μεγαλύτερο μέρος τους είναι μη δομημένα (unstructured) – εικόνες, ήχοι, βίντεο, κείμενα, συναισθηματική ανάλυση, log files – ή ημιδομημένα (semi-structured), όπως αρχεία XML ή JSON. Η επεξεργασία τους απαιτεί ευέλικτα και εξελιγμένα συστήματα, όπως NoSQL βάσεις δεδομένων (MongoDB, Cassandra), Hadoop, Spark, και εργαλεία data lakes.

- Εγκυρότητα (Veracity): Η αξιοπιστία των δεδομένων είναι κρίσιμη. Πολύ συχνά τα δεδομένα περιέχουν ασυνέπειες, ελλείψεις ή ακόμα και σφάλματα. Η διαδικασία καθαρισμού (data cleaning) και επαλήθευσης είναι απαραίτητη πριν τη χρήση τους, ειδικά σε ευαίσθητες εφαρμογές όπως η ιατρική ή ο τραπεζικός τομέας.
- Αξία (Value): Το σημαντικότερο από όλα είναι η αξία των δεδομένων – η δυνατότητα εξαγωγής ουσιαστικής και χρήσιμης πληροφορίας. Τα Big Data δεν είναι χρήσιμα απλώς και μόνο λόγω του μεγέθους τους, αλλά λόγω του τρόπου που μπορούν να μετατραπούν σε **γνώση**. Η επιχειρηματική ευφυΐα (Business Intelligence), η πρόβλεψη τάσεων, η βελτιστοποίηση παραγωγής ή μάρκετινγκ είναι μερικές μόνο από τις εφαρμογές της αξιοποίησης αυτών των δεδομένων.

1.2. Τεχνολογίες και Υποδομές Big Data

Η διαχείριση και ανάλυση των Μεγάλων Δεδομένων προϋποθέτει προηγμένες τεχνολογικές υποδομές, ικανές να ανταπεξέλθουν στον τεράστιο όγκο, την πολυπλοκότητα και τη συνεχή ροή της πληροφορίας. Η ανάγκη για αποδοτική αποθήκευση, γρήγορη ανάκτηση και ταχεία επεξεργασία των δεδομένων έχει οδηγήσει στην ευρεία ανάπτυξη και υιοθέτηση κατανεμημένων συστημάτων, τα οποία αποτελούν το θεμέλιο της αρχιτεκτονικής Big Data. Στο επίκεντρο αυτών των υποδομών βρίσκονται πλατφόρμες όπως το Hadoop και το Apache Spark, οι οποίες προσφέρουν δυνατότητες οριζόντιας κλιμάκωσης, κατανεμημένης επεξεργασίας και ανθεκτικότητας σε σφάλματα. Το Hadoop, ως πλατφόρμα ανοικτού κώδικα, επιτρέπει την αποθήκευση τεράστιων ποσοτήτων δεδομένων σε κατανεμημένους κόμβους, ενώ μέσω του συστήματος MapReduce προσφέρει έναν μηχανισμό επεξεργασίας που διαμοιράζει το υπολογιστικό φορτίο σε πολλαπλές μονάδες.

Σε πιο σύγχρονο πλαίσιο, το Apache Spark λειτουργεί ως εξέλιξη του Hadoop, παρέχοντας υποστήριξη για επεξεργασία εντός της κύριας μνήμης (in-memory processing), γεγονός που το καθιστά ταχύτερο και πιο αποτελεσματικό σε εφαρμογές που απαιτούν επαναλαμβανόμενη ανάλυση δεδομένων ή επεξεργασία σε πραγματικό χρόνο. Το Spark υποστηρίζει επίσης λειτουργίες μηχανικής μάθησης, ροής δεδομένων (streaming), γραφημάτων και άλλων εξελιγμένων μοντέλων επεξεργασίας, καθιστώντας το ιδιαίτερα δημοφιλές στην επιστημονική και επιχειρησιακή κοινότητα.

Παράλληλα με τις πλατφόρμες κατανεμημένης επεξεργασίας, σημαντικό ρόλο διαδραματίζουν οι βάσεις δεδομένων τύπου NoSQL. Σε αντίθεση με τις παραδοσιακές σχεσιακές βάσεις, οι NoSQL βάσεις, όπως η

MongoDB και η Cassandra, σχεδιάστηκαν για να διαχειρίζονται μεγάλα και ποικίλα σύνολα δεδομένων, στα οποία η δομή δεν είναι απαραίτητα σταθερή ή προκαθορισμένη. Οι συγκεκριμένες τεχνολογίες προσφέρουν μεγαλύτερη ευελιξία στην αποθήκευση, ειδικά όταν πρόκειται για δεδομένα με υψηλή ποικιλία μορφών (κείμενα, έγγραφα, JSON, logs), όπως συναντώνται σε εφαρμογές κοινωνικής δικτύωσης, συστήματα αισθητήρων και ψηφιακές πλατφόρμες περιεχομένου.

Τέλος, ένα σημαντικό κομμάτι του τεχνολογικού παζλ των Big Data αποτελεί το υπολογιστικό νέφος (cloud computing). Πλατφόρμες όπως το Amazon Web Services (AWS), το Microsoft Azure και το Google Cloud παρέχουν τη δυνατότητα για δυναμική και οικονομική αποθήκευση, ανάλυση και διάθεση δεδομένων σε παγκόσμια κλίμακα. Το μοντέλο «pay-as-you-go» επιτρέπει στις επιχειρήσεις να χρησιμοποιούν μόνο τους πόρους που χρειάζονται, μειώνοντας το κόστος επένδυσης σε υποδομές και εξοπλισμό. Επιπλέον, οι υπηρεσίες cloud ενσωματώνουν μηχανισμούς ασφάλειας, διαθεσιμότητας και αυτοματοποίησης, καθιστώντας τις ιδανικές για την ανάπτυξη ευέλικτων και επεκτάσιμων λύσεων Big Data.

1.3. Τομείς Εφαρμογής Μεγάλων Δεδομένων

Η ολοένα και μεγαλύτερη εξάρτηση των επιχειρήσεων από την πληροφορία, έχει οδηγήσει στη συνειδητοποίηση της ανάγκης ανάλυσης των δεδομένων που παράγονται κατά τη διάρκεια της λειτουργίας τους. Μέσα από την αξιοποίηση των Big Data, ενισχύεται τόσο η στρατηγική λήψης αποφάσεων όσο και η δυνατότητα πρόβλεψης μελλοντικών εξελίξεων. Ένα χαρακτηριστικό παράδειγμα είναι η μηχανή αναζήτησης της Google, η οποία όχι μόνο βελτιώνει συνεχώς τις δυνατότητές της βάσει των αναζητήσεων των χρηστών σε παγκόσμιο επίπεδο, αλλά προσαρμόζει και τα αποτελέσματα στις προσωπικές προτιμήσεις κάθε ατόμου, βασισμένη στο ιστορικό του και στις αλληλεπιδράσεις του στο διαδίκτυο. Από αυτό το γεγονός καθίσταται προφανές πως η σύγχρονη επιχειρηματικότητα καλείται να ακολουθήσει ένα μοντέλο περισσότερο επικεντρωμένο στον πελάτη. Με τον τρόπο αυτό, οι εταιρείες είναι σε θέση όχι μόνο να ενισχύσουν τα κέρδη τους προωθώντας εξατομικευμένα προϊόντα, αλλά και να προβλέπουν πιο αποτελεσματικά τις τάσεις της αγοράς, την οικονομική βιωσιμότητά τους και τις ανάγκες του καταναλωτικού κοινού σε βάθος χρόνου.

Η πρακτική εφαρμογή των τεχνικών ανάλυσης δεδομένων έχει επιφέρει αξιοσημείωτες βελτιώσεις σε διάφορους τομείς. Ένας από τους πρώτους τομείς που επηρεάστηκαν είναι αυτός των τηλεπικοινωνιών. Οι εταιρείες παροχής τέτοιων υπηρεσιών διαχειρίζονται τεράστιους όγκους δεδομένων σε πραγματικό χρόνο, που αφορούν σε δημογραφικά στοιχεία χρηστών, καταναλωτικές συνήθειες, προτιμήσεις σε εφαρμογές και τύπους υπηρεσιών. Η αξιοποίηση αυτών των δεδομένων επιτρέπει την ανάπτυξη

προσαρμοσμένων τιμολογιακών πολιτικών ανά χρήστη, βελτιστοποιώντας έτσι τη σχέση κόστους-οφέλους. Ταυτόχρονα, μέσω της ανάλυσης της ζήτησης και των προτιμήσεων, είναι δυνατή η αποτελεσματικότερη διαχείριση των τηλεπικοινωνιακών υποδομών, μειώνοντας προβλήματα συμφόρησης και επιτρέποντας την καλύτερη κατανομή πόρων. Ως αποτέλεσμα, οι πάροχοι εξασφαλίζουν υψηλότερα ποσοστά ικανοποίησης πελατών, αλλά και αυξημένη απόδοση των υπηρεσιών τους.

Σημαντική είναι επίσης η συνεισφορά των Big Data στον τομέα της υγείας, όπου η ανάγκη για προσωποποιημένες θεραπευτικές προσεγγίσεις είναι αυξανόμενη. Η πληθώρα δεδομένων που προκύπτει από τα ιατρικά ιστορικά, τα αποτελέσματα εξετάσεων, τις ηλεκτρονικές συνταγογραφήσεις και τα wearable devices, επιτρέπει τη διαμόρφωση ακριβέστερων διαγνωστικών μοντέλων και προβλέψεων. Οι γιατροί πλέον μπορούν να λαμβάνουν αποφάσεις βάσει στατιστικά τεκμηριωμένων προτύπων, ελαχιστοποιώντας την πιθανότητα λάθους. Επιπλέον, η δυνατότητα γεωγραφικού εντοπισμού και ανάλυσης δημογραφικών στοιχείων αποδείχθηκε καίρια κατά την πανδημία COVID-19, συμβάλλοντας στη δημιουργία διαδραστικών χαρτών διάδοσης και στην εφαρμογή τοπικών μέτρων πρόληψης.

Ένας αναπτυσσόμενος τομέας όπου τα Big Data έχουν δώσει ισχυρή ώθηση είναι οι μεταφορές. Η σύγχρονη έρευνα και τεχνολογία επικεντρώνεται πλέον στη δημιουργία έξυπνων και ασφαλών συστημάτων μετακίνησης. Η ενσωμάτωση αισθητήρων σε οχήματα, σε συνδυασμό με την ανάλυση δεδομένων σε πραγματικό χρόνο, επιτρέπει την πρόβλεψη και αποφυγή επικίνδυνων καταστάσεων όπως η υπνηλία του οδηγού ή η επικείμενη σύγκρουση. Αντίστοιχα, λύσεις όπως η δυναμική προσαρμογή κυκλοφοριακών ρυθμίσεων και η βελτιστοποίηση διαδρομών βασίζονται σε υποδομές που επεξεργάζονται τεράστιες ποσότητες πληροφοριών. Η συμβολή της ανάλυσης δεδομένων στην ελαχιστοποίηση των καθυστερήσεων, στην εξοικονόμηση ενέργειας και στη μείωση των εκπομπών ρύπων είναι καθοριστική, ειδικά σε μια εποχή όπου οι περιβαλλοντικοί στόχοι αποτελούν προτεραιότητα.

Στον τομέα της ναυτιλίας, οι νέες τεχνολογίες διαμορφώνουν ένα πιο αποδοτικό, αυτοματοποιημένο και οικολογικά φιλικό πρότυπο λειτουργίας. Πλοία εξοπλισμένα με συστήματα τηλεμετρίας, αισθητήρες και εργαλεία απομακρυσμένης παρακολούθησης μπορούν να συλλέγουν συνεχώς πληροφορίες σχετικές με τη λειτουργική τους κατάσταση, τις καιρικές συνθήκες ή την πορεία πλεύσης. Η αξιοποίηση αυτών των δεδομένων ενισχύει την πρόληψη τεχνικών προβλημάτων, τη βελτιστοποίηση κατανάλωσης καυσίμου και τη λήψη κρίσιμων αποφάσεων σε πραγματικό χρόνο. Ο συνδυασμός αυτών των τεχνολογιών οδηγεί

σταδιακά στην εδραίωση της μη επανδρωμένης, ενεργειακά αποδοτικής και προγνωστικά ελεγχόμενης ναυτιλίας του μέλλοντος.

Αναπόσπαστο κομμάτι της σύγχρονης ψηφιακής στρατηγικής αποτελεί και ο αθλητισμός, όπου η τεχνολογία της ανάλυσης δεδομένων χρησιμοποιείται για να αναδείξει αθλητικά ταλέντα, να βελτιστοποιήσει τα προγράμματα εκγύμνασης και να ενισχύσει την προετοιμασία έναντι αντιπάλων. Τα στατιστικά δεδομένα που συλλέγονται για παίκτες, αγώνες και προπονήσεις μεταφράζονται σε πληροφορία στρατηγικής αξίας για τις ομάδες. Οι αθλητικοί οργανισμοί αποκτούν πλέον τη δυνατότητα να αξιολογούν σε πραγματικό χρόνο την απόδοση και τη φυσική κατάσταση των παικτών, να περιορίζουν τον κίνδυνο τραυματισμών και να διαμορφώνουν μακροπρόθεσμα πλάνα επιτυχίας. Η οικονομική σημασία αυτής της αγοράς είναι τεράστια, με τις επενδύσεις σε τεχνολογίες ανάλυσης να παρουσιάζουν σταθερή αύξηση διεθνώς.

Ο οικονομικός τομέας και ιδιαίτερα το χρηματοπιστωτικό σύστημα αποτελεί ίσως τον τομέα με τη μεγαλύτερη διείσδυση των Big Data. Οι τράπεζες και οι χρηματοοικονομικοί οργανισμοί αναλύουν καθημερινά gigabyte δεδομένων που προκύπτουν από συναλλαγές, κινήσεις λογαριασμών, χρήσεις πιστωτικών καρτών και πληρωμών, με στόχο την παροχή ασφαλών, ταχύτερων και πιο προσωποποιημένων υπηρεσιών. Οι πληροφορίες αυτές χρησιμοποιούνται για την αναγνώριση καταναλωτικών προτύπων, την ανίχνευση οικονομικής απάτης, την πιστοληπτική αξιολόγηση πελατών και την κατηγοριοποίηση του πελατολογίου σε ομάδες υψηλού ή χαμηλού ρίσκου. Τα δεδομένα επιτρέπουν την αυτόματη λήψη αποφάσεων και την παροχή χρηματοοικονομικών προϊόντων με τρόπο που ικανοποιεί τις απαιτήσεις του πελάτη, ενώ ταυτόχρονα προστατεύει τα συμφέροντα του οργανισμού.

2. Τα Μεγάλα Δεδομένα στον Τραπεζικό Τομέα

Η αλματώδης πρόοδος της τεχνολογίας της πληροφορικής και των επικοινωνιών έχει επιδράσει καταλυτικά στον τρόπο λειτουργίας πολλών κλάδων της οικονομίας, μεταξύ αυτών και του τραπεζικού τομέα. Η ενσωμάτωση ψηφιακών τεχνολογιών στις τραπεζικές υπηρεσίες και η μετάβαση στην πλήρη ψηφιοποίηση των συναλλαγών έχουν οδηγήσει στη δημιουργία και συγκέντρωση τεράστιων ποσοτήτων δεδομένων. Τα δεδομένα αυτά αφορούν όχι μόνο στην εσωτερική λειτουργία των τραπεζικών οργανισμών, αλλά και στη συμπεριφορά των πελατών τους, καλύπτοντας ένα ευρύ φάσμα όπως κινήσεις λογαριασμών, ιστορικό συναλλαγών, είδη προϊόντων που χρησιμοποιούνται, αλλά και γεωγραφικά και δημογραφικά χαρακτηριστικά.

Η αξιοποίηση τεχνολογιών όπως το cloud computing, το blockchain και η τεχνητή νοημοσύνη (AI), έχει ανοίξει νέους δρόμους για την καινοτομία και την ανάπτυξη στον τραπεζικό χώρο. Οι τεχνολογίες αυτές προσφέρουν τη δυνατότητα διαχείρισης και ανάλυσης μεγάλων, ετερογενών συνόλων δεδομένων με τρόπο αποδοτικό, ασφαλή και κλιμακούμενο. Παρόλα αυτά, ο μεγάλος όγκος, η ταχύτητα και η ποικιλομορφία των τραπεζικών δεδομένων καθιστούν την πλήρη αξιοποίησή τους ένα περίπλοκο και απαιτητικό εγχείρημα, το οποίο προϋποθέτει την ανάπτυξη προηγμένων τεχνικών εξόρυξης γνώσης και στρατηγικού σχεδιασμού.

Στην παρούσα φάση, η τραπεζική βιομηχανία καλείται να επενδύσει στην ευρύτερη και βαθύτερη συλλογή δεδομένων από όσο το δυνατόν περισσότερες πηγές, πέρα από τα ήδη υπάρχοντα ιστορικά αρχεία. Η έμφαση μετατοπίζεται σταδιακά από τα δομημένα δεδομένα των εσωτερικών συστημάτων προς πιο δυναμικές και ανεξερεύνητες πηγές, όπως οι εφαρμογές Internet of Things (IoT), τα κοινωνικά δίκτυα, τα κινητά τηλέφωνα και οι κυβερνητικές πλατφόρμες ανοικτών δεδομένων. Η ανάλυση των στοιχείων σε πραγματικό χρόνο (real-time analytics) γίνεται επιτακτική, καθώς οι τράπεζες επιδιώκουν να παρέχουν εξατομικευμένες υπηρεσίες, να εντοπίζουν έγκαιρα κινδύνους και να λαμβάνουν αποφάσεις με άμεση ισχύ και ακρίβεια.

Η επέκταση των καναλιών εισροής δεδομένων φέρνει στην επιφάνεια νέες προκλήσεις, μεταξύ των οποίων ξεχωρίζει η ανάγκη επεξεργασίας μη δομημένων και ημιδομημένων πληροφοριών. Παραδοσιακά, τα τραπεζικά πληροφοριακά συστήματα ήταν σχεδιασμένα να χειρίζονται κυρίως δομημένα δεδομένα, όπως αριθμητικά πεδία ή πεδία καθορισμένων τιμών. Όμως, τα νέα δεδομένα – κείμενα, εικόνες, αρχεία ήχου και βίντεο, δεδομένα από συσκευές IoT, ή αναρτήσεις στα κοινωνικά μέσα – δεν είναι εύκολα ταξινομήσιμα και απαιτούν διαφορετικά εργαλεία ανάλυσης και επεξεργασίας.

Η πελατοκεντρική προσέγγιση, που αποτελεί πλέον το επίκεντρο της στρατηγικής κάθε χρηματοπιστωτικού ιδρύματος, ενισχύεται σημαντικά μέσω της αξιοποίησης των Big Data. Οι τράπεζες δεν αρκούνται πλέον σε στατιστικά μοντέλα βασισμένα σε γενικευμένες παραδοχές, αλλά χρησιμοποιούν την πληθώρα των διαθέσιμων πληροφοριών για να κατανοήσουν εις βάθος τη συμπεριφορά και τις ανάγκες κάθε μεμονωμένου πελάτη. Μέσα από την ανάλυση των κοινωνικών δικτύων, για παράδειγμα, μπορούν να εντοπίζουν νέες καταναλωτικές τάσεις και να αντιδρούν άμεσα με στοχευμένες προσφορές ή πακέτα προϊόντων. Η χρήση τεχνικών φυσικής γλώσσας (Natural Language Processing – NLP) τους επιτρέπει να συλλέγουν και να επεξεργάζονται feedback πελατών από πολλαπλές πλατφόρμες, ενώ οι αλγόριθμοι machine learning ενισχύουν την ικανότητα πρόβλεψης συμπεριφορών ή μελλοντικών ενεργειών.

Ένα πολύ χαρακτηριστικό παράδειγμα αξιοποίησης δεδομένων σε πραγματικό χρόνο αφορά την ανίχνευση ύποπτων συναλλαγών. Ας υποθέσουμε ότι σε κάποιο κατάστημα γίνεται προσπάθεια πληρωμής μέσω κάρτας και το αντίστοιχο POS ζητά έγκριση. Αν την ίδια στιγμή, με βάση γεωγραφικά δεδομένα που συλλέγονται από τον πάροχο κινητής τηλεφωνίας, ο κάτοχος της κάρτας φαίνεται να βρίσκεται σε εντελώς διαφορετική τοποθεσία, τότε το τραπεζικό σύστημα μπορεί να μπλοκάρει αμέσως την πληρωμή. Ένα τέτοιο σύστημα προστασίας όχι μόνο μειώνει τα περιστατικά απάτης, αλλά ενισχύει και το αίσθημα εμπιστοσύνης των πελατών προς την τράπεζα. Παράλληλα, δίνει τη δυνατότητα στο χρηματοπιστωτικό ίδρυμα να τοποθετηθεί στην αγορά ως τεχνολογικά προηγμένο και ασφαλές.

Η συνεργασία μεταξύ τραπεζών και παρόχων άλλων υπηρεσιών – όπως τηλεπικοινωνιακές εταιρείες, e-commerce platforms, κυβερνητικοί οργανισμοί – μπορεί να δημιουργήσει οικοσυστήματα δεδομένων στα οποία διαμοιράζονται πληροφορίες προς όφελος όλων των εμπλεκόμενων. Μια τέτοια διαλειτουργικότητα απαιτεί αυστηρούς κανόνες για την προστασία της ιδιωτικότητας και της ασφάλειας των δεδομένων, αλλά δημιουργεί ταυτόχρονα ευκαιρίες για καινοτόμα προϊόντα και υπηρεσίες.

Αξίζει να σημειωθεί πως ο τραπεζικός τομέας διαθέτει, ίσως περισσότερο από κάθε άλλο κλάδο, μία από τις πιο πλούσιες βάσεις δεδομένων που σχετίζονται με καταναλωτική συμπεριφορά, οικονομική δυνατότητα και προφίλ κινδύνου. Οι πληροφορίες που απορρέουν από αυτές τις βάσεις (όπως δημογραφικά δεδομένα, συναλλακτικά πρότυπα, συχνότητα χρήσης καρτών, περιοχές δραστηριότητας) αποτελούν πολύτιμο κεφάλαιο για την τράπεζα, εφόσον αξιοποιηθούν σωστά και με στρατηγική προσέγγιση.

Στο σημείο αυτό κρίνεται σκόπιμο να παρουσιαστούν συγκεκριμένες περιοχές εφαρμογής των Big Data στον τραπεζικό τομέα, όπου η εξόρυξη γνώσης από δεδομένα λειτουργεί ως θεμέλιο για τη βελτιστοποίηση των διαδικασιών, την ενίσχυση της εμπειρίας του πελάτη, και την ελαχιστοποίηση κινδύνων.

2.1. Ανίχνευση Απάτης (Fraud Detection)

Η ανίχνευση απάτης αποτελεί έναν από τους σημαντικότερους τομείς εφαρμογής των Μεγάλων Δεδομένων στον τραπεζικό κλάδο, καθώς σχετίζεται άμεσα με την προστασία τόσο των οικονομικών πόρων των πελατών όσο και της φήμης και αξιοπιστίας των ίδιων των τραπεζικών ιδρυμάτων. Με τον όρο οικονομική απάτη αναφερόμαστε σε κάθε ενέργεια που περιλαμβάνει εσκεμμένη παραπλάνηση με σκοπό την παράνομη οικονομική ωφέλεια. Στο πλαίσιο των τραπεζικών συναλλαγών, η απάτη μπορεί να

πάρει πολλές μορφές και να εκδηλωθεί μέσα από διαφορετικά κανάλια, προκαλώντας σημαντικές οικονομικές απώλειες τόσο στους πελάτες όσο και στους ίδιους τους οργανισμούς.

Η έκταση του προβλήματος γίνεται εμφανής από στατιστικά στοιχεία διεθνών μελετών. Χαρακτηριστικά, σύμφωνα με την έκθεση της Javelin Strategy & Research για το έτος 2020, το κόστος της χρηματοοικονομικής απάτης παγκοσμίως άγγιξε τα 56 δισεκατομμύρια δολάρια, γεγονός που καταδεικνύει την επείγουσα ανάγκη για την ανάπτυξη και εφαρμογή προηγμένων συστημάτων πρόληψης και εντοπισμού τέτοιων ενεργειών [17].

Ανάμεσα στις συνηθέστερες μορφές τραπεζικής απάτης συγκαταλέγονται οι παραβιάσεις που σχετίζονται με τη χρήση χρεωστικών και πιστωτικών καρτών, όπως είναι η απάτη από αντιγραφή (cloning) ή η ανεξουσιοδοτητή χρήση σε περιβάλλον φυσικό ή διαδικτυακό. Επιπλέον, περιλαμβάνονται απάτες επιταγών, πρόσβαση σε τραπεζικούς λογαριασμούς χωρίς εξουσιοδότηση – είτε μέσω online banking, είτε μέσω τηλεφωνικών υπηρεσιών – καθώς και περιπτώσεις όπου το ίδιο το θύμα, εξαιτίας παραπλανητικών πρακτικών, εγκρίνει την μεταφορά χρηματικών ποσών σε λογαριασμούς των δραστών.

Η τεχνολογική πρόοδος και η εξάπλωση των ψηφιακών μέσων συναλλαγής, ενώ έχουν ενισχύσει την ευκολία και την ταχύτητα στη διεκπεραίωση των τραπεζικών δραστηριοτήτων, ταυτόχρονα έχουν διευρύνει τις δυνατότητες των επιτήδειων να εκμεταλλευτούν κενά ασφαλείας. Στο πλαίσιο αυτό, οι τράπεζες αξιοποιούν δεδομένα από πληθώρα πηγών και τα αναλύουν μέσω τεχνικών μηχανικής μάθησης, αναγνώρισης προτύπων (pattern recognition), ανάλυσης ακραίων τιμών και άλλων μεθόδων εξόρυξης γνώσης.

Τα δεδομένα που χρησιμοποιούνται μπορεί να περιλαμβάνουν το ποσό της συναλλαγής, την τοποθεσία, το είδος του εμπόρου, τη συχνότητα της χρήσης της κάρτας, καθώς και άλλες μεταβλητές που αποτυπώνουν τη "φυσιολογική" συναλλακτική συμπεριφορά ενός πελάτη. Όταν ένα νέο γεγονός αποκλίνει σημαντικά από το αναμενόμενο πρότυπο, το σύστημα εντοπίζει αυτή την ανωμαλία και αποδίδει στη συναλλαγή έναν δείκτη κινδύνου. Εφόσον ξεπεραστεί το καθορισμένο όριο, η συναλλαγή μπλοκάρεται ή τίθεται σε προσωρινή αναστολή για περαιτέρω διερεύνηση.

Αξίζει να σημειωθεί πως η απάτη εξελίσσεται με ταχύτατους ρυθμούς και δεν εμφανίζεται πάντα με τον ίδιο τρόπο. Για παράδειγμα, ενώ οι απάτες μέσω επιταγών έχουν μειωθεί σημαντικά την τελευταία δεκαετία, οι περιπτώσεις απάτης μέσω διαδικτύου – όπως το phishing και η κατάχρηση στοιχείων ταυτοποίησης – έχουν αυξηθεί θεαματικά. Έτσι, η προσαρμοστικότητα και η συνεχής ανανέωση των μεθόδων ανίχνευσης καθίστανται επιβεβλημένες.

2.2. Προσωποποιημένη Προώθηση Προϊόντων (Personalized Marketing)

Η εξατομικευμένη προώθηση προϊόντων συνιστά μία από τις πρώτες και πλέον διαδεδομένες εφαρμογές των τεχνικών εξόρυξης γνώσης στον χώρο των επιχειρήσεων, και ιδιαίτερα στον τραπεζικό τομέα. Σκοπός της πρακτικής αυτής είναι η αξιοποίηση των δεδομένων που διαθέτει ένας οργανισμός για να προσδιορίσει με μεγαλύτερη ακρίβεια ποιοι πελάτες είναι πιθανότερο να ανταποκριθούν θετικά σε μία προωθητική ενέργεια. Μέσω αυτής της προσέγγισης, επιτυγχάνεται σημαντική μείωση του κόστους, καθώς περιορίζεται η ανάγκη για μαζικές, γενικευμένες καμπάνιες και αντικαθίστανται από στοχευμένες ενέργειες υψηλής αποτελεσματικότητας.

Στην πράξη, τα διαθέσιμα δεδομένα προέρχονται είτε από τις εσωτερικές βάσεις του οργανισμού – όπως είναι τα ιστορικά αγορών, τα προφίλ των πελατών και τα δημογραφικά τους χαρακτηριστικά – είτε από εξωτερικές πηγές, όπως τα μέσα κοινωνικής δικτύωσης και άλλες ψηφιακές πλατφόρμες, όπου καταγράφεται η συμπεριφορά και το ενδιαφέρον των χρηστών για διάφορα προϊόντα. Μέσα από τη συνδυαστική ανάλυση αυτών των πηγών πληροφοριών, οι τράπεζες και οι επιχειρήσεις γενικότερα μπορούν να αναπτύξουν πολύπλοκα μοντέλα πρόβλεψης αγοραστικής συμπεριφοράς.

Τα μοντέλα αυτά, στηριζόμενα σε τεχνικές μηχανικής μάθησης, μπορούν να εντοπίζουν τάσεις, μοτίβα και κρυμμένες συσχετίσεις μεταξύ χαρακτηριστικών πελατών και προηγούμενων αποφάσεών τους, επιτρέποντας την ακριβή πρόβλεψη της πιθανότητας ανταπόκρισης ενός ατόμου σε μια συγκεκριμένη υπηρεσία ή προϊόν. Η δυνατότητα αυτή δίνει τη δυνατότητα στις επιχειρήσεις να προσαρμόσουν τα μηνύματά τους, την τιμολόγηση, ακόμη και τη δομή των προϊόντων τους ανάλογα με το προφίλ κάθε χρήστη, βελτιώνοντας σημαντικά τόσο την εμπειρία του πελάτη όσο και την αποδοτικότητα των εμπορικών στρατηγικών.

Η εφαρμογή αυτής της προσέγγισης στον τραπεζικό τομέα αποδίδει εξαιρετικά αποτελέσματα. Οι τράπεζες, μέσα από τη συνεχή καταγραφή των συναλλακτικών συνηθειών των πελατών τους, τις καταθέσεις, τα δάνεια, τις πληρωμές λογαριασμών και τις κινήσεις λογαριασμών, είναι σε θέση να δημιουργήσουν ολοκληρωμένα πελατειακά προφίλ. Αυτά τα προφίλ αξιοποιούνται για τη δημιουργία εξατομικευμένων προσφορών, όπως π.χ. πρόταση για νέο λογαριασμό αποταμίευσης, αύξηση ορίου κάρτας ή προνομιακό επιτόκιο δανεισμού. Επιπλέον, με την κατάλληλη αξιολόγηση της αγοραστικής πρόθεσης του κάθε πελάτη, αποφεύγεται η υπερπροβολή άσχετων προϊόντων που ενδέχεται να οδηγήσουν σε δυσαρέσκεια ή κόπωση από την πλευρά του χρήστη.

Η προώθηση των προϊόντων μέσω τεχνικών εξόρυξης γνώσης μπορεί να κατηγοριοποιηθεί σε επιμέρους στρατηγικές. Μία από αυτές είναι η απόκτηση νέων πελατών, όπου μέσω της σύγκρισης των χαρακτηριστικών των υπάρχοντων πελατών που ανταποκρίθηκαν θετικά σε προηγούμενες καμπάνιες, επιχειρείται η προσέγγιση νέων ατόμων με παρόμοιο προφίλ. Άλλη σημαντική στρατηγική είναι η διατήρηση πελατών, κατά την οποία η ανάλυση δεδομένων επιτρέπει την πρόβλεψη αποχώρησης (churn prediction) και την έγκαιρη παρέμβαση με εξατομικευμένες προσφορές, ώστε να μειωθεί η απώλεια πελατών. Υπάρχει επίσης η δυνατότητα αποκλεισμού πελατών με χαμηλή αξία για τον οργανισμό, εφόσον διαπιστωθεί ότι συνδέονται με αυξημένο ρίσκο ή χαμηλή συνεισφορά στα έσοδα της επιχείρησης, μια τακτική που επιτρέπει την πιο ορθολογική κατανομή των διαθέσιμων πόρων μάρκετινγκ.

Τέλος, μια ιδιαίτερα χρήσιμη τεχνική είναι η λεγόμενη ανάλυση καλαθιού αγορών (market basket analysis), η οποία επιδιώκει να αποκαλύψει ποια προϊόντα συνήθως αγοράζονται μαζί. Η γνώση αυτών των συσχετίσεων δίνει στις τράπεζες τη δυνατότητα να προτείνουν αυτόματα επιπλέον υπηρεσίες σε πελάτες με βάση τις επιλογές άλλων με παρόμοια συμπεριφορά – για παράδειγμα, αν ένας πελάτης ανοίξει λογαριασμό μισθοδοσίας, μπορεί να του προταθεί επιπλέον πακέτο ασφάλισης ή προπληρωμένη κάρτα.

2.3. Διαχείριση Ρίσκου (Risk Management)

Η διαχείριση ρίσκου αποτελεί θεμέλιο της λειτουργίας κάθε τραπεζικού οργανισμού, καθώς σχετίζεται άμεσα με την ασφάλεια των κεφαλαίων, τη σταθερότητα των λειτουργιών και τη φερεγγυότητα του ιδρύματος. Η εισαγωγή και αξιοποίηση των Big Data έχει αναμορφώσει ριζικά τις παραδοσιακές μεθόδους εκτίμησης και ελέγχου του κινδύνου, προσφέροντας τη δυνατότητα δημιουργίας πιο ακριβών, προσαρμοστικών και αυτοματοποιημένων μοντέλων πρόβλεψης.

Η βασική λειτουργία των τραπεζών είναι η διαχείριση κεφαλαίων μέσω δανεισμού, γεγονός που από τη φύση του εμπεριέχει ρίσκο. Το κύριο ερώτημα που καλείται να απαντήσει κάθε τραπεζικό ίδρυμα είναι κατά πόσο ένας υποψήφιος δανειολήπτης είναι ικανός να τηρήσει τις υποχρεώσεις του. Σε αυτό το σημείο, τα δεδομένα παίζουν κομβικό ρόλο, καθώς η ανάλυσή τους οδηγεί στην αξιολόγηση του προφίλ του πελάτη και της πιστοληπτικής του ικανότητας.

Μέσω τεχνικών εξόρυξης γνώσης, οι τράπεζες επεξεργάζονται ένα πλήθος χαρακτηριστικών, τόσο ποσοτικών όσο και ποιοτικών, που αφορούν τον πελάτη. Δημογραφικά στοιχεία, επαγγελματικό υπόβαθρο, παλαιότερη συναλλακτική συμπεριφορά, ύψος εισοδήματος, προϋπάρχοντες δανειακές υποχρεώσεις, καθώς και ιστορικό πληρωμών, συντίθενται σε ένα ενιαίο σύνολο δεδομένων το οποίο

χρησιμοποιείται για την εκτίμηση του κινδύνου αθέτησης. Η ανάλυση αυτή επιτρέπει τον υπολογισμό της πιστοληπτικής βαθμολογίας του πελάτη, γνωστής και ως credit score, η οποία καταλήγει να είναι το κριτήριο βάσει του οποίου εγκρίνεται ή απορρίπτεται μια αίτηση δανείου [14].

Πέρα από τις παραδοσιακές μορφές χρηματοδότησης, όπως είναι τα στεγαστικά ή καταναλωτικά δάνεια, η διαχείριση ρίσκου εφαρμόζεται και σε προϊόντα μικρότερης κλίμακας, όπως είναι οι πιστωτικές κάρτες και τα υπεραναλήψιμα όρια λογαριασμών. Σε αυτές τις περιπτώσεις, η τράπεζα καλείται να αξιολογήσει την οικονομική συμπεριφορά του πελάτη σχεδόν σε πραγματικό χρόνο, ώστε να καθορίσει αν θα αυξήσει ή θα μειώσει τα διαθέσιμα κεφάλαια προς αυτόν, με τρόπο που να εξασφαλίζει την αποφυγή αθέτησης.

Η συμβολή των Big Data είναι καθοριστική, διότι η παραδοσιακή ανάλυση κινδύνου στηρίζεται συχνά σε ιστορικά δεδομένα και βασικούς δείκτες (π.χ. σχέση χρέους/εισοδήματος), χωρίς να λαμβάνει υπόψη της δυναμικές πληροφορίες. Αντίθετα, τα σύγχρονα συστήματα ανάλυσης δεδομένων μπορούν να ενσωματώσουν πλήθος μεταβλητών, όπως π.χ. αλλαγές στο επάγγελμα ή το εισόδημα, ασυνήθιστες συναλλαγές, νέες συνδρομές σε υπηρεσίες, μέχρι και δεδομένα από τα κοινωνικά δίκτυα που υποδηλώνουν αυξημένο ή μειωμένο οικονομικό ρίσκο. Τα μοντέλα αυτά μπορούν να προσαρμόζονται διαρκώς, επαναξιολογώντας την κατάσταση του κάθε πελάτη και βελτιώνοντας τις αποφάσεις που αφορούν στην επέκταση ή τον περιορισμό του πιστωτικού κινδύνου.

Επιπλέον, η διαχείριση ρίσκου δεν περιορίζεται αποκλειστικά στον δανεισμό. Οι τράπεζες αξιολογούν διαρκώς κινδύνους που σχετίζονται με επενδυτικές κινήσεις, ρευστότητα κεφαλαίων, συμμόρφωση με κανονιστικά πλαίσια (compliance risk), αλλά και λειτουργικά σφάλματα ή εξωτερικές απειλές (όπως κυβερνοεπιθέσεις). Η ολοκληρωμένη διαχείριση αυτών των ρίσκων προϋποθέτει συστήματα που μπορούν να αντλούν, να συγκρίνουν και να αξιολογούν δεδομένα από πολλαπλές πηγές, με στόχο τη διασφάλιση της επιχειρησιακής συνέχειας και την προστασία του οργανισμού.

Η χρήση τεχνικών μηχανικής μάθησης και ανάλυσης προγνωστικών μοντέλων έχει οδηγήσει στην ανάπτυξη έξυπνων συστημάτων που όχι μόνο παρακολουθούν τις μεταβλητές που σχετίζονται με τον κίνδυνο, αλλά προσαρμόζονται σε νέες συνθήκες, μαθαίνοντας από την ιστορική συμπεριφορά των πελατών. Αυτό σημαίνει ότι ένα τραπεζικό ίδρυμα μπορεί όχι απλώς να αξιολογεί εάν ένας πελάτης πληροί τις προϋποθέσεις για χορήγηση δανείου σήμερα, αλλά να προβλέπει τη μακροχρόνια συνέπεια του, βοηθώντας στην καλύτερη αποτίμηση της αξίας του (lifetime value) και της μακροπρόθεσμης αποδοτικότητας του οργανισμού.

2.4. Αξία του Χρόνου Ζωής του Πελάτη (Customer Lifetime Value – CLV)

Η έννοια της αξίας του χρόνου ζωής του πελάτη (Customer Lifetime Value – CLV) έχει αποκτήσει τα τελευταία χρόνια στρατηγική σημασία για τις επιχειρήσεις και ιδιαίτερα για τον τραπεζικό τομέα. Ο δείκτης αυτός δεν περιορίζεται σε μια στατική απεικόνιση της οικονομικής αξίας ενός πελάτη τη δεδομένη χρονική στιγμή, αλλά επιδιώκει να εκτιμήσει την προβλεπόμενη αξία που θα προσδώσει ένας πελάτης σε έναν οργανισμό κατά τη διάρκεια ολόκληρης της σχέσης τους. Η εκτίμηση αυτή δεν είναι απλώς χρήσιμη – είναι καθοριστική για τον σχεδιασμό προϊόντων, την ορθολογική κατανομή των πόρων και την ανάπτυξη αποτελεσματικών στρατηγικών marketing και εξυπηρέτησης.

Στον χρηματοπιστωτικό χώρο, η έννοια του CLV αποκτά ιδιαίτερη βαρύτητα λόγω της φύσης των προϊόντων που προσφέρονται και της μακροχρόνιας σχέσης που οι τράπεζες επιδιώκουν να διατηρήσουν με τους πελάτες τους. Η λογική της μέτρησης της αξίας δεν περιορίζεται μόνο στα έσοδα που προκύπτουν από άμεσες χρεώσεις ή τόκους. Συμπεριλαμβάνει παράγοντες όπως η συχνότητα και το είδος των συναλλαγών, η υιοθέτηση νέων προϊόντων, η πιθανότητα σύστασης της τράπεζας σε τρίτους, αλλά και η συνέπεια στην αποπληρωμή των υποχρεώσεων.

Η παραδοσιακή προσέγγιση επιχειρούσε να αξιολογήσει την αξία ενός πελάτη με βάση το ύψος των καταθέσεων ή το υπόλοιπο του λογαριασμού. Ωστόσο, αυτή η οπτική αποδεικνύεται ανεπαρκής, καθώς αγνοεί άλλες παραμέτρους, όπως η προθυμία για χρήση περισσότερων προϊόντων, η αφοσίωση στον οργανισμό, ή ακόμα και η επιρροή του πελάτη στο κοινωνικό του περιβάλλον. Μέσω της ανάλυσης μεγάλων δεδομένων, οι τράπεζες μπορούν να δημιουργήσουν ένα πολύπλοκο και δυναμικό προφίλ του πελάτη, λαμβάνοντας υπόψη ένα σύνολο χαρακτηριστικών που ξεπερνά κατά πολύ τα παραδοσιακά οικονομικά κριτήρια.

Η αξία του CLV ενισχύεται και από το γεγονός ότι, σύμφωνα με σχετικές μελέτες, το κόστος απόκτησης ενός νέου πελάτη είναι πολλαπλάσιο σε σχέση με το κόστος διατήρησης ενός ήδη υπάρχοντος. Η διατήρηση μιας μακροχρόνιας και υγιούς σχέσης με τους πελάτες δεν αποτελεί μόνο ένδειξη αξιοπιστίας, αλλά και μοχλό αύξησης της κερδοφορίας. Παράλληλα, το CLV αποτελεί εργαλείο πρόβλεψης: οι πελάτες με χαμηλή αναμενόμενη αξία μπορούν να ταυτοποιηθούν και να αποκλειστούν από ενέργειες υψηλού κόστους, ενώ εκείνοι με υψηλό δυνητικό όφελος μπορούν να αποτελέσουν στόχο στοχευμένων καμπανιών, προσφορών ή επιβραβεύσεων.

Η εφαρμογή του δείκτη CLV γίνεται μέσα από μοντέλα μηχανικής μάθησης που συνδυάζουν ιστορικά δεδομένα συναλλαγών, δημογραφικά στοιχεία, συνήθειες χρήσης, πρότυπα συμπεριφοράς, και άλλες

παραμέτρους. Τα μοντέλα αυτά προβλέπουν την πιθανότητα μελλοντικής αλληλεπίδρασης με τον οργανισμό, την αποδοτικότητα αυτής της αλληλεπίδρασης, και τελικά τη συνολική οικονομική συνεισφορά του πελάτη. Η πληροφορία αυτή μπορεί να χρησιμοποιηθεί τόσο για τον σχεδιασμό εξατομικευμένων στρατηγικών marketing όσο και για την παραμετροποίηση των προσφερόμενων προϊόντων ή υπηρεσιών.

Είναι επίσης σημαντικό να τονιστεί ότι η στρατηγική αξιολόγησης του CLV συνδυάζει την ανάγκη για διατήρηση των υφιστάμενων πελατών με την πρόβλεψη και προσέλκυση νέων. Οι επιχειρήσεις μπορούν να επιδιώξουν τη διεύρυνση του πελατολογίου τους στοχεύοντας σε άτομα με χαρακτηριστικά αντίστοιχα αυτών που έχουν ήδη υψηλό CLV. Έτσι, δεν αντιμετωπίζουν την απόκτηση και τη διατήρηση ως διακριτές και ανταγωνιστικές στρατηγικές, αλλά ως δύο αλληλένδετα σκέλη μιας ενιαίας στρατηγικής ανάπτυξης.

2.5. Τμηματοποίηση Πελατών (Customer Segmentation)

Η τμηματοποίηση πελατών αποτελεί έναν από τους πλέον θεμελιώδεις πυλώνες στρατηγικής στην ψηφιακή εποχή της τραπεζικής, καθώς επιτρέπει στους οργανισμούς να κατανοήσουν εις βάθος τη βάση των πελατών τους, να εντοπίσουν εσωτερικές διαφοροποιήσεις και να διαμορφώσουν στοχευμένες δράσεις μάρκετινγκ, εξυπηρέτησης και ανάπτυξης προϊόντων. Η έννοια αυτή περιγράφει τη διαδικασία κατά την οποία οι πελάτες χωρίζονται σε ομάδες (segments), οι οποίες παρουσιάζουν κοινά χαρακτηριστικά ή συμπεριφορές, καθιστώντας δυνατή την εξατομικευμένη προσέγγιση, αντί της μαζικής και γενικευμένης επικοινωνίας.

Η ανάλυση μεγάλων δεδομένων προσφέρει στη διαδικασία αυτή ουσιαστική ενίσχυση, επιτρέποντας την αξιοποίηση ενός ευρύτατου φάσματος πληροφοριών για τη διαμόρφωση τμημάτων που είναι ακριβέστερα, δυναμικά και πιο λειτουργικά από ποτέ. Παραδοσιακά, οι τράπεζες στηρίζονταν σε χαρακτηριστικά όπως η ηλικία, το φύλο, η γεωγραφική τοποθεσία ή το εισόδημα για να ομαδοποιήσουν τους πελάτες τους. Ωστόσο, η σημερινή πραγματικότητα απαιτεί πιο εξελιγμένες και πολυδιάστατες μεθόδους, στις οποίες συνυπολογίζονται δημογραφικά στοιχεία, ιστορικά συναλλαγών, συμπεριφορικά μοτίβα, ψηφιακή αλληλεπίδραση, αλλά και δεδομένα από τρίτους φορείς ή συνεργαζόμενα οικοσυστήματα.

Η τμηματοποίηση είναι απαραίτητη για τη βελτίωση του επιπέδου εξυπηρέτησης και την αύξηση της αποτελεσματικότητας των επικοινωνιακών στρατηγικών. Για παράδειγμα, ένας πελάτης που εμφανίζει περιοδικά αυξημένες συναλλαγές και συχνή χρήση mobile banking θα μπορούσε να ενταχθεί σε ένα δυναμικό προφίλ που αντιστοιχεί σε "τεχνολογικά ενεργούς πελάτες", οι οποίοι ενδεχομένως είναι πιο

δεκτικοί σε νέες ψηφιακές υπηρεσίες ή αυτόματες επενδυτικές προτάσεις. Αντίστοιχα, πελάτες με υψηλό υπόλοιπο και χαμηλή δραστηριότητα μπορεί να αποτελέσουν στόχο για ενημερωτικές καμπάνιες σχετικές με επενδυτικά προϊόντα, ή για προγράμματα πιστότητας.

Ένα ιδιαίτερα ενδιαφέρον υποσύνολο της τμηματοποίησης είναι η λεγόμενη ψυχογραφική ανάλυση (psychographic segmentation). Αυτή δεν βασίζεται αποκλειστικά σε αντικειμενικά δεδομένα αλλά επιχειρεί να εντοπίσει τις στάσεις, τις αξίες, τις προθέσεις και τα ενδιαφέροντα των πελατών. Μέσω αυτής της προσέγγισης, η οποία αντλεί στοιχεία από τη συμπεριφορά στα ψηφιακά κανάλια, την κίνηση στο Internet banking ή ακόμη και τη γλώσσα που χρησιμοποιείται σε φόρμες επικοινωνίας, οι τράπεζες μπορούν να δημιουργήσουν πιο ολοκληρωμένα και συναισθηματικά συνδεδεμένα μοντέλα πελατών. Έτσι, πελάτες με διαφορετικά οικονομικά χαρακτηριστικά, αλλά παρόμοια κίνητρα ή στόχους, μπορεί να ενταχθούν στο ίδιο λειτουργικό τμήμα.

Η τμηματοποίηση δεν είναι απλώς εργαλείο ταξινόμησης αλλά λειτουργεί ως βάση για μια σειρά κρίσιμων στρατηγικών αποφάσεων. Επιτρέπει την καλύτερη πρόβλεψη της αποδοτικότητας συγκεκριμένων πελατών, την ανάπτυξη πιο στοχευμένων προϊόντων, τη μείωση του κινδύνου αποχώρησης πελατών (churn) και την αύξηση των διασταυρούμενων πωλήσεων (cross-selling). Επιπλέον, προσφέρει τη δυνατότητα εντοπισμού προβληματικών ή λιγότερο αποδοτικών ομάδων πελατών, επιτρέποντας την αποδοτικότερη κατανομή των διαθέσιμων πόρων υποστήριξης και επικοινωνίας.

Στο πλαίσιο του τραπεζικού τομέα, η ύπαρξη πολλών διαφορετικών τμημάτων πελατών επιτρέπει την εξατομίκευση της εμπειρίας χρήστη. Οι προωθητικές ενέργειες προσαρμόζονται στα ιδιαίτερα χαρακτηριστικά της κάθε ομάδας, οδηγώντας σε υψηλότερα ποσοστά ανταπόκρισης και αυξημένα επίπεδα ικανοποίησης. Παράλληλα, οι στρατηγικές πιστότητας και ανταμοιβής πελατών σχεδιάζονται με βάση τη μακροπρόθεσμη αξία του κάθε τμήματος, εξασφαλίζοντας ότι οι πελάτες που προσφέρουν τη μεγαλύτερη αξία λαμβάνουν και την ανάλογη προσοχή και υποστήριξη.

2.6. Συστήματα Συστάσεων (Recommendation Systems)

Τα συστήματα συστάσεων αποτελούν μία από τις πιο καινοτόμες και δυναμικές εφαρμογές των τεχνικών εξόρυξης γνώσης και μηχανικής μάθησης στο πλαίσιο της τραπεζικής τεχνολογίας. Πρόκειται για μηχανισμούς που σχεδιάζονται ώστε να προτείνουν στους χρήστες προϊόντα, υπηρεσίες ή ενέργειες που ταιριάζουν στο προφίλ, τις ανάγκες ή τις προτιμήσεις τους. Η χρήση τους δεν περιορίζεται πλέον σε πλατφόρμες ηλεκτρονικού εμπορίου και streaming, αλλά έχει βρει πρόσφορο έδαφος και στον

χρηματοοικονομικό τομέα, ενισχύοντας τη στοχευμένη εξυπηρέτηση και την προσωποποιημένη εμπειρία πελάτη.

Η βασική αρχή λειτουργίας των συστημάτων συστάσεων έγκειται στην ικανότητά τους να επεξεργάζονται μεγάλα και ετερογενή σύνολα δεδομένων που σχετίζονται με το ιστορικό του χρήστη, τις συναλλαγές του, την τοποθεσία, τις συνήθειες χρήσης ψηφιακών υπηρεσιών και τη συμπεριφορά παρόμοιων πελατών. Με την εφαρμογή κατάλληλων αλγορίθμων, όπως collaborative filtering, content-based filtering ή υβριδικών μοντέλων, τα συστήματα αυτά μπορούν να προβλέψουν με αξιοπιστία ποια προϊόντα είναι πιο πιθανό να ενδιαφέρουν τον χρήστη ή να του φανούν χρήσιμα σε συγκεκριμένες χρονικές στιγμές.

Στην τραπεζική, τα recommendation systems χρησιμοποιούνται ήδη σε διάφορες μορφές. Για παράδειγμα, μέσα από την ανάλυση του ιστορικού των συναλλαγών με πιστωτικές κάρτες, μπορούν να παραχθούν προσωποποιημένες προσφορές για αγορές από συνεργαζόμενα καταστήματα, λαμβάνοντας υπόψη τη γεωγραφική τοποθεσία του πελάτη ή τις κατηγορίες στις οποίες συχνά ξοδεύει. Επιπλέον, όταν ένας πελάτης έχει προφίλ που δείχνει ενδιαφέρον για επενδυτικά προϊόντα ή ασφαλιστικές υπηρεσίες, το σύστημα μπορεί να προτείνει αντίστοιχες λύσεις – είτε αυτόματα μέσω Internet ή mobile banking, είτε κατά την προσωπική επικοινωνία με σύμβουλο.

Η αξιοποίηση των recommendation engines δεν περιορίζεται στην ενίσχυση των πωλήσεων. Παίζουν σημαντικό ρόλο και στη βελτιστοποίηση της εμπειρίας χρήστη, καθώς μειώνουν τον χρόνο αναζήτησης και διευκολύνουν την απόφαση επιλογής ανάμεσα σε πολλά προϊόντα. Η δυνατότητα αυτή είναι ιδιαίτερως χρήσιμη σε πολύπλοκες ή υψηλού κόστους αποφάσεις, όπως η επιλογή στεγαστικού δανείου, όπου οι ανάγκες διαφέρουν σημαντικά από πελάτη σε πελάτη. Μέσα από τη συλλογή δεδομένων σχετικά με την εισοδηματική ικανότητα, τον οικογενειακό προγραμματισμό και τις γεωγραφικές προτιμήσεις, μπορούν να προταθούν στον πελάτη λύσεις που ικανοποιούν τόσο τα οικονομικά του δεδομένα όσο και τις προσωπικές του προτιμήσεις.

Αξιοσημείωτο είναι ότι τα συστήματα αυτά μπορούν να λειτουργήσουν και προληπτικά, αποτρέποντας, για παράδειγμα, την απόκτηση προϊόντων που ο χρήστης δεν χρειάζεται ή δεν μπορεί να υποστηρίξει. Με αυτόν τον τρόπο, η τράπεζα προστατεύει τόσο το συμφέρον του πελάτη όσο και το δικό της, διατηρώντας τη σχέση εμπιστοσύνης και μειώνοντας δυνητικούς πιστωτικούς κινδύνους.

Η υλοποίηση συστημάτων συστάσεων απαιτεί προηγμένες υποδομές ανάλυσης δεδομένων, αλλά και μηχανισμούς διαρκούς εκπαίδευσης των αλγορίθμων (model training) ώστε να μπορούν να ανταποκρίνονται στις αλλαγές της καταναλωτικής συμπεριφοράς. Η επιτυχία τους εξαρτάται σε μεγάλο

βαθμό από την ποιότητα και την επικαιρότητα των δεδομένων, καθώς και από την ετοιμότητα του οργανισμού να ενσωματώσει τις συστάσεις στις διαδικασίες λήψης αποφάσεων και εξυπηρέτησης πελατών.

Τέλος, είναι σημαντικό να διασφαλίζεται η διαφάνεια, η προστασία των προσωπικών δεδομένων και η συμμόρφωση με τους κανονισμούς, καθώς τα συστήματα αυτά λειτουργούν στη βάση ευαίσθητης πληροφορίας. Ο σεβασμός στην ιδιωτικότητα των πελατών είναι κρίσιμος για την αποδοχή και τη χρήση αυτών των τεχνολογιών από το ευρύ κοινό.

3. Εξόρυξη Δεδομένων (Data Mining)

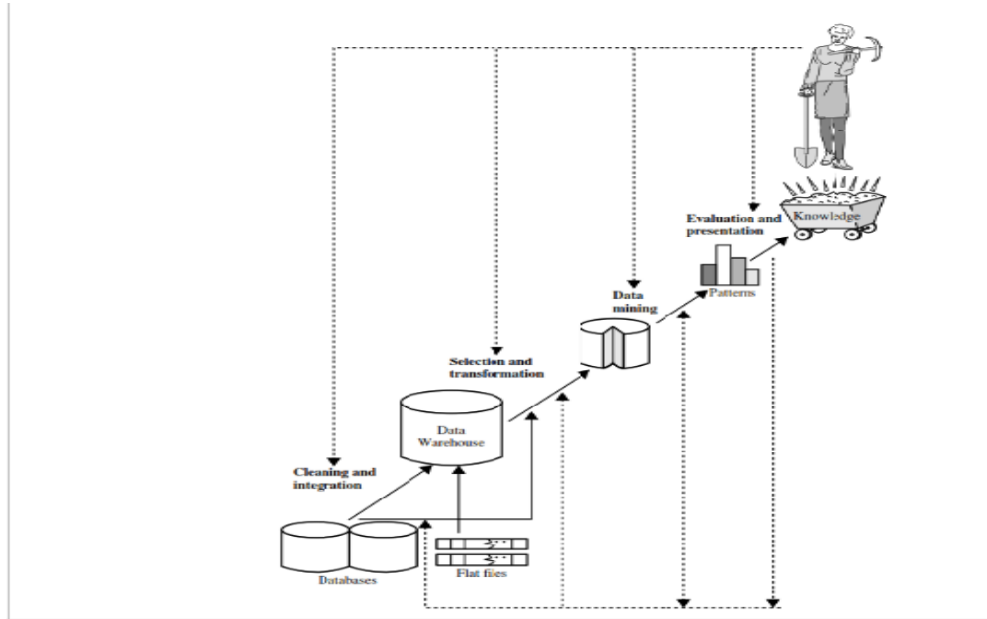
Η εξόρυξη δεδομένων (Data Mining), ή αλλιώς ανακάλυψη γνώσης από δεδομένα (Knowledge Discovery from Data – KDD), αποτελεί ένα διεπιστημονικό πεδίο που συνδυάζει στοιχεία από τη στατιστική, την τεχνητή νοημοσύνη, τη μηχανική μάθηση και τις βάσεις δεδομένων, με σκοπό την αποκάλυψη χρήσιμων και μη εμφανών πληροφοριών από μεγάλα σύνολα δεδομένων. Σε έναν κόσμο όπου η παραγωγή δεδομένων είναι διαρκής και εκθετική, η ικανότητα μετατροπής αυτής της πληροφορίας σε γνώση συνιστά συγκριτικό πλεονέκτημα για οργανισμούς, επιχειρήσεις και ερευνητές.

Η ανάγκη για αποτελεσματική ανάλυση δεδομένων προέκυψε από την αδυναμία των παραδοσιακών μεθόδων στατιστικής να ανταπεξέλθουν σε δεδομένα με υψηλή διαστατικότητα, ποικιλία και όγκο. Η εξόρυξη δεδομένων εστιάζει όχι μόνο στην ανάλυση των δεδομένων αλλά και στην αυτοματοποιημένη αναγνώριση προτύπων, σχέσεων και τάσεων, οι οποίες μπορούν να αξιοποιηθούν για σκοπούς πρόβλεψης, ταξινόμησης, συσταδοποίησης, σύστασης ή περιγραφής.

Η τυπική διαδικασία εξόρυξης γνώσης περιλαμβάνει τα ακόλουθα στάδια:

- Καθαρισμός δεδομένων (data cleaning), για τη διαχείριση ελλιπών ή θορυβωδών εγγραφών.
- Ολοκλήρωση και ενοποίηση δεδομένων από ετερογενείς πηγές.
- Επιλογή και μετασχηματισμός χαρακτηριστικών, ώστε να εξαχθούν μόνο τα σημαντικά και να μετατραπούν σε κατάλληλη μορφή για ανάλυση.
- Εφαρμογή αλγορίθμων εξόρυξης, όπως τεχνικές ταξινόμησης, συσταδοποίησης ή συσχετίσεων.

- Ερμηνεία και οπτικοποίηση αποτελεσμάτων, με στόχο την εξαγωγή συμπερασμάτων και τη λήψη αποφάσεων.



Εικόνα 2: Διαγραμματική απεικόνιση της διαδικασίας εξόρυξης γνώσης (KDD), από το αρχικό σύνολο δεδομένων έως την παραγωγή πληροφορίας. Πηγή: Morgan Kaufmann (Han et al., *Data Mining: Concepts and Techniques*).

Οι πιο διαδεδομένες τεχνικές εξόρυξης περιλαμβάνουν την εποπτευόμενη μάθηση (supervised learning), όπως η ταξινόμηση και η παλινδρόμηση· τη μη εποπτευόμενη μάθηση (unsupervised learning), όπως η συσταδοποίηση και η ανάλυση συσχετίσεων· καθώς και υβριδικές μεθόδους, που συνδυάζουν διαφορετικά μοντέλα και τεχνικές για βελτιωμένη απόδοση.

Η εξόρυξη δεδομένων βρίσκει εφαρμογή σε ένα ευρύ φάσμα τομέων, όπως:

- Οικονομικά και τραπεζική: για την ανίχνευση απάτης, την πιστοληπτική αξιολόγηση και το στοχευμένο μάρκετινγκ.
- Υγεία: για την πρόβλεψη νοσημάτων και την προσωποποιημένη ιατρική.
- Λιαν εμπόριο: για την ανάλυση καλαθιού αγορών και τη διαχείριση αποθεμάτων.
- Κοινωνικά δίκτυα και τεχνολογίες περιεχομένου: για την κατανόηση συμπεριφοράς χρηστών και τη δημιουργία συστημάτων συστάσεων.

Η αξία της εξόρυξης γνώσης δεν περιορίζεται στην τεχνική πρόβλεψη· ενισχύει τη στρατηγική σκέψη, αυτοματοποιεί κρίσιμες διαδικασίες και ενδυναμώνει τη λήψη αποφάσεων με δεδομένα και όχι με εικασίες.

3.1. Στάδια της Διαδικασίας KDD

Η διαδικασία της εξόρυξης γνώσης από δεδομένα (Knowledge Discovery in Databases – KDD) δεν συνίσταται απλώς στην εφαρμογή ενός αλγορίθμου, αλλά περιλαμβάνει μια σειρά από αλληλοεξαρτώμενα βήματα, τα οποία οδηγούν από την ακατέργαστη πληροφορία στην παραγωγή τεκμηριωμένης γνώσης. Παρότι κατά την εισαγωγή του κεφαλαίου έγινε μια συνοπτική παρουσίαση του πλαισίου της KDD, εδώ θα επιχειρηθεί μια πιο εστιασμένη ανάλυση των επιμέρους σταδίων, αναδεικνύοντας τις τεχνικές προσεγγίσεις, τα εμπόδια που συναντώνται στην πράξη, καθώς και τη σημασία κάθε φάσης για την επιτυχία της συνολικής διαδικασίας.

Το πρώτο βήμα, η επιλογή και συλλογή των δεδομένων, αφορά τον προσδιορισμό των πηγών πληροφορίας που είναι σχετικές με το υπό μελέτη φαινόμενο. Αν και μπορεί να φαίνεται τυπικό, στην πραγματικότητα πρόκειται για ένα κρίσιμο στάδιο, καθώς από την ποιότητα και την καταλληλότητα των δεδομένων αυτών εξαρτάται η εγκυρότητα ολόκληρης της διαδικασίας. Στον χρηματοοικονομικό τομέα, για παράδειγμα, δεδομένα μπορούν να προέρχονται από συναλλαγές, ερωτηματολόγια, ιστορικά πιστώσεων ή ακόμα και από εξωτερικές βάσεις δεδομένων, όπως δημογραφικά μητρώα ή κοινωνικά δίκτυα.

Ακολουθεί ο καθαρισμός των δεδομένων, ένα στάδιο που παραδοσιακά απαιτεί σημαντικό χρόνο και εξειδίκευση. Τα ακατέργαστα δεδομένα συνήθως περιέχουν ελλείψεις, αντιφάσεις, ακραίες τιμές ή και λανθασμένες εγγραφές. Οι στρατηγικές που εφαρμόζονται περιλαμβάνουν είτε τη διόρθωση (όταν είναι δυνατή), είτε τη διαγραφή προβληματικών παρατηρήσεων, είτε τη συμπλήρωση τιμών βάσει μεθόδων παρεμβολής ή στατιστικής εκτίμησης. Στην πράξη, η σωστή εκτέλεση αυτής της φάσης μπορεί να βελτιώσει θεαματικά την ακρίβεια ενός μοντέλου, ιδιαίτερα σε προβλήματα όπου τα δεδομένα συλλέγονται αυτόματα και με χαμηλή εποπτεία.

Η ενοποίηση και ο μετασχηματισμός δεδομένων αφορά την ενοποίηση διαφορετικών πινάκων ή αρχείων σε μια ενιαία βάση, καθώς και την προετοιμασία των μεταβλητών προς ανάλυση. Σε αυτό το σημείο, μπορεί να εφαρμοστούν τεχνικές κανονικοποίησης (όπως min-max scaling ή standardization), μετατροπής κατηγορικών μεταβλητών σε δυαδικές (one-hot encoding), χειρισμού μεταβλητών υψηλής διαστατικότητας ή ακόμα και τεχνικές μείωσης διαστάσεων, όπως PCA. Η επιλογή τεχνικής εξαρτάται

τόσο από τη φύση των δεδομένων όσο και από τις απαιτήσεις του μοντέλου που πρόκειται να χρησιμοποιηθεί. Στα σύγχρονα εργαλεία μηχανικής μάθησης, η προεπεξεργασία αυτή συνήθως ενσωματώνεται σε pipelines, διασφαλίζοντας συνέπεια στην αναπαραγωγή των πειραμάτων.

Το στάδιο της εξόρυξης δεδομένων, που αποτελεί την «καρδιά» της διαδικασίας, περιλαμβάνει την εφαρμογή των κατάλληλων αλγορίθμων για την ανακάλυψη προτύπων, σχέσεων ή μοντέλων πρόβλεψης. Εδώ γίνεται η διάκριση μεταξύ επιβλεπόμενων και μη επιβλεπόμενων τεχνικών, αναλόγως αν υπάρχει διαθέσιμη ετικέτα-στόχος. Η επιλογή του κατάλληλου αλγορίθμου (όπως δέντρα απόφασης, λογιστική παλινδρόμηση, δίκτυα νευρώνων ή συστάδες τύπου k-means) επηρεάζεται από παράγοντες όπως η φύση των δεδομένων, το μέγεθος του συνόλου, η ερμηνευσιμότητα του αποτελέσματος και η υπολογιστική αποδοτικότητα. Ιδιαίτερη προσοχή πρέπει να δίνεται στην αποφυγή υπερπροσαρμογής (overfitting), ενώ συχνά απαιτείται η χρήση τεχνικών διασταυρούμενης επικύρωσης (cross-validation) ή η βελτιστοποίηση υπερπαραμέτρων μέσω grid search ή Bayesian optimization.

Η ερμηνεία και αξιολόγηση των αποτελεσμάτων είναι ένα στάδιο το οποίο συνδέει την τεχνική ανάλυση με τη λήψη αποφάσεων. Τα μοντέλα πρέπει να αξιολογούνται όχι μόνο ως προς την απόδοσή τους (μέσω μετρικών όπως accuracy, precision, recall, F1-score ή ROC-AUC), αλλά και ως προς τη λειτουργικότητα και σημασία τους στο εκάστοτε επιχειρησιακό ή ερευνητικό πλαίσιο. Για παράδειγμα, σε μια εφαρμογή πρόβλεψης αθέτησης πληρωμών, ενδέχεται να έχει μεγαλύτερη σημασία το recall για την κατηγορία των πελατών υψηλού κινδύνου, παρά η συνολική ακρίβεια. Επιπλέον, μοντέλα όπως το Random Forest ή το Gradient Boosting παρέχουν χρήσιμα εργαλεία για την εξαγωγή «σημαντικότητας χαρακτηριστικών» (feature importance), ενισχύοντας τη διαφάνεια της ανάλυσης.

Το τελικό στάδιο της διαδικασίας είναι η παρουσίαση και ενσωμάτωση της γνώσης. Αφορά τη μετάφραση των τεχνικών αποτελεσμάτων σε μορφές που είναι κατανοητές και αξιοποιήσιμες από τους τελικούς αποδέκτες. Αυτό μπορεί να περιλαμβάνει διαγράμματα, πίνακες, συνοπτικές αναφορές ή και διαδραστικές εφαρμογές που επιτρέπουν την εξερεύνηση των αποτελεσμάτων από μη ειδικούς. Η αποτυχία στη σωστή παρουσίαση συχνά συνεπάγεται και αποτυχία στη λήψη σωστών αποφάσεων, ανεξαρτήτως της ποιότητας του μοντέλου.

Η διαδικασία KDD, με τα διακριτά αλλά αλληλοσυμπληρούμενα στάδιά της, επιτρέπει τη δομημένη και μεθοδική αξιοποίηση δεδομένων σε πολύπλοκα περιβάλλοντα. Η επιτυχία της εξαρτάται όχι μόνο από την τεχνική αρτιότητα των επιμέρους σταδίων, αλλά και από τη συνοχή, τη στρατηγική σκέψη και τη συνεχή αξιολόγηση της καταλληλότητας των μεθόδων στο πλαίσιο εφαρμογής. Επομένως, η εξόρυξη

γνώσης συνιστά ένα εργαλείο ευφυούς ανάλυσης, που συνδέει την τεχνολογία με τη λήψη τεκμηριωμένων αποφάσεων.

3.2. Τεχνικές Μάθησης στην Εξόρυξη Δεδομένων

Η μηχανική μάθηση (machine learning) αποτελεί έναν από τους βασικούς μηχανισμούς με τους οποίους η εξόρυξη γνώσης καθίσταται εφικτή. Μέσω της μηχανικής μάθησης, οι υπολογιστές αποκτούν την ικανότητα να αναγνωρίζουν πρότυπα, να λαμβάνουν αποφάσεις ή να προβλέπουν μελλοντικές καταστάσεις βασιζόμενοι αποκλειστικά σε δεδομένα και όχι σε ρητά προγραμματισμένους κανόνες. Εντασσόμενη στο ευρύτερο πλαίσιο της τεχνητής νοημοσύνης, η μηχανική μάθηση λειτουργεί ως κινητήριος δύναμη της εξόρυξης δεδομένων, ειδικά όταν οι σχέσεις ανάμεσα στις μεταβλητές δεν είναι άμεσα ορατές ή όταν ο όγκος και η πολυπλοκότητα των δεδομένων υπερβαίνει τις δυνατότητες παραδοσιακής ανάλυσης.

Η βασική αρχή πίσω από τη μηχανική μάθηση είναι η ανάπτυξη μοντέλων που μπορούν να «μάθουν» από παρατηρούμενα δεδομένα, δηλαδή να αναγνωρίζουν συστηματικά μοτίβα και σχέσεις και να τα αξιοποιούν για να κάνουν προβλέψεις ή να γενικεύουν σε νέα, άγνωστα δεδομένα. Τα μοντέλα αυτά εκπαιδεύονται, αξιολογούνται και στη συνέχεια χρησιμοποιούνται σε πραγματικά σενάρια, από την αναγνώριση εικόνων και φωνής, έως την πρόβλεψη χρηματοοικονομικής συμπεριφοράς ή την ανίχνευση απάτης. Η επιτυχία της μηχανικής μάθησης εξαρτάται σε μεγάλο βαθμό από την ποιότητα των δεδομένων, την επιλογή των κατάλληλων χαρακτηριστικών, και τη φύση του εκάστοτε προβλήματος.

Αναλόγως της διαθεσιμότητας εποπτευόμενης πληροφόρησης κατά τη φάση εκπαίδευσης, η μηχανική μάθηση διακρίνεται σε τρεις κύριες κατηγορίες: την επιβλεπόμενη μάθηση (supervised learning), τη μη επιβλεπόμενη μάθηση (unsupervised learning) και την ενισχυτική μάθηση (reinforcement learning). Αυτές οι προσεγγίσεις δεν ανταγωνίζονται μεταξύ τους αλλά απευθύνονται σε διαφορετικού τύπου προβλήματα, διαμορφώνοντας ένα πλούσιο φάσμα τεχνικών που καλύπτουν το σύνολο σχεδόν των αναγκών ανάλυσης δεδομένων.

Η επιβλεπόμενη μάθηση είναι η πιο ευρέως χρησιμοποιούμενη προσέγγιση. Βασίζεται στην ύπαρξη μιας εξαρτημένης μεταβλητής ή στόχου, της οποίας οι τιμές είναι γνωστές στο σύνολο εκπαίδευσης. Ο αλγόριθμος προσπαθεί να μάθει τη σχέση μεταξύ των εισόδων (χαρακτηριστικών) και της εξόδου (στόχου), έτσι ώστε να μπορεί να προβλέψει τη σωστή τιμή σε νέα δεδομένα. Η διαδικασία αυτή μοιάζει με τον παραδοσιακό στατιστικό συσχετισμό, με τη διαφορά ότι εδώ δίνεται έμφαση στην προγνωστική ισχύ και στη γενίκευση του μοντέλου. Οι πιο χαρακτηριστικές εφαρμογές της επιβλεπόμενης μάθησης

είναι η ταξινόμηση (classification), όπου προβλέπεται μια διακριτή κατηγορία, και η παλινδρόμηση (regression), όπου η έξοδος είναι συνεχής τιμή.

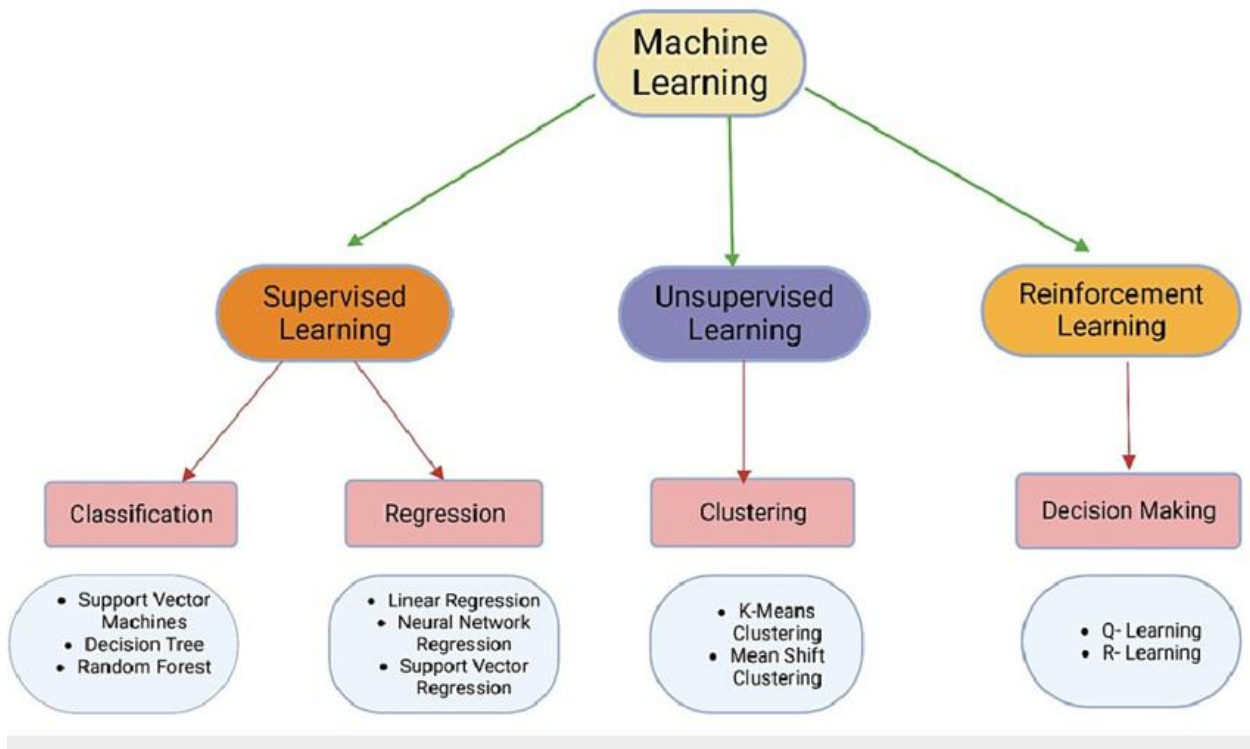
Η ταξινόμηση, για παράδειγμα, εφαρμόζεται σε περιπτώσεις όπου επιδιώκεται η κατηγοριοποίηση ενός αντικειμένου ή γεγονότος, όπως η αξιολόγηση ενός πελάτη ως «υψηλού» ή «χαμηλού» πιστωτικού κινδύνου, ή η πρόβλεψη εάν μια συναλλαγή είναι κανονική ή ύποπτη. Από την άλλη, η παλινδρόμηση χρησιμοποιείται για προβλήματα όπως η πρόβλεψη της αξίας ενός περιουσιακού στοιχείου ή της ζήτησης σε ένα προϊόν, στηριζόμενη σε ιστορικά δεδομένα. Στην επιβλεπόμενη μάθηση, η αξιολόγηση των μοντέλων γίνεται συνήθως μέσω δεικτών όπως το ποσοστό ορθών προβλέψεων (accuracy), ο δείκτης F1, και η περιοχή κάτω από την καμπύλη ROC (AUC-ROC), που επιτρέπει την εκτίμηση της ισορροπίας μεταξύ θετικών και αρνητικών προβλέψεων.

Σε αντιδιαστολή, η μη επιβλεπόμενη μάθηση εφαρμόζεται όταν τα δεδομένα δεν συνοδεύονται από μεταβλητή στόχο. Σε αυτές τις περιπτώσεις, το ζητούμενο είναι να εντοπιστούν υποκείμενα πρότυπα ή δομές χωρίς καθοδήγηση. Πρόκειται για ιδιαίτερα χρήσιμη τεχνική σε διερευνητικές αναλύσεις ή όταν δεν υπάρχουν εκ των προτέρων γνωστές ετικέτες. Η συσταδοποίηση (clustering) αποτελεί την πιο αναγνωρίσιμη μορφή μη επιβλεπόμενης μάθησης, όπου επιδιώκεται η ομαδοποίηση των παρατηρήσεων σε κλάσεις με υψηλή ενδο-ομοιογένεια και χαμηλή δια-ομοιογένεια. Οι τεχνικές αυτές είναι εξαιρετικά χρήσιμες σε περιβάλλοντα όπου απαιτείται ανακάλυψη τμημάτων πελατών (customer segmentation), ανάλυση συμπεριφοράς ή εύρεση ασυνήθιστων προτύπων (anomaly detection).

Η εξόρυξη κανόνων συσχέτισης, μια επίσης μη επιβλεπόμενη μέθοδος, χρησιμοποιείται ευρέως στην ανάλυση συναλλακτικών δεδομένων για την εξαγωγή σχέσεων της μορφής «αν A τότε B», όπως για παράδειγμα «αν ένας πελάτης αγοράσει γάλα, υπάρχει πιθανότητα να αγοράσει και δημητριακά». Οι τεχνικές αυτές είναι ιδιαίτερα χρήσιμες στο λιανεμπόριο και στις online πλατφόρμες για σκοπούς σύστασης προϊόντων και αύξησης διασταυρούμενων πωλήσεων.

Η ενισχυτική μάθηση αποτελεί μια πιο εξειδικευμένη μορφή μάθησης, στην οποία το μοντέλο (ή πράκτορας) αλληλεπιδρά με το περιβάλλον και βελτιώνει τη συμπεριφορά του μέσα από ακολουθία ενεργειών και αποτελεσμάτων. Αντί να μαθαίνει από στατικά δεδομένα, αξιολογεί τη χρησιμότητα των ενεργειών του βάσει της «ανταμοιβής» που λαμβάνει. Αν και λιγότερο διαδεδομένη στην παραδοσιακή εξόρυξη δεδομένων, η ενισχυτική μάθηση έχει αξιοσημείωτη ανάπτυξη σε πεδία όπως η ρομποτική, τα συστήματα σύστασης σε πραγματικό χρόνο και η χρηματιστηριακή διαπραγμάτευση (algorithmic

trading).



Εικόνα 3: Κατηγοριοποίηση της Μηχανικής Μάθησης με βάση το είδος της μάθησης και τους τύπους εφαρμογών, Πηγή: ResearchGate.

Η επιλογή μεταξύ αυτών των τύπων μάθησης εξαρτάται από τη φύση του προβλήματος, τη διαθεσιμότητα ετικετών, το επίπεδο πολυπλοκότητας των δεδομένων και τον τελικό στόχο του αναλυτή ή του οργανισμού. Η τεχνολογική εξέλιξη, σε συνδυασμό με την αυξανόμενη διαθεσιμότητα δεδομένων, ενισχύει τη χρήση και τον συνδυασμό των τεχνικών αυτών με τρόπους που επιτρέπουν τη δυναμική και προσαρμοστική ανάλυση σε πλήθος εφαρμογών.

3.3. Μοντέλα Μηχανικής Μάθησης

Καθώς οι εφαρμογές της μηχανικής μάθησης εξαπλώνονται σε ένα συνεχώς διευρυνόμενο φάσμα πεδίων, η ανάγκη επιλογής κατάλληλων αλγορίθμων ανάλογα με τη φύση του προβλήματος καθίσταται κρίσιμη. Η ποικιλομορφία των διαθέσιμων τεχνικών, οι διαφορές ως προς τις υποθέσεις, τη δομή, και τον τρόπο εκμάθησης, καθιστούν απαραίτητη την κατανόηση των βασικών αρχών που διέπουν κάθε μοντέλο. Σε ένα περιβάλλον όπου η πολυπλοκότητα των δεδομένων αυξάνεται και η ακρίβεια των προβλέψεων αποκτά επιχειρησιακή σημασία, η επιλογή της κατάλληλης μεθοδολογίας δεν αποτελεί απλώς τεχνική

απόφαση, αλλά κομβικό σημείο για την αξιοπιστία και τη χρηστικότητα των αποτελεσμάτων. Για τον λόγο αυτό, αξίζει να εξεταστούν ορισμένες από τις πιο διαδεδομένες προσεγγίσεις, ώστε να αναδειχθούν τα βασικά τους χαρακτηριστικά, τα πλεονεκτήματα, καθώς και οι περιορισμοί τους.

3.3.1. Random Forest Classifier

Ο αλγόριθμος Τυχαίων Δασών (Random Forest) αποτελεί μια από τις πιο ισχυρές τεχνικές μηχανικής μάθησης για προβλήματα ταξινόμησης και παλινδρόμησης, βασισμένος στη φιλοσοφία των μεθόδων συνόλου (ensemble learning). Πρόκειται για μια τεχνική που συνδυάζει την απόδοση πολλαπλών δέντρων απόφασης (decision trees), με στόχο τη βελτίωση της γενικευσιμότητας και την αύξηση της ακρίβειας των προβλέψεων.

Η βασική ιδέα πίσω από τον Random Forest είναι η κατασκευή ενός μεγάλου αριθμού δέντρων απόφασης, καθένα από τα οποία εκπαιδεύεται σε διαφορετικό τυχαίο υποσύνολο του συνόλου εκπαίδευσης (μέσω bootstrap sampling) [3]. Κατά τη διαδικασία κατασκευής κάθε δέντρου, σε κάθε κόμβο του δέντρου δεν εξετάζονται όλα τα διαθέσιμα χαρακτηριστικά για τον διαχωρισμό, αλλά επιλέγεται τυχαία ένα μικρότερο υποσύνολο. Αυτός ο συνδυασμός δειγματοληψίας τόσο σε επίπεδο παρατηρήσεων όσο και χαρακτηριστικών έχει ως αποτέλεσμα τη δημιουργία ποικιλίας (diversity) μεταξύ των δέντρων, κάτι που μειώνει τη συσχέτιση μεταξύ τους και ενισχύει την ακρίβεια του συνολικού μοντέλου.

Η διαδικασία εκπαίδευσης του μοντέλου δεν απαιτεί κλάδεμα των δέντρων, καθώς η αυξημένη διασπορά ελέγχεται από τη συλλογική ψήφο του δάσους. Για προβλήματα ταξινόμησης, η τελική απόφαση προκύπτει από την πλειοψηφική ψήφο των δέντρων (majority voting), ενώ για προβλήματα παλινδρόμησης λαμβάνεται υπόψη ο μέσος όρος των επιμέρους προβλέψεων.

Ένα από τα βασικά πλεονεκτήματα του Random Forest είναι η ανθεκτικότητά του στην υπερπροσαρμογή (overfitting), που συχνά παρατηρείται σε μεμονωμένα δέντρα απόφασης. Επιπλέον, το μοντέλο προσφέρει υψηλή ακρίβεια χωρίς να απαιτείται έντονη ρύθμιση υπερπαραμέτρων και μπορεί να διαχειριστεί δεδομένα με υψηλή διαστατικότητα ή ελλείψεις. Ένα ακόμη ιδιαίτερα χρήσιμο χαρακτηριστικό είναι η δυνατότητα εκτίμησης της σχετικής σημασίας των χαρακτηριστικών (feature importance), γεγονός που διευκολύνει την κατανόηση της συμβολής κάθε μεταβλητής στην πρόβλεψη.

Παρά τα πλεονεκτήματά του, ο Random Forest παρουσιάζει και ορισμένους περιορισμούς. Η αυξημένη υπολογιστική απαίτηση κατά την εκπαίδευση και την πρόβλεψη αποτελεί έναν από αυτούς, ειδικά όταν ο αριθμός των δέντρων ή το μέγεθος των δεδομένων είναι μεγάλο. Επίσης, παρά το γεγονός ότι αποτελεί

πιο διαφανή επιλογή σε σχέση με πιο σύνθετες τεχνικές, η ερμηνευσιμότητα του μοντέλου μειώνεται λόγω της φύσης του ως σύνολο πολλών δέντρων, γεγονός που καθιστά δύσκολη την παρακολούθηση της λογικής λήψης απόφασης.

Η απόδοση του αλγορίθμου επηρεάζεται σημαντικά από τις υπερπαραμέτρους του. Ο αριθμός των δέντρων (`n_estimators`), το πλήθος χαρακτηριστικών που εξετάζονται σε κάθε κόμβο (`max_features`), το μέγιστο επιτρεπόμενο βάθος των δέντρων (`max_depth`), καθώς και οι ελάχιστες απαιτήσεις σε δείγματα για διαχωρισμό (`min_samples_split`) ή φύλλωμα (`min_samples_leaf`) καθορίζουν την ισορροπία μεταξύ πολυπλοκότητας, ταχύτητας και απόδοσης. Η ρύθμιση αυτών των παραμέτρων μπορεί να πραγματοποιηθεί με τεχνικές όπως το Grid Search ή το Randomized Search, σε συνδυασμό με διασταυρούμενη επικύρωση (`cross-validation`), προκειμένου να βρεθεί ο συνδυασμός που μεγιστοποιεί την ακρίβεια της πρόβλεψης. [4]

Συνολικά, ο Random Forest αποτελεί μια ισχυρή και σταθερή επιλογή για προβλήματα εποπτευόμενης μάθησης, προσφέροντας ισορροπία μεταξύ ακρίβειας και αξιοπιστίας, χωρίς να απαιτεί υπερβολική τεχνική πολυπλοκότητα ή έντονη ρύθμιση, ενώ ταυτόχρονα παραμένει ανθεκτικός σε θόρυβο και εξαιρετικά ευέλικτος στην προσαρμογή σε διαφορετικά χαρακτηριστικά των δεδομένων.

3.3.2. Logistic Regression

Η Logistic Regression είναι ένα από τα πιο βασικά αλλά και ιδιαίτερα χρήσιμα μοντέλα μηχανικής μάθησης για προβλήματα δυαδικής ταξινόμησης. Αν και ονομάζεται «παλινδρόμηση», στην πραγματικότητα πρόκειται για μέθοδο ταξινόμησης, καθώς δεν προβλέπει συνεχή τιμή, αλλά πιθανότητα για την ένταξη ενός δείγματος σε μία από δύο διακριτές κατηγορίες.

Η κεντρική ιδέα του μοντέλου είναι η χρήση ενός γραμμικού συνδυασμού των εισόδων και η μετατροπή του σε πιθανότητα μέσω της συνάρτησης sigmoid. Συγκεκριμένα, ο γραμμικός συνδυασμός ορίζεται ως:

$$z = w^T x + b$$

όπου:

- x είναι το διάνυσμα των χαρακτηριστικών,
- w είναι τα αντίστοιχα βάρη,
- b είναι η προκατάληψη (bias).

Η συνάρτηση sigmoid εφαρμόζεται στη μεταβλητή z και δίνει ως έξοδο την πιθανότητα το δείγμα να ανήκει στην κλάση 1:

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

Το αποτέλεσμα της sigmoid είναι ένας αριθμός μεταξύ 0 και 1. Αν αυτός ξεπερνά ένα καθορισμένο κατώφλι (συνήθως 0.5), το δείγμα ταξινομείται στην κατηγορία 1· αλλιώς, στην κατηγορία 0.

Η εκπαίδευση του μοντέλου πραγματοποιείται μέσω της ελαχιστοποίησης της λογιστικής απώλειας (log-loss), η οποία «τιμωρεί» τις λανθασμένες προβλέψεις με έναν βαθμό ανάλογο του πόσο μακριά βρίσκονται από την αληθινή κλάση. Ο τύπος της log-loss δίνεται ως:

$$L = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

όπου:

- N είναι ο συνολικός αριθμός των δειγμάτων,
- y_i είναι η πραγματική κλάση του δείγματος i ,
- \hat{y}_i η προβλεπόμενη πιθανότητα από το μοντέλο.

Η βελτιστοποίηση της log-loss γίνεται συνήθως μέσω παραλλαγών του αλγορίθμου gradient descent, επιτρέποντας την εύρεση των βέλτιστων παραμέτρων w και b .

Το μοντέλο της λογιστικής παλινδρόμησης είναι απλό και υπολογιστικά αποδοτικό. Προσφέρει σημαντικά πλεονεκτήματα όταν οι συσχετίσεις μεταξύ των χαρακτηριστικών και της εξαρτημένης μεταβλητής είναι γραμμικές, και όταν τα δεδομένα είναι ανεξάρτητα και χωρίς ισχυρές αλληλεπιδράσεις. Ένα ακόμα σημαντικό πλεονέκτημα είναι η ερμηνευσιμότητα: οι συντελεστές του μοντέλου επιτρέπουν την κατανόηση της επίδρασης κάθε χαρακτηριστικού στην πιθανότητα πρόβλεψης [5], διευκολύνοντας έτσι την ανάλυση των αποτελεσμάτων, ειδικά σε τομείς όπου η διαφάνεια είναι σημαντική (όπως η οικονομία ή η ιατρική).

Ωστόσο, η λογιστική παλινδρόμηση εμφανίζει περιορισμούς. Το μοντέλο βασίζεται στην υπόθεση γραμμικότητας μεταξύ των εξαρτημένων και ανεξάρτητων μεταβλητών, κάτι που δεν ισχύει σε πιο σύνθετες περιπτώσεις. Επίσης, είναι ευαίσθητο σε πολυσυγγραμμικότητα (multicollinearity) και δεν

αποδίδει καλά σε προβλήματα όπου οι κατηγορίες δεν είναι διαχωρίσιμες γραμμικά. Σε τέτοιες περιπτώσεις, ενδέχεται να απαιτείται η χρήση μετασχηματισμών των χαρακτηριστικών ή η προσθήκη πολυωνυμικών όρων.

Η απόδοση του μοντέλου μπορεί να ενισχυθεί με την ενσωμάτωση τακτικών κανονικοποίησης (regularization), όπως η L1 (Lasso) ή L2 (Ridge), που περιορίζουν την πολυπλοκότητα του μοντέλου και αποτρέπουν την υπερπροσαρμογή.

Παρά τους περιορισμούς της, η λογιστική παλινδρόμηση παραμένει ένα πολύτιμο εργαλείο, ειδικά όταν η σαφήνεια και η απλότητα ερμηνείας αποτελούν προτεραιότητα, ή όταν τα δεδομένα πληρούν τις βασικές προϋποθέσεις γραμμικότητας και ισορροπημένων κλάσεων. [5]

3.3.3. Gradient Boost Classifier

Ο Gradient Boosting Classifier είναι ένας από τους πιο ισχυρούς αλγόριθμους μηχανικής μάθησης, ιδιαίτερα χρήσιμος για προβλήματα ταξινόμησης και παλινδρόμησης. Ανήκει στην κατηγορία των ensemble methods και βασίζεται στη μέθοδο boosting, η οποία δημιουργεί ένα σύνολο αδύναμων μοντέλων (συνήθως δέντρων απόφασης) και τα συνδυάζει με τρόπο που μεγιστοποιεί την ακρίβεια της πρόβλεψης. Σε αντίθεση με άλλες τεχνικές ensemble, όπως το bagging, η μέθοδος boosting λειτουργεί με αλληλοδιαδοχικά βήματα, όπου κάθε νέο μοντέλο προσπαθεί να διορθώσει τα σφάλματα του προηγούμενου [4]. Στην περίπτωση του Gradient Boosting, η διόρθωση αυτή γίνεται με την ελαχιστοποίηση μιας συνάρτησης κόστους μέσω της μεθόδου του γραμμικού κατήφορου (gradient descent).

Ο Gradient Boosting Classifier χρησιμοποιεί δέντρα απόφασης ως βασικά μοντέλα, που συχνά είναι περιορισμένου βάθους (weak learners), ώστε να επικεντρώνονται στη διόρθωση συγκεκριμένων σφαλμάτων. Η βασική ιδέα είναι ότι κάθε νέο δέντρο προβλέπει τα υπολειπόμενα λάθη του συνδυαστικού μοντέλου. Αυτό επιτυγχάνεται μέσω της ελαχιστοποίησης της συνάρτησης απώλειας που σχετίζεται με το πρόβλημα. Για προβλήματα ταξινόμησης, η συνάρτηση απώλειας μπορεί να είναι η log-loss, ενώ για παλινδρόμηση χρησιμοποιείται συχνά το μέσο τετραγωνικό σφάλμα (mean squared error). Το γεγονός ότι το Gradient Boosting προσπαθεί να βελτιστοποιήσει τη συνάρτηση απώλειας σε κάθε βήμα το καθιστά εξαιρετικά αποτελεσματικό σε προβλήματα με πολύπλοκες σχέσεις.

Ένα από τα σημαντικότερα πλεονεκτήματα του είναι η υψηλή του απόδοση σε προβλήματα με ανισόρροπα δεδομένα και σύνθετες σχέσεις μεταξύ χαρακτηριστικών και στόχου. Είναι ιδιαίτερα ανθεκτικό σε outliers, καθώς η προσαρμογή του μοντέλου γίνεται σε βήματα που επικεντρώνονται στα

πιο δύσκολα δείγματα. Επίσης, το μοντέλο Gradient Boosting παρέχει ενσωματωμένη δυνατότητα αξιολόγησης της σημασίας των χαρακτηριστικών (feature importance), επιτρέποντας στον χρήστη να κατανοήσει καλύτερα ποια χαρακτηριστικά επηρεάζουν περισσότερο τις προβλέψεις.

Παρά την υψηλή του απόδοση, παρουσιάζει και ορισμένα μειονεκτήματα. Αρχικά, είναι υπολογιστικά απαιτητικό, καθώς κάθε νέο δέντρο εξαρτάται από τα σφάλματα του προηγούμενου, γεγονός που αυξάνει τον χρόνο εκπαίδευσης, ειδικά σε μεγάλα datasets. Επίσης, μπορεί να παρουσιάσει υπερπροσαρμογή (overfitting) εάν το μοντέλο δεν ρυθμιστεί σωστά. Παράμετροι όπως ο αριθμός των δέντρων (`n_estimators`), το βάθος των δέντρων (`max_depth`), η ταχύτητα εκμάθησης (`learning_rate`) και το ελάχιστο πλήθος δειγμάτων ανά κόμβο (`min_samples_split`) πρέπει να επιλέγονται προσεκτικά μέσω τεχνικών όπως η Grid Search ή η Random Search. Η σωστή ρύθμιση αυτών των παραμέτρων είναι κρίσιμη για την αποφυγή υπερπροσαρμογής και για τη βελτίωση της γενίκευσης του μοντέλου.

Ο Gradient Boosting Classifier έχει βρει εφαρμογή σε πολλούς τομείς, όπως τα χρηματοοικονομικά (πρόβλεψη πιστωτικού κινδύνου), την υγεία (διάγνωση ασθενειών), το μάρκετινγκ (πρόβλεψη συμπεριφοράς πελατών) και άλλες περιοχές που απαιτούν υψηλή ακρίβεια στις προβλέψεις. Στη σύγχρονη εποχή, βελτιώσεις και παραλλαγές του Gradient Boosting, όπως τα μοντέλα XGBoost και LightGBM, έχουν ενισχύσει την ταχύτητα και την απόδοσή του, καθιστώντας το μια από τις πρώτες επιλογές για προβλήματα όπου απαιτείται μέγιστη ακρίβεια [9].

Συνοψίζοντας, ο Gradient Boosting Classifier είναι ένα ισχυρό εργαλείο μηχανικής μάθησης που προσφέρει εξαιρετική ακρίβεια, ευελιξία και ανθεκτικότητα. Παρόλο που απαιτεί προσεκτική ρύθμιση των υπερπαραμέτρων και έχει αυξημένο υπολογιστικό κόστος, η ικανότητά του να διαχειρίζεται σύνθετα datasets και να αποδίδει υψηλά επίπεδα ακρίβειας το καθιστά αναντικατάστατο σε πολλές εφαρμογές.

3.4. Μετρικές Αξιολόγησης

Η αξιολόγηση της απόδοσης ενός μοντέλου μηχανικής μάθησης αποτελεί καθοριστικό βήμα στην ανάπτυξη και επιλογή της βέλτιστης μεθοδολογίας. Στο πλαίσιο της ταξινόμησης, ιδιαίτερα σε προβλήματα δυαδικής φύσης, είναι αναγκαίο να υιοθετηθούν κατάλληλες μετρικές που αποτυπώνουν με ακρίβεια τη συμπεριφορά του μοντέλου έναντι των διαθέσιμων δεδομένων. Η απλή χρήση της ακρίβειας (accuracy) δεν είναι πάντοτε επαρκής, ειδικά σε περιπτώσεις όπου παρατηρείται έντονη ανισορροπία μεταξύ των κατηγοριών της μεταβλητής στόχου.

Στον υπολογισμό των μετρικών αξιολόγησης χρησιμοποιούνται οι παρακάτω βασικοί όροι, που προκύπτουν από τον πίνακα Confusion Matrix, γνωστό και ως πίνακας λάθους (error matrix):

- TP (True Positives): Πλήθος δειγμάτων που ανήκουν στην θετική κλάση και προβλέφθηκαν σωστά ως θετικά.
- TN (True Negatives): Πλήθος δειγμάτων που ανήκουν στην αρνητική κλάση και προβλέφθηκαν σωστά ως αρνητικά.
- FP (False Positives): Πλήθος δειγμάτων που ανήκουν στην αρνητική κλάση αλλά προβλέφθηκαν εσφαλμένα ως θετικά.
- FN (False Negatives): Πλήθος δειγμάτων που ανήκουν στην θετική κλάση αλλά προβλέφθηκαν εσφαλμένα ως αρνητικά.

Η ακρίβεια (accuracy) ορίζεται ως το ποσοστό των σωστά ταξινομημένων δειγμάτων επί του συνόλου των παρατηρήσεων. Αν και πρόκειται για μια ευρέως διαδεδομένη μετρική, μπορεί να οδηγήσει σε παραπλανητικά συμπεράσματα όταν η μία κλάση υπερισχύει σημαντικά της άλλης. Σε τέτοιες περιπτώσεις, ένα μοντέλο μπορεί να εμφανίζει υψηλή ακρίβεια προβλέποντας συνεχώς την πλειοψηφούσα κλάση, χωρίς ουσιαστική προγνωστική ικανότητα. Η ακρίβεια, ορίζεται από τον παρακάτω τύπο:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

Για την αποτύπωση της ευαισθησίας του μοντέλου απέναντι στην κατηγορία ενδιαφέροντος, χρησιμοποιούνται η ανάκληση (recall) και η ακρίβεια θετικών προβλέψεων (precision). Η ανάκληση ορίζεται ως το κλάσμα των σωστών προβλέψεων ανα κλάση προς το άθροισμα της γραμμής που σχετίζεται με την συγκεκριμένη κλάση. Δηλαδή οι σωστές προβλέψεις δια τον πραγματικό αριθμό των δειγμάτων για την συγκεκριμένη κλάση.

$$Recall = \frac{TP}{TP + FN}$$

Αντίστοιχα, η ακρίβεια ορίζεται ως το κλάσμα των σωστών προβλέψεων για την συγκεκριμένη κλάση προς το άθροισμα της στήλης που σχετίζεται με την συγκεκριμένη κλάση. Δηλαδή οι σωστές προβλέψεις δια τον αριθμό που ο αλγόριθμος προέβλεψε για την συγκεκριμένη κλάση.

$$Precision = \frac{TP}{TP + FP}$$

Ο συνδυασμός των δύο αυτών μεγεθών συνοψίζεται μέσω του F1-score, το οποίο αποτελεί τον αρμονικό μέσο όρο της ακρίβειας και της ανάκλησης, προσφέροντας μια ισορροπημένη εκτίμηση της απόδοσης, ιδιαίτερα όταν υπάρχει ασυμμετρία στις κλάσεις ή όταν τόσο τα ψευδώς θετικά όσο και τα ψευδώς αρνητικά είναι κρίσιμα.

$$F1 - score = 2 \cdot \frac{Precision \cdot Recall}{Precision + Recall}$$

Εξίσου σημαντική είναι η καμπύλη ROC (Receiver Operating Characteristic), η οποία απεικονίζει τη σχέση μεταξύ του ποσοστού αληθώς θετικών (True Positive Rate) και του ποσοστού ψευδώς θετικών (False Positive Rate) για διαφορετικά κατώφλια ταξινόμησης. Το εμβαδόν κάτω από την καμπύλη ROC (AUC) λειτουργεί ως συνοπτική μέτρηση της ικανότητας του μοντέλου να διακρίνει μεταξύ των κλάσεων.

$$ROC - AUC = \int_0^1 TPR(x) dx$$

Τιμές AUC κοντά στο 1 υποδηλώνουν υψηλή διακριτική ικανότητα, ενώ τιμές κοντά στο 0.5 υποδηλώνουν απόδοση αντίστοιχη με τυχαία πρόβλεψη.

Η κατάλληλη επιλογή μετρικής εξαρτάται σε μεγάλο βαθμό από τις ιδιαιτερότητες του προβλήματος. Σε εφαρμογές με έντονη ανισορροπία, όπως η πρόβλεψη αθέτησης πληρωμών ή η ανταπόκριση σε προωθητικές ενέργειες, η εστίαση αποκλειστικά στην ακρίβεια ενδέχεται να οδηγήσει σε λανθασμένα συμπεράσματα. Αντίθετα, μετρικές όπως το F1-score και το ROC-AUC προσφέρουν πιο αντιπροσωπευτική εικόνα της επίδοσης του μοντέλου και αξιοποιούνται ευρέως για τη συγκριτική αξιολόγηση διαφορετικών αλγορίθμων.

4. Ηθικά Ζητήματα και Κανονιστικές Προκλήσεις στη Χρήση Τεχνητής Νοημοσύνης στον Τραπεζικό Τομέα

Η χρήση τεχνικών εξόρυξης γνώσης και αλγορίθμων τεχνητής νοημοσύνης στον χρηματοπιστωτικό τομέα δημιουργεί, εκτός από σημαντικές ευκαιρίες, και μία σειρά από ηθικά και κανονιστικά ζητήματα. Τα οικονομικά δεδομένα, ως δεδομένα που αφορούν σε προσωπική και οικονομική κατάσταση των ατόμων,

χαρακτηρίζονται από υψηλή ευαισθησία. Η αυτοματοποίηση αποφάσεων μέσω αλγορίθμων μπορεί να επηρεάσει άμεσα την πρόσβαση των πολιτών σε βασικές υπηρεσίες, όπως η παροχή δανείων, η ασφάλιση ή η χρηματοδότηση.

4.1. Ζητήματα Μεροληψίας και Διακρίσεων (Algorithmic Bias)

Η ανάπτυξη και ευρεία χρήση τεχνητής νοημοσύνης στον τραπεζικό τομέα έχει αναδείξει το ζήτημα της μεροληψίας των αλγορίθμων ως έναν από τους σοβαρότερους κινδύνους που συνδέονται με την εφαρμογή τέτοιων συστημάτων στη λήψη χρηματοπιστωτικών αποφάσεων. Η μεροληψία, είτε προέρχεται από τα ίδια τα δεδομένα είτε από τον τρόπο με τον οποίο σχεδιάζονται και εκπαιδεύονται τα μοντέλα, μπορεί να οδηγήσει σε άνισες ή άδικες αποφάσεις που ενισχύουν υφιστάμενες κοινωνικές και οικονομικές ανισότητες. Στην περίπτωση των τραπεζών, όπου οι αποφάσεις αφορούν άμεσα τη χορήγηση δανείων, τη διαμόρφωση πιστωτικών ορίων και την πρόσβαση σε οικονομικά προϊόντα, οι επιπτώσεις αυτής της μεροληψίας είναι εξαιρετικά σημαντικές, αφού μπορούν να στερήσουν σε ευρείες ομάδες πληθυσμού το δικαίωμα στην ισότιμη οικονομική συμμετοχή [11].

Η βασικότερη πηγή της αλγοριθμικής μεροληψίας είναι η ποιότητα και η σύνθεση των δεδομένων που χρησιμοποιούνται για την εκπαίδευση των μοντέλων. Καθώς οι αλγόριθμοι μηχανικής μάθησης βασίζονται στην ανάλυση ιστορικών δεδομένων για να εντοπίσουν πρότυπα συμπεριφοράς, οποιαδήποτε προϋπάρχουσα κοινωνική ή οικονομική ανισότητα αποτυπώνεται στα δεδομένα αυτά και αναπαράγεται στο μοντέλο. Για παράδειγμα, αν στο παρελθόν η πρόσβαση σε στεγαστικά ή επιχειρηματικά δάνεια ήταν περιορισμένη για συγκεκριμένες κοινωνικές ομάδες, όπως οι γυναίκες, οι μετανάστες ή τα άτομα νεαρής ηλικίας, τότε τα ιστορικά δεδομένα δανειοδότησης θα περιέχουν ήδη αυτές τις αποκλίσεις. Εάν τα δεδομένα αυτά εισαχθούν αυτούσια στην εκπαίδευση του μοντέλου, ο αλγόριθμος θα «μάθει» ότι οι αιτήσεις δανείων από άτομα των ομάδων αυτών συνοδεύονται από αυξημένο πιστωτικό κίνδυνο, όχι απαραίτητα λόγω της πραγματικής φερεγγυότητάς τους, αλλά επειδή παραδοσιακά το τραπεζικό σύστημα είχε δείξει μεγαλύτερη επιφυλακτικότητα απέναντί τους.

Η μεροληψία ενδέχεται επίσης να προκύψει μέσω μεταβλητών που δρουν έμμεσα ως υποκατάστατα ευαίσθητων κοινωνικών χαρακτηριστικών. Για παράδειγμα, η ταχυδρομική διεύθυνση, η περιοχή κατοικίας ή το είδος επαγγέλματος μπορεί να συνδέονται στενά με κοινωνικοοικονομική κατάσταση, εθνοτική καταγωγή ή επίπεδο εκπαίδευσης. Αν και τα δεδομένα αυτά δεν θεωρούνται ευαίσθητα από νομικής άποψης, η συσχέτισή τους με προστατευόμενα κοινωνικά χαρακτηριστικά σημαίνει ότι η χρήση τους μπορεί να επιτείνει φαινόμενα έμμεσης διάκρισης. Ακόμη και όταν οι ευαίσθητες μεταβλητές

αφαιρούνται από τα datasets, οι proxy variables αρκούν για να αναπαράγουν την προκατάληψη. Ο αλγόριθμος, μέσα από στατιστικές συσχετίσεις, μπορεί να καταλήγει σε παρόμοια συμπεράσματα αποκλεισμού, παρά την απουσία ρητά διακριτών κριτηρίων.

Επιπλέον, οι τεχνικές παράμετροι του τρόπου εκπαίδευσης των μοντέλων μπορεί να ενισχύσουν το πρόβλημα. Τα σύνολα εκπαίδευσης πολύ συχνά εμφανίζουν έντονη ανισορροπία μεταξύ των κλάσεων. Στις περισσότερες τραπεζικές εφαρμογές, οι δανειολήπτες που αθετούν τις υποχρεώσεις τους αποτελούν ποσοτικά μικρή ομάδα. Όταν η μειονοτική κλάση των επισφαλών πελατών δεν αντιπροσωπεύεται επαρκώς, τα μοντέλα δυσκολεύονται να εκπαιδευτούν σωστά ως προς τη συμπεριφορά αυτής της κατηγορίας, επιδεινώνοντας τις προβλέψεις για άτομα που ενδέχεται να ανήκουν σε ήδη ευάλωτες κοινωνικές ομάδες. Επιπλέον, τα φαινόμενα υπερπροσαρμογής των μοντέλων στα δεδομένα εκπαίδευσης μπορεί να οδηγήσουν στην ενίσχυση των υφιστάμενων στατιστικών συσχετίσεων, μεταφέροντας τις υπάρχουσες κοινωνικές στρεβλώσεις αυτούσιες στο σύστημα αποφάσεων.

Η ύπαρξη αλγοριθμικής μεροληψίας δεν είναι απλώς ένα θεωρητικό πρόβλημα. Έχει άμεσες πρακτικές συνέπειες για τους δανειολήπτες και την κοινωνία ευρύτερα. Άτομα με ουσιαστική οικονομική φερεγγυότητα μπορεί να απορρίπτονται αδικαιολόγητα από πιστωτικά προϊόντα, επειδή ανήκουν σε ομάδες που παραδοσιακά είχαν χαμηλότερη πρόσβαση στη χρηματοδότηση. Η άνιση μεταχείριση δεν περιορίζεται μόνο στις εγκρίσεις δανείων αλλά μπορεί να αφορά και τη διαμόρφωση διαφορετικών πιστωτικών ορίων, επιτοκίων ή όρων αποπληρωμής, με δυσανάλογο κόστος για τους λιγότερο προνομιούχους. Έτσι, τα ίδια τα συστήματα τεχνητής νοημοσύνης, αντί να εξομαλύνουν ανισότητες, κινδυνεύουν να διευρύνουν το χάσμα πρόσβασης σε οικονομικές ευκαιρίες.

Πέρα από τις κοινωνικές συνέπειες, η μεροληψία δημιουργεί και σημαντικούς κινδύνους για τις ίδιες τις τράπεζες. Η χρήση μεροληπτικών αλγορίθμων εκθέτει τα ιδρύματα σε κανονιστικές παραβάσεις και σε παραβίαση του πλαισίου ισότητας στην οικονομική μεταχείριση πολιτών. Διεθνώς, ήδη ισχύουν και ενισχύονται ρυθμιστικά πλαίσια που επιβάλλουν τη διασφάλιση της μη διάκρισης στις αυτοματοποιημένες χρηματοοικονομικές αποφάσεις, όπως η νομοθεσία Equal Credit Opportunity Act στις Ηνωμένες Πολιτείες ή οι κατευθυντήριες οδηγίες της Ευρωπαϊκής Ένωσης στο πλαίσιο του GDPR και των συζητήσεων για τον Ευρωπαϊκό AI Act. Στο ευρωπαϊκό πλαίσιο, ο σεβασμός της ισότητας στην αυτοματοποιημένη λήψη αποφάσεων θεωρείται αναπόσπαστο κομμάτι της αρχής της δίκαιης επεξεργασίας προσωπικών δεδομένων.

Η αντιμετώπιση της μεροληψίας στις αλγοριθμικές διαδικασίες αξιολόγησης δανειοληπτών απαιτεί πολυδιάστατη προσέγγιση. Αρχικά, είναι αναγκαία η συνειδητή αξιολόγηση των συνόλων δεδομένων πριν την εκπαίδευση, ώστε να εντοπίζονται ενδείξεις προκατάληψης σε συγκεκριμένα χαρακτηριστικά ή πληθυσμιακές ομάδες. Τα δεδομένα πρέπει να αναλύονται όχι μόνο ως προς τη συνολική τους ποιότητα αλλά και σε επίπεδο αντιπροσωπευτικότητας διαφορετικών κοινωνικών και δημογραφικών υποομάδων. Στη συνέχεια, κατά την επιλογή των χαρακτηριστικών που τροφοδοτούν το μοντέλο, πρέπει να αποφεύγεται η χρήση μεταβλητών που σχετίζονται έμμεσα με ευαίσθητα κριτήρια, ακόμη και αν φαινομενικά αποτελούν χρήσιμες πηγές πρόβλεψης.

Πέρα όμως από την επιλογή μεταβλητών, κρίσιμη είναι και η υιοθέτηση στατιστικών δεικτών μέτρησης της αλγοριθμικής μεροληψίας, ώστε το fairness να ελέγχεται ως ποσοτικοποιήσιμη διάσταση της αξιολόγησης των μοντέλων. Έτσι, μέσω δεικτών όπως η ισότητα ευκαιριών (equal opportunity), η ισότητα πιθανότητας (demographic parity) ή η εξισορροπημένη ακρίβεια μεταξύ των ομάδων (equalized odds), μπορεί να διαπιστώνεται αν το μοντέλο παράγει συστηματικές διαφορές στην αξιολόγηση διαφορετικών κοινωνικών ομάδων. Παράλληλα, αναπτύσσονται συνεχώς καινοτόμες τεχνικές διόρθωσης bias, οι οποίες εφαρμόζονται είτε κατά το στάδιο της εκπαίδευσης (pre-processing), είτε κατά τη διαδικασία μοντελοποίησης (in-processing), είτε ακόμη και μετά την παραγωγή των αποτελεσμάτων (post-processing), στοχεύοντας στη μείωση των διακρίσεων χωρίς να θυσιάζεται η προγνωστική ικανότητα του μοντέλου.

Τέλος, είναι σημαντικό να τονιστεί ότι η διαχείριση του algorithmic bias δεν αποτελεί αποκλειστικά τεχνικό ζήτημα. Απαιτεί την εμπλοκή της ανώτερης διοίκησης, της νομικής υπηρεσίας, των μονάδων κανονιστικής συμμόρφωσης και εσωτερικού ελέγχου, καθώς και τη θεσμοθέτηση πολιτικών ηθικής διακυβέρνησης της τεχνητής νοημοσύνης στο σύνολο του οργανισμού. Η θεσμοθέτηση λογοδοσίας για την ηθική χρήση των αλγορίθμων αποτελεί πλέον βασικό στοιχείο υπεύθυνης τραπεζικής πρακτικής, εντασσόμενο σε ευρύτερα προγράμματα ESG (Environmental, Social, Governance) και κοινωνικής ευθύνης των οργανισμών.

Με τον τρόπο αυτό, η πρόοδος της τεχνητής νοημοσύνης στον χρηματοοικονομικό κλάδο μπορεί να συνδυαστεί με τον σεβασμό στις αρχές της ισότητας και της κοινωνικής δικαιοσύνης, αποτρέποντας τον κίνδυνο μετατροπής των αλγορίθμων από εργαλεία καινοτομίας σε μηχανισμούς διαίωξης κοινωνικών αποκλεισμών.

4.2. Διαφάνεια και Ερμηνευσιμότητα Μοντέλων (Explainability)

Η διαφάνεια και η ερμηνευσιμότητα των μοντέλων τεχνητής νοημοσύνης αποτελούν έναν από τους πλέον σημαντικούς άξονες προβληματισμού για την ασφαλή και δίκαιη εφαρμογή τους στον τραπεζικό τομέα. Καθώς οι αλγόριθμοι αποκτούν αυξανόμενο ρόλο στη λήψη αποφάσεων που αφορούν την πιστοληπτική ικανότητα των δανειοληπτών, τη διαμόρφωση επιτοκίων και την αξιολόγηση επενδυτικών κινδύνων, η ανάγκη κατανόησης του τρόπου λειτουργίας τους από όλους τους εμπλεκόμενους γίνεται κρίσιμη [10]. Σε αντίθεση με τα παραδοσιακά συστήματα κανόνων, όπου οι παράγοντες λήψης απόφασης είναι ρητά διατυπωμένοι, τα συστήματα μηχανικής μάθησης, και ιδίως τα σύγχρονα σύνθετα μοντέλα όπως τα βαθιά νευρωνικά δίκτυα και οι μέθοδοι ensemble, λειτουργούν με τρόπο συχνά αδιαφανή, καθιστώντας δυσχερή την κατανόηση των λόγων που οδήγησαν σε μια συγκεκριμένη απόφαση.

Το πρόβλημα της ερμηνευσιμότητας δεν είναι αποκλειστικά τεχνικό. Αφορά άμεσα τα δικαιώματα των πολιτών, την υποχρέωση λογοδοσίας των χρηματοπιστωτικών οργανισμών και τη δυνατότητα παρέμβασης των εποπτικών αρχών. Για τον δανειολήπτη, η απόρριψη μιας αίτησης δανείου χωρίς επαρκή αιτιολόγηση μπορεί να συνιστά παραβίαση του δικαιώματός του για ενημέρωση και διαφάνεια στη λήψη αποφάσεων που επηρεάζουν την οικονομική του ζωή. Για τον τραπεζικό οργανισμό, η αδυναμία εξήγησης των αποφάσεων δυσκολεύει τον εσωτερικό έλεγχο, την κανονιστική συμμόρφωση και τη διαχείριση νομικών και φήμης κινδύνων.

Η δυσκολία στην ερμηνεία προκύπτει από τη φύση των ίδιων των μοντέλων. Οι σύγχρονοι αλγόριθμοι, προκειμένου να επιτύχουν υψηλή προγνωστική ακρίβεια, αξιοποιούν πολύπλοκους μηχανισμούς βελτιστοποίησης που δεν επιτρέπουν εύκολη αναγωγή της τελικής απόφασης σε συγκεκριμένα, εύκολα κατανοητά κριτήρια. Τα βαθιά νευρωνικά δίκτυα, για παράδειγμα, αποτελούνται από πολυεπίπεδες δομές τεχνητών νευρώνων, όπου η τελική απόφαση προκύπτει από εκατομμύρια παραμέτρους και ενδογενείς σχέσεις μεταξύ των εισόδων και των εξόδων του μοντέλου. Παρομοίως, τα μοντέλα ensemble, που βασίζονται στον συνδυασμό πολλών απλών μοντέλων σε μια ενιαία απόφαση, δημιουργούν ένα πλέγμα επιμέρους κανόνων τόσο σύνθετο που καθιστά δύσκολη την παρακολούθηση της διαδρομής που ακολουθεί η κάθε εισροή για να φτάσει στην τελική εκτίμηση.

Η έλλειψη διαφάνειας μπορεί να λειτουργήσει αποθαρρυντικά ακόμη και σε επίπεδο διοίκησης τραπεζικών οργανισμών. Στελέχη που λαμβάνουν στρατηγικές αποφάσεις για την εφαρμογή τέτοιων συστημάτων συχνά εκφράζουν επιφυλάξεις όταν δεν είναι σε θέση να κατανοήσουν επαρκώς τον τρόπο

λειτουργίας τους. Η δυσκολία κατανόησης του τρόπου με τον οποίο οι αλγόριθμοι αξιολογούν κινδύνους μπορεί να περιορίσει τη διάχυση της καινοτομίας, ακριβώς επειδή ενισχύει την αίσθηση απώλειας ελέγχου επί των διαδικασιών.

Παράλληλα, οι ρυθμιστικές αρχές τόσο στην Ευρωπαϊκή Ένωση όσο και διεθνώς αναγνωρίζουν την ανάγκη ενίσχυσης της διαφάνειας ως προϋπόθεση υπεύθυνης χρήσης τεχνητής νοημοσύνης. Στο πλαίσιο του Ευρωπαϊκού Κανονισμού για την Προστασία Δεδομένων (GDPR), το δικαίωμα του υποκειμένου των δεδομένων να λαμβάνει ουσιαστική πληροφόρηση για τη λογική που ακολουθείται στις αυτοματοποιημένες αποφάσεις αποτελεί ήδη θεμελιώδη διάταξη. Παράλληλα, στο υπό διαμόρφωση Ευρωπαϊκό AI Act προβλέπεται ρητά η ανάγκη για explainability και auditability σε αλγορίθμους υψηλού κινδύνου, στους οποίους περιλαμβάνονται και τα συστήματα αξιολόγησης πιστωτικής ικανότητας.

Για την αντιμετώπιση του προβλήματος έχουν αναπτυχθεί, ιδίως την τελευταία πενταετία, πληθώρα μεθόδων ερμηνείας των αποφάσεων των αλγορίθμων. Οι τεχνικές αυτές, γνωστές διεθνώς ως explainable AI (XAI), στοχεύουν στην ανάδειξη των παραγόντων που επηρέασαν την απόφαση του μοντέλου, καθιστώντας πιο διαφανή τη λειτουργία του χωρίς να απαιτείται απλοποίηση του ίδιου του μοντέλου. Μεταξύ των πλέον διαδεδομένων εργαλείων συγκαταλέγονται οι μέθοδοι SHAP και LIME, οι οποίες βασίζονται σε διαφορετικές στατιστικές προσεγγίσεις για την ποσοτικοποίηση της συνεισφοράς κάθε μεταβλητής στην τελική πρόβλεψη.

Η μέθοδος SHAP στηρίζεται σε έννοιες από τη θεωρία παιγνίων και επιτρέπει τη μέτρηση της ακριβούς συμβολής κάθε εισροής στη διαμόρφωση του αποτελέσματος, παρέχοντας τόσο ατομικού όσο και συνολικού επιπέδου επεξηγήσεις. Από την άλλη πλευρά, η προσέγγιση LIME δημιουργεί τοπικά μοντέλα απλοποίησης γύρω από κάθε μεμονωμένη πρόβλεψη, προσφέροντας έτσι εύκολη κατανόηση της συγκεκριμένης απόφασης σε επίπεδο μεμονωμένου πελάτη. Οι μέθοδοι αυτές δεν αντικαθιστούν την αρχιτεκτονική του μοντέλου αλλά λειτουργούν παράλληλα, προσφέροντας στους αναλυτές και τους εποπτικούς φορείς τη δυνατότητα να παρακολουθούν πώς το μοντέλο κατέληξε σε κάθε απόφαση.

Παρά την πρόοδο στην ανάπτυξη εργαλείων explainability, εξακολουθούν να υφίστανται προκλήσεις στην καθολική εφαρμογή τους. Ορισμένα πολύπλοκα συστήματα εξακολουθούν να παράγουν αποφάσεις που είναι δύσκολο να ερμηνευτούν σε πλήρως κατανοητή μορφή για μη εξειδικευμένους χρήστες. Η ισορροπία μεταξύ της ακρίβειας του μοντέλου και της ερμηνευσιμότητάς του παραμένει κρίσιμο ζητούμενο, ιδίως σε περιβάλλοντα υψηλής κανονιστικής απαίτησης όπως οι τράπεζες. Συχνά τίθεται το ερώτημα αν είναι προτιμότερο να υιοθετούνται απλούστερα μοντέλα, ελαφρώς λιγότερο ακριβή αλλά

απολύτως ερμηνεύσιμα, έναντι σύνθετων μοντέλων που επιτυγχάνουν βέλτιστη ακρίβεια αλλά είναι πρακτικά αδιαφανή.

Η επιλογή δεν είναι πάντοτε εύκολη, ιδίως όταν οι επιπτώσεις από μια λανθασμένη εκτίμηση πιστωτικού κινδύνου μπορεί να είναι σημαντικές τόσο για την τράπεζα όσο και για τον πελάτη. Σε κάθε περίπτωση, ο τραπεζικός οργανισμός οφείλει να σταθμίσει τη σημασία της ακρίβειας έναντι της διαφάνειας υπό το πρίσμα της κανονιστικής συμμόρφωσης, της προστασίας των δικαιωμάτων των πελατών και της διασφάλισης της εταιρικής του φήμης.

Η ερμηνευσιμότητα των αλγορίθμων δεν αποτελεί απλώς τεχνικό ζήτημα εσωτερικής σημασίας, αλλά συνδέεται άμεσα με τη διατήρηση της κοινωνικής αποδοχής και της εμπιστοσύνης του κοινού στη χρήση τεχνητής νοημοσύνης στο χρηματοπιστωτικό σύστημα. Σε έναν τομέα όπου οι αποφάσεις επηρεάζουν τη δυνατότητα των πολιτών να συμμετέχουν ενεργά στην οικονομική ζωή, η διαφάνεια λειτουργεί ως θεμέλιο υπεύθυνης τεχνολογικής εξέλιξης και μακροπρόθεσμης βιωσιμότητας της καινοτομίας.

4.3. Προστασία Προσωπικών Δεδομένων και Κανονιστική Συμμόρφωση

Η ανάπτυξη συστημάτων τεχνητής νοημοσύνης και εξόρυξης γνώσης στον τραπεζικό τομέα, παράλληλα με τα πλεονεκτήματα που προσφέρει, εγείρει σοβαρά ζητήματα προστασίας προσωπικών δεδομένων και συμμόρφωσης με τα εκάστοτε κανονιστικά πλαίσια. Οι τραπεζικές δραστηριότητες βασίζονται σε μεγάλο βαθμό στην ανάλυση πληθώρας προσωπικών και οικονομικών στοιχείων των πελατών, τα οποία συχνά εμπεριέχουν ευαίσθητες πληροφορίες που αφορούν το εισόδημα, την επαγγελματική κατάσταση, την καταναλωτική συμπεριφορά, το ιστορικό συναλλαγών, ακόμη και τον τρόπο αποπληρωμής των οικονομικών υποχρεώσεων. Η επεξεργασία αυτών των δεδομένων από συστήματα τεχνητής νοημοσύνης καθιστά απαραίτητη την τήρηση ενός αυστηρού πλαισίου προστασίας της ιδιωτικότητας, καθώς οποιαδήποτε παραβίαση θα μπορούσε να πλήξει σοβαρά τόσο τα δικαιώματα των πολιτών όσο και την αξιοπιστία του χρηματοπιστωτικού συστήματος.

Σε αντίθεση με άλλους τομείς, όπου τα δεδομένα είναι συχνά ανώνυμα ή σχετικά απρόσωπα, οι τραπεζικές συναλλαγές συνθέτουν ένα εξαιρετικά αναλυτικό προφίλ για κάθε πελάτη. Μέσω της συνδυαστικής ανάλυσης διαφορετικών πηγών πληροφορίας, οι αλγόριθμοι αποκτούν τη δυνατότητα να δημιουργούν προβλέψεις και εκτιμήσεις που φτάνουν σε υψηλό βαθμό ακρίβειας ως προς την οικονομική συμπεριφορά των ατόμων. Η παρακολούθηση ροών συναλλαγών, η συσχέτιση δαπανών με τοποθεσίες ή χρονικά πρότυπα και η αξιολόγηση της σταθερότητας εισοδήματος συνθέτουν ένα ψηφιακό αποτύπωμα που μπορεί να χρησιμοποιηθεί σε πλήθος επιχειρησιακών αποφάσεων. Ωστόσο, όσο

μεγαλύτερη είναι η επεξεργασία αυτών των δεδομένων, τόσο ενισχύεται και η ανάγκη αυστηρής προστασίας τους.

Η Ευρωπαϊκή Ένωση έχει ήδη διαμορφώσει ένα από τα αυστηρότερα νομοθετικά πλαίσια προστασίας προσωπικών δεδομένων παγκοσμίως, με τον Γενικό Κανονισμό για την Προστασία Δεδομένων (GDPR) να θέτει σαφή όρια και υποχρεώσεις σε κάθε οργανισμό που επεξεργάζεται δεδομένα πολιτών εντός της επικράτειάς της. Στο πλαίσιο του τραπεζικού τομέα, οι αρχές του GDPR αποκτούν ιδιαίτερη βαρύτητα, καθώς η βάση των τραπεζικών υπηρεσιών στηρίζεται σχεδόν αποκλειστικά στη συλλογή και επεξεργασία τέτοιων δεδομένων. Ο κανονισμός επιβάλλει τη διασφάλιση ότι η συλλογή των δεδομένων γίνεται για συγκεκριμένους και σαφείς σκοπούς, με ρητή συγκατάθεση των πελατών, ενώ παράλληλα παρέχει στους πολίτες δικαιώματα πρόσβασης, διόρθωσης και διαγραφής των προσωπικών τους στοιχείων. [12]

Ιδιαίτερη έμφαση δίνεται στην αρχή της ελαχιστοποίησης των δεδομένων, σύμφωνα με την οποία κάθε οργανισμός οφείλει να συλλέγει και να επεξεργάζεται μόνο το απολύτως απαραίτητο σύνολο πληροφοριών για την επίτευξη του εκάστοτε επιχειρησιακού σκοπού. Στην περίπτωση των αλγορίθμων, αυτό συνεπάγεται ότι πρέπει να αποφεύγεται η ανεξέλεγκτη συγκέντρωση μεγάλου όγκου δεδομένων χωρίς επαρκή αιτιολόγηση, ακόμη και όταν αυτά θεωρούνται χρήσιμα για τη βελτίωση της ακρίβειας των μοντέλων. Η αυξημένη διαθεσιμότητα δεδομένων λόγω των τεχνολογιών big data δεν αίρει τις υποχρεώσεις του οργανισμού ως προς την αναγκαιότητα της επεξεργασίας.

Παράλληλα, το άρθρο 22 του GDPR εισάγει ειδική ρύθμιση για την αυτοματοποιημένη λήψη αποφάσεων, η οποία καλύπτει ευθέως τις τραπεζικές πρακτικές που βασίζονται σε αλγορίθμους αξιολόγησης πιστοληπτικής ικανότητας [13]. Σύμφωνα με αυτή τη διάταξη, κάθε πρόσωπο διατηρεί το δικαίωμα να μην υπόκειται σε απόφαση που λαμβάνεται αποκλειστικά μέσω αυτοματοποιημένης επεξεργασίας, εφόσον αυτή έχει έννομες συνέπειες ή τον επηρεάζει σημαντικά. Στο πλαίσιο αυτό, οι τράπεζες οφείλουν να εξασφαλίζουν μηχανισμούς ανθρώπινης παρέμβασης και αξιολόγησης όταν πρόκειται για κρίσιμες αποφάσεις χορήγησης ή απόρριψης δανείων, ακριβώς για να προστατευτεί ο δανειολήπτης από ενδεχόμενα σφάλματα ή άδικες κρίσεις του αλγορίθμου.

Πέρα από τον GDPR, σε διεθνές επίπεδο διαμορφώνεται σταδιακά ένα ευρύτερο κανονιστικό πλαίσιο αναφορικά με την προστασία δεδομένων και τη χρήση τεχνητής νοημοσύνης. Η Ευρωπαϊκή Ένωση προχωρά ήδη στην υιοθέτηση του Ευρωπαϊκού Κανονισμού για την Τεχνητή Νοημοσύνη (AI Act), ο οποίος αναμένεται να θέσει πρόσθετους κανόνες για τις εφαρμογές υψηλού ρίσκου, στις οποίες εντάσσονται τα χρηματοοικονομικά συστήματα αξιολόγησης. Ο νέος κανονισμός δίνει ιδιαίτερη έμφαση

στην αρχή της διαφάνειας, της λογοδοσίας, της ασφάλειας και του σεβασμού των θεμελιωδών δικαιωμάτων των πολιτών, ενισχύοντας την προστασία τους από τις ενδεχόμενες αδυναμίες ή αδικίες των αυτόματων συστημάτων λήψης αποφάσεων.

Για τους τραπεζικούς οργανισμούς, η συμμόρφωση με τα παραπάνω κανονιστικά πλαίσια δεν αποτελεί απλώς νομική υποχρέωση αλλά και στρατηγικό ζήτημα διαχείρισης κινδύνου και εταιρικής φήμης. Η οποιαδήποτε παραβίαση προσωπικών δεδομένων μπορεί να οδηγήσει σε σημαντικές διοικητικές κυρώσεις, βαριά οικονομικά πρόστιμα και απώλεια εμπιστοσύνης εκ μέρους των πελατών και της κοινωνίας. Επιπλέον, σε περιβάλλον έντονου ανταγωνισμού, η διασφάλιση της ιδιωτικότητας και η υπεύθυνη διαχείριση των προσωπικών δεδομένων αποτελούν πλέον παράγοντες διαφοροποίησης και προσέλκυσης πελατείας που εκτιμά την αξιοπιστία και την ηθική λειτουργία του οργανισμού.

Η προστασία προσωπικών δεδομένων στο πλαίσιο ανάπτυξης αλγορίθμων τεχνητής νοημοσύνης δεν εξαντλείται μόνο στο στάδιο της συλλογής και επεξεργασίας, αλλά επεκτείνεται σε ολόκληρο τον κύκλο ζωής του συστήματος. Περιλαμβάνει την ασφαλή αποθήκευση, την προστασία κατά τη μεταφορά δεδομένων, την εσωτερική διαβάθμιση πρόσβασης στο προσωπικό, τη διατήρηση αρχείων επεξεργασίας (logs) και την ύπαρξη μηχανισμών audit που διασφαλίζουν την ιχνηλασιμότητα κάθε ενέργειας στο σύστημα. Παράλληλα, απαιτείται η συνεχής παρακολούθηση της συμμόρφωσης, καθώς οι τεχνολογίες εξελίσσονται και οι κίνδυνοι μεταβάλλονται.

Η εδραίωση κουλτούρας σεβασμού των προσωπικών δεδομένων εντός του οργανισμού αποτελεί αναπόσπαστο μέρος της σύγχρονης τραπεζικής διακυβέρνησης. Η επένδυση σε εκπαιδευτικά προγράμματα, η ενίσχυση των μονάδων κανονιστικής συμμόρφωσης και η διαρκής συνεργασία με τις αρμόδιες εποπτικές αρχές αποτελούν κρίσιμες ενέργειες για την αποτελεσματική θωράκιση των τραπεζικών συστημάτων απέναντι στους αυξανόμενους κινδύνους παραβίασης ιδιωτικότητας.

Σε αυτό το πλαίσιο, η αρμονική συνύπαρξη της τεχνητής νοημοσύνης με την προστασία της ιδιωτικής ζωής δεν αποτελεί αντικρουόμενη συνθήκη αλλά πρόκληση ορθής ισορροπίας. Η υπεύθυνη διαχείριση των δεδομένων, με σεβασμό στα θεμελιώδη δικαιώματα των πολιτών, μπορεί να αποτελέσει το θεμέλιο για βιώσιμη και κοινωνικά αποδεκτή καινοτομία στον τραπεζικό τομέα.

4.4. Δίκαιη Πρόσβαση σε Χρηματοοικονομικές Υπηρεσίες

Η χρήση τεχνητής νοημοσύνης στον τραπεζικό τομέα δεν επηρεάζει μόνο την ακρίβεια ή την αποδοτικότητα των χρηματοοικονομικών αποφάσεων, αλλά αγγίζει βαθύτερα ζητήματα ισότητας και

κοινωνικής δικαιοσύνης. Η δίκαιη πρόσβαση σε τραπεζικές υπηρεσίες αποτελεί διαχρονικά κεντρικό ζήτημα δημόσιας πολιτικής, καθώς η οικονομική ενσωμάτωση συνδέεται άμεσα με την ευημερία των ατόμων και τη σταθερότητα του κοινωνικού ιστού. Οι χρηματοπιστωτικές υπηρεσίες, όπως η πίστωση, η αποταμίευση, οι ασφαλίσσεις και οι επενδύσεις, αποτελούν βασικά εργαλεία άσκησης οικονομικών δικαιωμάτων. Όταν η πρόσβαση σε αυτά τα αγαθά περιορίζεται, οι επιπτώσεις δεν αφορούν μόνο το μεμονωμένο άτομο αλλά επηρεάζουν τη συνολική κοινωνική ανάπτυξη και διευρύνουν τις κοινωνικές ανισότητες.

Η ενσωμάτωση αλγορίθμων στην αξιολόγηση πιστοληπτικής ικανότητας, στον καθορισμό όρων δανειοδότησης ή στην εκτίμηση ασφαλιστικών κινδύνων, εφόσον δεν εφαρμοστεί με προσοχή, μπορεί να οδηγήσει σε νέες μορφές αποκλεισμού, ακόμη και όταν αυτές δεν προέρχονται από συνειδητές επιλογές του χρηματοπιστωτικού οργανισμού. Σε πολλές περιπτώσεις, τα μοντέλα τεχνητής νοημοσύνης, λειτουργώντας βάσει ιστορικών δεδομένων και στατιστικών συσχετίσεων, ενδέχεται να κρίνουν άτομα ή ομάδες ως υψηλού κινδύνου αποκλειστικά λόγω των χαρακτηριστικών του περιβάλλοντος στο οποίο ανήκουν, παραβλέποντας την ατομική τους πιστοληπτική συμπεριφορά ή δυναμική.

Χαρακτηριστικό παράδειγμα αυτού του φαινομένου συνιστά η εκτίμηση κινδύνου βάσει γεωγραφικής περιοχής. Περιοχές με χαμηλότερο μέσο εισόδημα ή ιστορικό οικονομικών δυσχερειών συχνά καταγράφονται σε βάσεις δεδομένων ως ζώνες υψηλότερης πιστωτικής επισφάλειας. Αν ο αλγόριθμος βασιστεί ανεπεξέργαστα σε τέτοια δεδομένα, τότε οι κάτοικοι των περιοχών αυτών ενδέχεται να αξιολογούνται συλλογικά ως επισφαλείς πελάτες, ανεξαρτήτως των προσωπικών οικονομικών τους στοιχείων. Αντίστοιχα, επαγγελματικές ομάδες με μεγαλύτερη εργασιακή ευθραυστότητα, όπως οι ελεύθεροι επαγγελματίες ή οι αυτοαπασχολούμενοι, μπορεί να εμφανίζονται ως στατιστικά υψηλότερου κινδύνου, με αποτέλεσμα τον περιορισμό της πρόσβασής τους σε προϊόντα πίστωσης ή την επιβολή αυστηρότερων όρων.

Το ζήτημα επιτείνεται σε κοινωνίες με ήδη υψηλά επίπεδα οικονομικής ανισότητας, όπου οι ιστορικοί αποκλεισμοί έχουν αφήσει έντονα αποτυπώματα στις βάσεις δεδομένων. Ο κίνδυνος εδώ έγκειται στη διαίωνηση ενός φαύλου κύκλου: οι ομάδες που αντιμετώπιζαν δυσκολίες πρόσβασης στο παρελθόν συνεχίζουν να απορρίπτονται και στο παρόν, όχι λόγω της τρέχουσας ατομικής τους κατάστασης αλλά λόγω της στατιστικής εικόνας του πληθυσμιακού τους δείγματος. Έτσι, η τεχνητή νοημοσύνη, αντί να λειτουργεί ως εργαλείο άρσης των ανισοτήτων, κινδυνεύει να λειτουργήσει ως μηχανισμός συντήρησης και εδραίωσης κοινωνικών αποκλεισμών.

Η δίκαιη πρόσβαση σε τραπεζικές υπηρεσίες δεν αφορά μόνο την ηθική διάσταση των οικονομικών συναλλαγών αλλά συνδέεται άμεσα και με τη χρηματοπιστωτική σταθερότητα. Ένα τραπεζικό σύστημα που αποκλείει συστηματικά ευρείες κοινωνικές ομάδες περιορίζει το εύρος της πελατειακής του βάσης, μειώνει την κυκλοφορία κεφαλαίων στην οικονομία και περιορίζει τις ευκαιρίες χρηματοδότησης νέων επιχειρηματικών δραστηριοτήτων. Παράλληλα, η άνιση πρόσβαση ενισχύει το αίσθημα κοινωνικής αδικίας, γεγονός που σε βάθος χρόνου μπορεί να υπονομεύσει τη θεσμική αξιοπιστία των χρηματοπιστωτικών οργανισμών και να ενισχύσει την κοινωνική πόλωση.

Αναγνωρίζοντας αυτούς τους κινδύνους, οι εποπτικές αρχές διεθνώς εντείνουν τις προσπάθειές τους για την προώθηση του fairness στις αλγοριθμικές διαδικασίες αξιολόγησης. Στο πλαίσιο αυτό, επιβάλλεται η υιοθέτηση πρακτικών συνεχούς παρακολούθησης και αξιολόγησης των μοντέλων ως προς την επίδρασή τους σε διαφορετικές κοινωνικές ομάδες. Οι οργανισμοί καλούνται να αναλύουν συστηματικά τα αποτελέσματα των συστημάτων τους, εντοπίζοντας ενδείξεις συστηματικής άνισης μεταχείρισης σε υποομάδες πληθυσμού και να λαμβάνουν διορθωτικά μέτρα όταν αυτό απαιτείται.

Η δίκαιη πρόσβαση συνδέεται και με την ευρύτερη στρατηγική βιωσιμότητας και υπευθυνότητας των τραπεζικών οργανισμών. Στο πλαίσιο των διεθνών προτύπων ESG, το κοινωνικό σκέλος της εταιρικής διακυβέρνησης αποκτά κεντρικό ρόλο. Οι τράπεζες που επιδεικνύουν υψηλό επίπεδο κοινωνικής ευαισθησίας και διασφαλίζουν την ισότιμη πρόσβαση σε χρηματοπιστωτικές υπηρεσίες ενισχύουν τη φήμη τους, διαφοροποιούνται θετικά στην αγορά και κερδίζουν την εμπιστοσύνη τόσο των πελατών όσο και των επενδυτών.

Στην ουσία, το fairness στις αλγοριθμικές τραπεζικές πρακτικές δεν μπορεί να αντιμετωπίζεται ως τεχνική βελτιστοποίηση του μοντέλου. Αποτελεί ευρύτερο ζήτημα κουλτούρας, διοικητικής ευθύνης και εταιρικής στρατηγικής. Η δίκαιη πρόσβαση σε χρηματοπιστωτικές υπηρεσίες οφείλει να αποτελεί σταθερή προτεραιότητα, ανεξαρτήτως της τεχνολογικής πολυπλοκότητας που μπορεί να συνοδεύει την υλοποίηση των συστημάτων. Μόνο μέσα από την ενσωμάτωση της ισότητας και της κοινωνικής ευθύνης στον πυρήνα της τεχνολογικής και επιχειρησιακής λειτουργίας μπορεί η τραπεζική τεχνητή νοημοσύνη να υπηρετήσει πραγματικά την κοινωνική πρόοδο και τη βιώσιμη ανάπτυξη.

4.5. Προτάσεις Ορθής Διαχείρισης Ηθικών Κινδύνων

Η αναγνώριση των ηθικών προκλήσεων που συνοδεύουν την ανάπτυξη και εφαρμογή συστημάτων τεχνητής νοημοσύνης στον τραπεζικό τομέα δεν αρκεί από μόνη της για την ουσιαστική διαχείρισή τους. Απαιτείται η διαμόρφωση ενός ολοκληρωμένου πλαισίου στρατηγικών, διαδικασιών και οργανωτικών

πρακτικών, ικανών να προλαμβάνουν, να εντοπίζουν και να αντιμετωπίζουν εγκαίρως τους πιθανούς ηθικούς κινδύνους που ενέχουν τα αλγοριθμικά συστήματα. Η ορθή διαχείριση αυτών των κινδύνων δεν αφορά αποκλειστικά τα τεχνικά τμήματα ανάπτυξης, αλλά οφείλει να διαπερνά οριζόντια όλη την οργανωτική δομή των χρηματοπιστωτικών ιδρυμάτων.

Καταρχάς, το θεμέλιο της ορθής ηθικής διακυβέρνησης είναι η εγκαθίδρυση ξεκάθαρων εσωτερικών πολιτικών για την υπεύθυνη χρήση της τεχνητής νοημοσύνης. Οι πολιτικές αυτές πρέπει να προσδιορίζουν σαφώς τις αρχές και τις αξίες που διέπουν την ανάπτυξη, τη λειτουργία και την παρακολούθηση των συστημάτων AI εντός του οργανισμού. Σε αυτό το πλαίσιο, κεντρική θέση κατέχουν η αρχή της δικαιοσύνης, ο σεβασμός της ιδιωτικότητας, η διαφάνεια, η λογοδοσία και η συμμόρφωση με το κανονιστικό πλαίσιο. Οι οργανισμοί οφείλουν να διατυπώνουν επίσημα κώδικες δεοντολογίας που να καλύπτουν το σύνολο του κύκλου ζωής των συστημάτων, από τη συλλογή των δεδομένων μέχρι την παραγωγική εφαρμογή και την αξιολόγηση των επιδόσεων.

Παράλληλα, είναι απαραίτητο να διασφαλιστεί η ενεργή συμμετοχή των μονάδων κανονιστικής συμμόρφωσης και νομικών υπηρεσιών ήδη από τα πρώτα στάδια ανάπτυξης νέων αλγορίθμων. Η στενή συνεργασία μεταξύ των τεχνικών ομάδων και των ρυθμιστικών οργάνων εντός του οργανισμού διασφαλίζει ότι οι νομικές υποχρεώσεις και οι κανονιστικές απαιτήσεις λαμβάνονται υπόψη κατά τον σχεδιασμό των συστημάτων, μειώνοντας τον κίνδυνο μεταγενέστερων παραβάσεων ή διορθωτικών παρεμβάσεων.

Ιδιαίτερη σημασία έχει η καθιέρωση ανεξάρτητων διαδικασιών εσωτερικού ελέγχου και εποπτείας για την αξιολόγηση της ηθικής συμμόρφωσης των αλγορίθμων. Οι εσωτερικοί έλεγχοι δεν πρέπει να περιορίζονται μόνο σε ζητήματα τεχνικής απόδοσης, αλλά να εξετάζουν συστηματικά το ενδεχόμενο ύπαρξης αλγοριθμικής μεροληψίας, άνισης μεταχείρισης ή ανεπαρκούς διαφάνειας στις αποφάσεις. Η ύπαρξη μονίμων επιτροπών ηθικής επεξεργασίας δεδομένων, που λειτουργούν με διατομεακή σύνθεση επιστημόνων, νομικών, ειδικών κανονιστικής συμμόρφωσης και διοικητικών στελεχών, μπορεί να προσφέρει πρόσθετα επίπεδα ελέγχου και διασφάλισης.

Επιπλέον, οι οργανισμοί οφείλουν να υιοθετήσουν και να ενσωματώσουν συστηματικές πρακτικές ποσοτικής μέτρησης ηθικών δεικτών στα πλαίσια αξιολόγησης των μοντέλων τους. Οι δείκτες αυτοί περιλαμβάνουν μετρικές fairness, ισότητας ευκαιριών, διαφοροποίησης επιπτώσεων σε ευάλωτες ομάδες πληθυσμού, καθώς και δείκτες explainability που αξιολογούν τον βαθμό κατανόησης των αποφάσεων από τους χρήστες. Μέσα από τη συνεχή παρακολούθηση αυτών των δεικτών σε πραγματικό

περιβάλλον λειτουργίας, οι οργανισμοί μπορούν να εντοπίζουν εγκαίρως δυσλειτουργίες ή απόκλιση από τα ηθικά πρότυπα που έχουν θέσει.

Αναπόσπαστο στοιχείο της ηθικής διαχείρισης αποτελεί και η διαρκής εκπαίδευση του προσωπικού σε ζητήματα τεχνητής νοημοσύνης και δεοντολογίας δεδομένων. Η ηθική διαχείριση δεν μπορεί να επαφίεται αποκλειστικά στους ειδικούς της τεχνολογίας ή στους νομικούς συμβούλους. Όλα τα διοικητικά στελέχη, οι αναλυτές, οι υπεύθυνοι προϊόντων και οι λειτουργοί εξυπηρέτησης πελατών πρέπει να διαθέτουν βασική κατανόηση των δυνατοτήτων, των κινδύνων και των ηθικών διακυβευμάτων που συνοδεύουν την εφαρμογή των νέων συστημάτων. Μέσα από συστηματικά εκπαιδευτικά προγράμματα και workshops, η ευαισθητοποίηση σε θέματα ethical AI μπορεί να καταστεί αναπόσπαστο στοιχείο της καθημερινής λειτουργίας του οργανισμού.

Ιδιαίτερη πρόκληση αποτελεί επίσης η ανάγκη διαφάνειας προς το ίδιο το κοινό και τους πελάτες του οργανισμού. Οι χρηματοπιστωτικοί οργανισμοί οφείλουν να παρέχουν ουσιαστική και κατανοητή πληροφόρηση προς τους πολίτες για το πώς χρησιμοποιούνται τα προσωπικά τους δεδομένα, ποιοι αλγόριθμοι επεξεργάζονται τις αιτήσεις τους και με ποια κριτήρια λαμβάνονται οι αποφάσεις που τους αφορούν. Η ενίσχυση της εμπιστοσύνης του κοινού στις τεχνολογίες τεχνητής νοημοσύνης συνδέεται άρρηκτα με την διαφάνεια και τη διασφάλιση λογοδοσίας από την πλευρά των τραπεζών.

Τέλος, η αποτελεσματική διαχείριση ηθικών κινδύνων απαιτεί συνεχή συνεργασία και διάλογο με τις εποπτικές αρχές και τους νομοθέτες, ιδίως σε ένα ρυθμιστικό περιβάλλον που εξελίσσεται δυναμικά λόγω της ταχύτητας των τεχνολογικών εξελίξεων. Οι τράπεζες καλούνται να συμμετέχουν ενεργά στη διαμόρφωση των νέων κανονιστικών πλαισίων, παρέχοντας εμπειρική γνώση και συμβάλλοντας στην ισορροπημένη διατύπωση κανόνων που προστατεύουν αφενός τα δικαιώματα των πολιτών και αφετέρου την καινοτόμο ανάπτυξη της χρηματοπιστωτικής τεχνολογίας.

5. Βαθιά Μηχανική Μάθηση (Deep Learning) σε Οικονομικά και Τραπεζικά Δεδομένα

Η βαθιά μηχανική μάθηση (Deep Learning) αποτελεί μια προηγμένη υποκατηγορία της μηχανικής μάθησης, με κύριο χαρακτηριστικό τη χρήση πολυεπίπεδων νευρωνικών δικτύων για την αναπαράσταση και εκμάθηση περίπλοκων σχέσεων στα δεδομένα [7]. Οι τεχνικές της βαθιάς μάθησης έχουν βρει εφαρμογή σε τομείς όπως η όραση υπολογιστών, η φωνητική αναγνώριση, η φυσική γλώσσα, αλλά και

όλο και περισσότερο στα οικονομικά και τραπεζικά δεδομένα, λόγω της αυξημένης διαθεσιμότητας ιστορικών χρονοσειρών και δεδομένων συμπεριφοράς πελατών.

Σε αντίθεση με τις παραδοσιακές μεθόδους μηχανικής μάθησης, όπου απαιτείται η προσεκτική επιλογή και μη αυτόματη εξαγωγή χαρακτηριστικών (feature engineering), τα νευρωνικά δίκτυα βαθιάς μάθησης έχουν την ικανότητα να μαθαίνουν αυτόματα πολυεπίπεδες αναπαραστάσεις των δεδομένων, μειώνοντας την ανάγκη για εξωτερική παρέμβαση.

5.1. Νευρωνικά Δίκτυα και Εφαρμογές σε Τραπεζικά και Οικονομικά Δεδομένα

Η πρόοδος της βαθιάς μηχανικής μάθησης (Deep Learning) συνδέεται στενά με την αυξανόμενη διαθεσιμότητα μεγάλων συνόλων δεδομένων και με την εξέλιξη της υπολογιστικής ισχύος. Στο πλαίσιο των τραπεζικών και οικονομικών εφαρμογών, η δυνατότητα αξιοποίησης ιστορικών δεδομένων πελατών, συναλλαγών, πιστωτικών κινδύνων και προτιμήσεων ανοίγει νέους ορίζοντες για την εφαρμογή βαθιών νευρωνικών αρχιτεκτονικών. Παρότι παραδοσιακές τεχνικές μηχανικής μάθησης προσφέρουν ικανοποιητικά αποτελέσματα, έχουν περιορισμούς, ιδίως όταν τα δεδομένα είναι υψηλής πολυπλοκότητας, μη γραμμικά ή δεν έχουν υποστεί εκτενή προεπεξεργασία.

Τα «ρηχά» μοντέλα, όπως τα decision trees ή η logistic regression, βασίζονται σε προκαθορισμένα χαρακτηριστικά, τα οποία απαιτούν συχνά χειροκίνητο σχεδιασμό (feature engineering). Στην περίπτωση μεγάλων τραπεζικών datasets, αυτή η διαδικασία γίνεται χρονοβόρα και ενίοτε ανεπαρκής, καθώς δεν αποτυπώνει τις πολυδιάστατες αλληλεπιδράσεις των δεδομένων. Η βαθιά μάθηση έρχεται να καλύψει αυτό το κενό, καθώς μπορεί να εντοπίζει αφανή μοτίβα χωρίς να απαιτείται ανθρώπινη παρέμβαση για τον ορισμό των χαρακτηριστικών.

Αρχιτεκτονικές όπως τα Convolutional Neural Networks (CNNs) βρίσκουν εφαρμογή σε χρηματοοικονομικά δεδομένα που σχετίζονται με χρονοσειρές ή μοτίβα συναλλαγών, όπως η ανάλυση συνεχόμενων καταναλωτικών συμπεριφορών, εντοπισμός ασυνήθιστων αναλήψεων ή συχνών δαπανών σε χρονικές ακολουθίες. Παράλληλα, τα Transformers, που αποτέλεσαν επανάσταση στην Επεξεργασία Φυσικής Γλώσσας (NLP), χρησιμοποιούνται πλέον στην ανάλυση σχολίων πελατών, αιτήσεων δανείων, αλλά και στην αυτόματη επεξεργασία κειμένων από συμβάσεις ή πολιτικές πίστωσης.

Σημαντική είναι επίσης η συμβολή της βαθιάς μάθησης στη βελτιστοποίηση της απόδοσης στρατηγικών πίστωσης ή διαχείρισης πελατών, μέσω μεθόδων ενισχυτικής μάθησης (Reinforcement Learning). Αντί να

εφαρμόζονται στατικά μοντέλα αξιολόγησης ρίσκου, χρησιμοποιούνται μοντέλα που «μαθαίνουν» συνεχώς από την αλληλεπίδραση με τον πελάτη – για παράδειγμα, στο πλαίσιο προσαρμοστικών συστημάτων τιμολόγησης ή στη διαχείριση πιστωτικών ορίων, όπου το μοντέλο μαθαίνει πότε και πώς να εγκρίνει αυξήσεις ή να μειώνει τα διαθέσιμα όρια.

Ένα από τα πλέον καινοτόμα πεδία είναι αυτό των Graph Neural Networks (GNNs). Οι τραπεζικοί οργανισμοί μπορούν να μοντελοποιήσουν το σύνολο των συναλλαγών και των σχέσεων πελατών ως έναν γράφο, όπου οι κόμβοι είναι οι πελάτες και οι ακμές οι συναλλαγές ή οι κοινοί προμηθευτές. Έτσι, η ανάλυση δεν βασίζεται μόνο στα μεμονωμένα χαρακτηριστικά του πελάτη, αλλά και στη θέση του στο συνολικό «οικοσύστημα». Αυτό επιτρέπει την ανίχνευση κυκλωμάτων απάτης (fraud rings), τη βελτιστοποίηση δικτύων συνεργασίας (π.χ. με εμπόρους ή μεσολαβητές), ή ακόμη και την αναγνώριση συμπεριφορών υψηλού κινδύνου που επηρεάζονται από τη «γειτονιά» ενός χρήστη στο δίκτυο.

Ιδιαίτερη σημασία αποκτούν τα inductive GNNs, τα οποία έχουν τη δυνατότητα να παράγουν έγκυρες αναπαραστάσεις ακόμη και για νέους πελάτες ή συναλλαγές που δεν έχουν παρατηρηθεί στο παρελθόν. Αυτό είναι καθοριστικό για την αναγνώριση ρίσκου σε νέα προϊόντα ή στην αξιολόγηση πιστοληπτικής ικανότητας χωρίς εκτενές ιστορικό (π.χ. νέοι χρήστες χωρίς credit history).

Γενικά, τα νευρωνικά δίκτυα προσφέρουν τη δυνατότητα μοντελοποίησης πολυπαραγοντικών σχέσεων και αξιοποίησης μεγάλων, ανομοιογενών συνόλων δεδομένων, χωρίς την ανάγκη για χειροκίνητη δημιουργία χαρακτηριστικών. Ειδικά στον τραπεζικό τομέα, όπου οι αλληλεπιδράσεις είναι σύνθετες, οι παράγοντες ρίσκου πολυδιάστατοι, και η συμπεριφορά των πελατών εξελίσσεται διαρκώς, τα βαθιά δίκτυα μπορούν να προσφέρουν δυναμικά και προσαρμοζόμενα μοντέλα πρόβλεψης και απόφασης.

Ωστόσο, παρά τα πλεονεκτήματα, η υιοθέτηση της βαθιάς μάθησης στις χρηματοοικονομικές εφαρμογές συνοδεύεται και από προκλήσεις, όπως η ερμηνευσιμότητα των μοντέλων (black-box issue), οι αυξημένες υπολογιστικές απαιτήσεις και η ανάγκη συμμόρφωσης με αυστηρά κανονιστικά πλαίσια (compliance, auditing, fairness). Παρ' όλα αυτά, η ενσωμάτωση αυτών των τεχνικών αποτελεί μια από τις πιο υποσχόμενες κατευθύνσεις για τον μελλοντικό σχεδιασμό έξυπνων, αυτοματοποιημένων και ασφαλών τραπεζικών υπηρεσιών.

5.2. Δομικά Στοιχεία Νευρωνικών Δικτύων σε Τραπεζικά Συστήματα

Η κατανόηση των δομικών στοιχείων ενός νευρωνικού δικτύου είναι απαραίτητη για την αξιολόγηση της χρησιμότητάς του στην επεξεργασία τραπεζικών και οικονομικών δεδομένων. Τα παρακάτω θεμελιώδη

μέρη περιγράφουν τον τρόπο λειτουργίας αυτών των μοντέλων, καθώς και πώς αξιοποιούνται σε πραγματικά προβλήματα, όπως η πρόβλεψη αθέτησης πληρωμών ή η ανίχνευση απάτης.

- Τεχνητός Νευρώνας

Ο βασικός «δομικός λίθος» κάθε νευρωνικού δικτύου είναι ο τεχνητός νευρώνας. Ουσιαστικά, δέχεται μία ομάδα εισόδων, τις οποίες σταθμίζει με συγκεκριμένα βάρη, υπολογίζει το αθροιστικό αποτέλεσμα και εφαρμόζει μια μη γραμμική συνάρτηση ενεργοποίησης στην έξοδο. Στο πλαίσιο οικονομικών προβλημάτων, οι εισοδοί μπορεί να είναι μεταβλητές όπως εισόδημα, ιστορικό πληρωμών, υπόλοιπα καρτών, ή αριθμός συναλλαγών.

- Συνάρτηση Ενεργοποίησης (Activation Function)

Οι συναρτήσεις ενεργοποίησης είναι αυτές που προσδίδουν μη γραμμικότητα στο δίκτυο, επιτρέποντας του να μάθει πολύπλοκες σχέσεις μεταξύ μεταβλητών. Στην οικονομική ανάλυση, αυτές οι μη γραμμικότητες είναι σημαντικές, καθώς οι σχέσεις μεταξύ εισοδημάτων, κινδύνου και πιστωτικής συμπεριφοράς συχνά δεν είναι γραμμικές. Συνήθεις επιλογές περιλαμβάνουν τις ReLU, tanh και softmax (ιδίως στην έξοδο για προβλήματα classification, όπως «Θα αθετήσει ο πελάτης την πληρωμή;»).

- Επίπεδο Εισόδου (Input Layer)

Πρόκειται για το πρώτο στρώμα του δικτύου, το οποίο δέχεται τις αριθμητικές ή κατηγορικές μεταβλητές ως εισόδους. Για παράδειγμα, σε ένα τραπεζικό dataset, μπορεί να περιλαμβάνει μεταβλητές όπως ηλικία, εισόδημα, αριθμός λογαριασμών ή ύπαρξη δανείων.

- Κρυμμένα Επίπεδα (Hidden Layers)

Τα κρυμμένα επίπεδα είναι αυτά στα οποία γίνεται η εκμάθηση χαρακτηριστικών. Σε αντίθεση με τις παραδοσιακές μεθόδους, εδώ το δίκτυο μπορεί να ανακαλύψει αφανή πρότυπα ή «συνδυαστικά χαρακτηριστικά» που δεν είναι άμεσα εμφανή στους αναλυτές. Για παράδειγμα, μπορεί να συνδυάσει τη συμπεριφορά αποπληρωμής με την εποχικότητα δαπανών για να εντοπίσει πρόβλημα ρευστότητας.

- Επίπεδο Εξόδου (Output Layer)

Το τελικό επίπεδο παράγει την προβλεπόμενη έξοδο. Στα περισσότερα οικονομικά προβλήματα πρόβλεψης, πρόκειται για ένα πρόβλημα δυαδικής ταξινόμησης (π.χ. default ή όχι), και χρησιμοποιείται η softmax ή sigmoid για τη μετατροπή της εξόδου σε πιθανότητες.

- Εκπαίδευση και Backpropagation

Η διαδικασία εκπαίδευσης βασίζεται στον αλγόριθμο της οπισθοδιάδοσης σφάλματος (backpropagation), ο οποίος υπολογίζει το σφάλμα ανάμεσα στην πρόβλεψη και την πραγματική τιμή και προσαρμόζει τα βάρη αναλόγως. Στον τραπεζικό τομέα, αυτό ισοδυναμεί με τη «μάθηση» του δικτύου για το πώς οι μεταβλητές σχετίζονται με τον κίνδυνο αθέτησης ή την αποδοχή μιας προσφοράς marketing.

- Ρυθμός Μάθησης (Learning Rate)

Ο ρυθμός μάθησης καθορίζει πόσο γρήγορα προσαρμόζονται τα βάρη σε κάθε βήμα της εκπαίδευσης. Στις χρηματοοικονομικές εφαρμογές, όπου τα δεδομένα μπορεί να είναι ευαίσθητα και πολύπλοκα, η σωστή ρύθμιση αυτής της παραμέτρου είναι κρίσιμη για να αποφευχθεί είτε υπερεκπαίδευση είτε πολύ αργή σύγκλιση.

- Αλγόριθμοι Βελτιστοποίησης (Optimizers)

Ο πιο διαδεδομένος αλγόριθμος είναι ο Stochastic Gradient Descent (SGD). Σε οικονομικά δεδομένα μεγάλου όγκου, χρησιμοποιούνται συχνά και εκδοχές όπως ο Mini-Batch SGD, που επιτρέπει την εκπαίδευση με υποσύνολα δεδομένων. Άλλοι δημοφιλείς βελτιστοποιητές είναι οι Adam και RMSprop, που τροποποιούν δυναμικά τον ρυθμό μάθησης για καλύτερη απόδοση.

- Εποχές (Epochs)

Κάθε «εποχή» αντιστοιχεί σε μία πλήρη διέλευση του μοντέλου από το σύνολο των παρατηρήσεων. Σε κάθε εποχή, το δίκτυο βελτιώνει τα βάρη του βάσει της πληροφορίας που αντλεί από το σφάλμα.

- Τακτοποίηση (Regularization)

Για να αποφεύγεται το φαινόμενο της υπερπροσαρμογής (overfitting) — κάτι συχνό σε τραπεζικά προβλήματα με ανισόρροπες κατηγορίες (π.χ. λίγοι πελάτες σε default) — χρησιμοποιούνται τεχνικές όπως το L2 regularization, το Dropout και το early stopping. Ειδικά το Dropout επιβάλλει τυχαία απενεργοποίηση κόμβων κατά την εκπαίδευση, ενισχύοντας τη γενίκευση.

5.3. Βασικοί Αλγόριθμοι Βαθιάς Μάθησης σε Οικονομικά και Τραπεζικά Δεδομένα

Παρόλο που η παρούσα εργασία δεν εφαρμόζει πρακτικά μοντέλα βαθιάς μάθησης, αξίζει να παρουσιαστούν οι βασικές κατηγορίες αλγορίθμων που χρησιμοποιούνται ευρέως σήμερα, καθώς έχουν

ήδη βρει σημαντική εφαρμογή στον χρηματοοικονομικό και τραπεζικό τομέα. Η αυξανόμενη διαθεσιμότητα δεδομένων σε αυτούς τους τομείς, σε συνδυασμό με την ανάγκη πρόβλεψης, κατηγοριοποίησης και ανίχνευσης ρίσκου, καθιστούν τη βαθιά μάθηση ένα πολλά υποσχόμενο εργαλείο. Παρακάτω παρουσιάζονται τέσσερις βασικές οικογένειες αλγορίθμων:

- Συνελκτικά Νευρωνικά Δίκτυα (Convolutional Neural Networks - CNNs)

Τα CNNs έχουν συνδεθεί κυρίως με προβλήματα οπτικής αναγνώρισης, ωστόσο η εφαρμογή τους έχει επεκταθεί και σε δομημένα δεδομένα όπως οι χρονοσειρές τραπεζικών συναλλαγών. Στο χρηματοοικονομικό πλαίσιο, μπορούν να χρησιμοποιηθούν για την αναγνώριση μοτίβων συμπεριφοράς, για παράδειγμα σε σειρές μηνιαίων πληρωμών ή σε διαδοχικές ενέργειες πελατών. Η ικανότητά τους να εντοπίζουν τοπικά μοτίβα και συσχετίσεις σε διαδοχικά δεδομένα τα καθιστά κατάλληλα για την πρόβλεψη αθέτησης πληρωμών ή την ανίχνευση ανωμαλιών σε δαπάνες.

- Γενετικοί Ανταγωνιστικοί Αλγόριθμοι (Generative Adversarial Networks - GANs)

Οι GANs είναι γνωστοί για τη χρήση τους σε δημιουργία εικόνας ή περιεχομένου, ωστόσο βρίσκουν και εφαρμογή στη δημιουργία συνθετικών οικονομικών δεδομένων. Αυτό μπορεί να φανεί ιδιαίτερα χρήσιμο σε περιπτώσεις όπου υπάρχει ανισορροπία κλάσεων (όπως σε default datasets), καθώς μπορούν να χρησιμοποιηθούν για data augmentation, δημιουργώντας τεχνητά παραδείγματα σπάνιων κατηγοριών (π.χ. προβληματικοί δανειολήπτες). Για παράδειγμα, στη μελέτη των Fiore et al. (2019), οι GANs εφαρμόστηκαν για τη δημιουργία συνθετικών δεδομένων πιστωτικής συμπεριφοράς, με στόχο τη βελτίωση της απόδοσης ταξινομητών σε credit scoring προβλήματα. Το συνθετικό dataset βελτίωσε τη γενίκευση μοντέλων, όπως Random Forest και MLP, ιδιαίτερα σε σενάρια με χαμηλή εκπροσώπηση μη εξυπηρετούμενων δανείων. Επιπλέον, σε εφαρμογές ανίχνευσης ξεπλύματος χρήματος (AML), οι Jurgovsky et al. (2021) ανέπτυξαν GAN-based συστήματα για την ενίσχυση των training sets σε δεδομένα συναλλαγών, εντοπίζοντας πολύπλοκα μοτίβα που δύσκολα ανιχνεύονται με κλασικά rule-based μοντέλα. Επίσης, έχουν αξιοποιηθεί σε μελέτες για ανώνυμη παραγωγή τραπεζικών δεδομένων, προστατεύοντας την ιδιωτικότητα.

- Αυτοκωδικοποιητές (Autoencoders)

Οι Autoencoders είναι αλγόριθμοι μη επιβλεπόμενης μάθησης, ιδανικοί για μείωση διαστάσεων και ανίχνευση ανωμαλιών. Σε τραπεζικά περιβάλλοντα μπορούν να χρησιμοποιηθούν για συμπίεση πολύπλοκων χαρακτηριστικών (όπως 25+ χρηματοοικονομικές μεταβλητές πελάτη) σε χαμηλότερη διάσταση, διατηρώντας τη βασική πληροφορία. Επιπλέον, σε εφαρμογές fraud detection, η αδυναμία

του μοντέλου να ανακατασκευάσει ένα «ύποπτο» παράδειγμα μπορεί να υποδείξει ότι πρόκειται για ανωμαλία.

- Επαναλαμβανόμενα Νευρωνικά Δίκτυα (Recurrent Neural Networks - RNNs)

Τα RNNs είναι κατάλληλα για δεδομένα που ακολουθούν χρονική αλληλουχία, κάτι που είναι πολύ συνηθισμένο σε οικονομικά περιβάλλοντα. Για παράδειγμα, η συμπεριφορά πληρωμών, οι μεταβολές υπολοίπων ή η σειρά τραπεζικών ενεργειών ενός πελάτη μέσα στον μήνα αποτελούν χρονοσειρές. Οι πιο εξελιγμένες εκδοχές των RNNs, όπως τα LSTM και GRU [8], είναι ικανές να διατηρούν «μνήμη» για μακροχρόνιες εξαρτήσεις και μπορούν να χρησιμοποιηθούν για πρόβλεψη πιθανότητας default, ανάλυση κινδύνου ή ακόμα και πρόβλεψη ανταπόκρισης σε καμπάνιες marketing.

Σε σύγκριση με παραδοσιακά στατιστικά μοντέλα πρόβλεψης χρονοσειρών, όπως τα ARIMA ή Holt-Winters, τα RNNs και οι εκτεταμένες εκδοχές τους προσφέρουν σημαντικό πλεονέκτημα στη μοντελοποίηση μακροχρόνιων εξαρτήσεων και στη διαχείριση πολύπλοκων χρονικών σχέσεων, χωρίς να απαιτείται αυστηρή σταθερότητα ή κανονικότητα των δεδομένων.

Παράλληλα, η προσέγγιση sequence-to-sequence μπορεί να υποστηρίξει προηγμένες εφαρμογές πρόβλεψης από ιστορικά οικονομικά δεδομένα, αξιοποιώντας μηχανισμούς encoder-decoder και attention, για πιο ακριβή και ευέλικτη ανάλυση.

Η επιλογή μεταξύ παραδοσιακών (shallow) και βαθιών (deep) μοντέλων εξαρτάται σε μεγάλο βαθμό από τη φύση και την πολυπλοκότητα του προβλήματος, καθώς και από τα χαρακτηριστικά των διαθέσιμων δεδομένων. Σε περιπτώσεις όπου η πληροφορία είναι περιορισμένη, καλά δομημένη και το ζητούμενο σχετικά απλό — όπως η πρόβλεψη αποδοχής μιας προσφοράς ή η εκτίμηση πιθανότητας καθυστέρησης πληρωμής — τα παραδοσιακά μοντέλα επαρκούν και προσφέρουν ερμηνευσιμότητα και ταχύτητα. Αντίθετα, σε πιο πολύπλοκα σενάρια, όπου εμπλέκονται ακολουθίες γεγονότων, υψηλής διάστασης χαρακτηριστικά ή ανάγκη για εντοπισμό αφανών μοτίβων, οι αλγόριθμοι βαθιάς μάθησης υπερτερούν σημαντικά, αξιοποιώντας στο έπακρο τα πλεονεκτήματα που προσφέρει η μάθηση από τα ίδια τα δεδομένα. Σε τέτοιες περιπτώσεις, η επένδυση σε υπολογιστική ισχύ και εξειδικευμένη παραμετροποίηση αντισταθμίζεται από την αυξημένη ακρίβεια και ευελιξία των αποτελεσμάτων.

6. Πειραματικό Μέρος / Υλοποίηση

6.1. Περιγραφή Συνόλων Δεδομένων

Στο πλαίσιο της παρούσας εργασίας αξιοποιούνται δύο διαφορετικά σύνολα δεδομένων, τα οποία είναι:

1. Bank Marketing Dataset (<https://archive.ics.uci.edu/dataset/222/bank+marketing>) [15]
2. Default of Credit Card Clients Dataset (<https://archive.ics.uci.edu/dataset/350/default+of+credit+card+clients>) [16]

Το πρώτο σύνολο δεδομένων περιέχει πληροφορίες από καμπάνιες τηλεφωνικών πωλήσεων που έγιναν από μια τράπεζα και χρησιμοποιείται συχνά για την πρόβλεψη της ανταπόκρισης των πελατών (δηλαδή, αν θα γίνουν πελάτες ή όχι). Ακολουθεί μια αναλυτική περιγραφή του συνόλου δεδομένων:

Αριθμός Δειγμάτων (instances): 45.211 δείγματα

Αριθμός Χαρακτηριστικών (features): 17 (συμπεριλαμβανομένης της μεταβλητής στόχου)

Χαρακτηριστικά Δεδομένων:

- age: Ηλικία (ακέραιος αριθμός).
- job: Επάγγελμα (κατηγορηματικό):
 - Τιμές: management, technician, entrepreneur, blue-collar, unknown, retired, admin., services, self-employed, unemployed, κ.λπ.
- marital: Οικογενειακή κατάσταση (κατηγορηματικό):
 - Τιμές: married, single, divorced.
- education: Εκπαίδευση (κατηγορηματικό):
 - Τιμές: tertiary, secondary, unknown, primary.
- default: Υπάρχει δάνειο σε αθέτηση; (ναι/όχι):
 - Τιμές: no, yes.
- balance: Υπόλοιπο λογαριασμού (αριθμητικό).
- housing: Έχει δάνειο κατοικίας; (ναι/όχι):
 - Τιμές: yes, no.
- loan: Έχει προσωπικό δάνειο; (ναι/όχι):
 - Τιμές: yes, no.
- contact: Μέθοδος επικοινωνίας (κατηγορηματικό):

- Τιμές: unknown, cellular, telephone.
- day: Ημέρα του μήνα κατά την οποία έγινε η επαφή (αριθμητικό).
- month: Μήνας της τελευταίας επαφής (κατηγορηματικό):
 - Τιμές: may, jun, jul, aug, oct, nov, dec, jan, feb, mar, κ.λπ.
- duration: Διάρκεια της τελευταίας επαφής σε δευτερόλεπτα (αριθμητικό).
- campaign: Αριθμός επαφών κατά τη διάρκεια αυτής της καμπάνιας (αριθμητικό).
- rdays: Αριθμός ημερών από την τελευταία επαφή σε προηγούμενη καμπάνια (αριθμητικό, -1 αν δεν υπάρχει επαφή).
- previous: Αριθμός επαφών σε προηγούμενη καμπάνια (αριθμητικό).
- routcome: Αποτέλεσμα της προηγούμενης καμπάνιας (κατηγορηματικό):
 - Τιμές: may, jun, jul, aug, oct, nov, dec, jan, feb, mar, κ.λπ.
- y: Μεταβλητή στόχος - Αν ο πελάτης αποδέχθηκε την πρόταση της τράπεζας ή όχι (κατηγορηματικό):
 - Τιμές: yes, no.

Το δεύτερο σύνολο δεδομένων περιέχει δεδομένα που αφορούν την πρόβλεψη αθέτησης πληρωμών από πελάτες πιστωτικών καρτών. Ακολουθούν οι λεπτομέρειες του dataset:

Αριθμός Δειγμάτων (instances): 30.000 δείγματα

Αριθμός Χαρακτηριστικών (features): 25 (συμπεριλαμβανομένης της μεταβλητής στόχου)

Χαρακτηριστικά Δεδομένων:

- ID: Αναγνωριστικός αριθμός κάθε πελάτη (αριθμητικό).
- LIMIT_BAL: Πιστωτικό όριο του πελάτη (σε δολάρια Ταϊβάν), (αριθμητικό).
- SEX: Φύλο του πελάτη (κατηγορηματικό):
 - Τιμές: 1: άνδρας, 2: γυναίκα.
- EDUCATION: Επίπεδο εκπαίδευσης (κατηγορηματικό):
 - Τιμές: 1: Μεταπτυχιακός τίτλος, 2: Πανεπιστημιακός τίτλος, 3: Δευτεροβάθμια εκπαίδευση, 4,5,6: Άλλο.
- MARRIAGE: Οικογενειακή κατάσταση (κατηγορηματικό):
 - Τιμές: 1: Παντρεμένος, 2: Ελεύθερος, 3: Άλλο.
- AGE: Ηλικία του πελάτη (σε έτη), (αριθμητικό).
- PAY_0: Κατάσταση πληρωμής τον Σεπτέμβριο (αριθμητικό):

- Τιμές: -1: Πλήρης εξόφληση, 0: Πληρωμή στην ώρα της, 1, 2, ...: Καθυστερήση 1, 2, ... μηνών.
- PAY_2: Κατάσταση πληρωμής τον Αύγουστο (αριθμητικό).
- PAY_3: Κατάσταση πληρωμής τον Ιούλιο (αριθμητικό).
- PAY_4: Κατάσταση πληρωμής τον Ιούνιο (αριθμητικό).
- PAY_5: Κατάσταση πληρωμής τον Μάιο (αριθμητικό).
- PAY_6: Κατάσταση πληρωμής τον Απρίλιο (αριθμητικό).
- BILL_AMT1: Ποσό χρέωσης τον Σεπτέμβριο (σε δολάρια Ταϊβάν), (αριθμητικό).
- BILL_AMT2: Ποσό χρέωσης τον Αύγουστο (σε δολάρια Ταϊβάν), (αριθμητικό).
- BILL_AMT3: Ποσό χρέωσης τον Ιούλιο (σε δολάρια Ταϊβάν), (αριθμητικό).
- BILL_AMT4: Ποσό χρέωσης τον Ιούνιο (σε δολάρια Ταϊβάν), (αριθμητικό).
- BILL_AMT5: Ποσό χρέωσης τον Μάιο (σε δολάρια Ταϊβάν), (αριθμητικό).
- BILL_AMT6: Ποσό χρέωσης τον Απρίλιο (σε δολάρια Ταϊβάν), (αριθμητικό).
- PAY_AMT1: Ποσό πληρωμής τον Σεπτέμβριο (σε δολάρια Ταϊβάν), (αριθμητικό).
- PAY_AMT2: Ποσό πληρωμής τον Αύγουστο (σε δολάρια Ταϊβάν), (αριθμητικό).
- PAY_AMT3: Ποσό πληρωμής τον Ιούλιο (σε δολάρια Ταϊβάν), (αριθμητικό).
- PAY_AMT4: Ποσό πληρωμής τον Ιούνιο (σε δολάρια Ταϊβάν), (αριθμητικό).
- PAY_AMT5: Ποσό πληρωμής τον Μάιο (σε δολάρια Ταϊβάν), (αριθμητικό).
- PAY_AMT6: Ποσό πληρωμής τον Απρίλιο (σε δολάρια Ταϊβάν), (αριθμητικό).
- default.payment.next.month: Μεταβλητή στόχος - Αν ο πελάτης καθυστέρησε την πληρωμή ή όχι (κατηγορηματικό).
 - Τιμές: 0: Καμία αθέτηση πληρωμής, 1: Αθέτηση πληρωμής.

Η επιλογή αυτών των δύο datasets εξυπηρετεί τον σκοπό της παρούσας εργασίας, ο οποίος είναι η διερεύνηση της εφαρμογής τεχνικών εξόρυξης γνώσης σε πραγματικά οικονομικά δεδομένα και η σύγκριση των αποτελεσμάτων υπό διαφορετικές συνθήκες. Η σύγκριση καλύπτει όχι μόνο διαφορετικούς τύπους προβλημάτων (απόκριση σε μάρκετινγκ έναντι αθέτησης πληρωμών), αλλά και διαφορετικές δομές και μεταβλητές, προσφέροντας πλούσιο υλικό για μελέτη και εξαγωγή χρήσιμων συμπερασμάτων.

6.2. Προ-επεξεργασία Δεδομένων

Στο dataset «Bank Marketing» αρχικά γίνεται έλεγχος για κενές τιμές με την `isnull().sum()`, χωρίς όμως να παρατηρούνται σημαντικά κενά δεδομένα. Στη συνέχεια, οι τιμές της στήλης Y μετατρέπονται από

κατηγορικές (yes,no) σε αριθμητικές (1,0), ώστε να μπορεί να χρησιμοποιηθεί η στήλη από τα μοντέλα μηχανικής μάθησης. Τα χαρακτηριστικά (X) διαχωρίζονται από τη στήλη Y. Έπειτα, οι κατηγορικές μεταβλητές μετατρέπονται σε αριθμητικές με One-Hot Encoding. Χρησιμοποιείται η StandardScaler για να κανονικοποιηθούν οι αριθμητικές μεταβλητές. Αυτή η διαδικασία βοηθά στη σταθεροποίηση των μεγεθών των χαρακτηριστικών. Τα δεδομένα χωρίζονται σε σύνολα εκπαίδευσης και δοκιμών σε αναλογία 80-20 χρησιμοποιώντας `train_test_split`. Για τη διαχείριση ανισορροπίας των δεδομένων εφαρμόζονται η Smote (Oversampling), η οποία δημιουργεί συνθετικά δείγματα για την κατηγορία μειονότητας, βασισμένη στους υπάρχοντες γείτονες (k-Nearest Neighbors), καθώς και η Random Undersampling, η οποία μειώνει τον αριθμό των δειγμάτων της κατηγορίας πλειονότητας, διατηρώντας την ισορροπία με την κατηγορία μειονότητας.

Στο dataset «Default of Credit Card Clients» η προ-επεξεργασία των δεδομένων ξεκινά με την αφαίρεση της πρώτης στήλης ως μη χρήσιμη. Στη συνέχεια, όλα τα δεδομένα μετατρέπονται σε αριθμητικά με τη χρήση της `pd.to_numeric`. Ακόμη, οποιεσδήποτε εγγραφές με κενά δεδομένα αφαιρούνται. Όμοια με το dataset «Bank Marketing», όλα τα αριθμητικά χαρακτηριστικά κανονικοποιούνται χρησιμοποιώντας την StandardScaler, όπως επίσης τα δεδομένα χωρίζονται σε σύνολα εκπαίδευσης και δοκιμών σε αναλογία 80-20. Για την ανισορροπία των δεδομένων, εφαρμόζονται και πάλι η Smote (Oversampling) για την αύξηση των δεδομένων της μειονότητας και η Random Undersampling, για τη μείωση των δεδομένων της πλειονότητας.

Η υλοποίηση των διαδικασιών προεπεξεργασίας, καθώς και το σύνολο της πειραματικής ανάλυσης που ακολουθεί στην παρακάτω ενότητα, πραγματοποιούνται με τη χρήση της γλώσσας Python. Συγκεκριμένα, χρησιμοποιούνται οι βιβλιοθήκες pandas και numpy για τη διαχείριση και τον μετασχηματισμό των δεδομένων, η scikit-learn για την εφαρμογή των αλγορίθμων ταξινόμησης και την αξιολόγηση της απόδοσής τους, η imbalanced-learn για τεχνικές εξισορρόπησης τάξεων, και οι matplotlib και seaborn για την οπτικοποίηση αποτελεσμάτων και συσχετίσεων.

6.3. Αποτελέσματα – Σύγκριση Αποτελεσμάτων

Σε αυτή την ενότητα, τα αποτελέσματα παρουσιάζονται μέσα από γραφήματα, πίνακες και βασικές μετρικές αξιολόγησης, όπως η ακρίβεια (Accuracy), ο συντελεστής ROC-AUC και οι καμπύλες ROC. Η ανάλυση επικεντρώνεται στη συγκριτική απόδοση των αλγορίθμων Random Forest, Logistic Regression και Gradient Boosting, αναδεικνύοντας τα πλεονεκτήματα και τα μειονεκτήματα κάθε τεχνικής ανά dataset. Μέσω αυτής της προσέγγισης, επιχειρείται η κατανόηση του τρόπου με τον οποίο οι μέθοδοι

μηχανικής μάθησης μπορούν να αξιοποιηθούν σε διαφορετικά οικονομικά δεδομένα, αλλά και να υποστηρίξουν την εξαγωγή χρήσιμων συμπερασμάτων.

Η ενότητα ολοκληρώνεται με την επισκόπηση της σημασίας κάθε χαρακτηριστικού για την πρόβλεψη, όπως προκύπτει από τους αλγορίθμους που προσφέρουν δυνατότητα ερμηνείας των αποτελεσμάτων. Με αυτό τον τρόπο, δίνεται έμφαση όχι μόνο στην αξιολόγηση της ακρίβειας, αλλά και στη χρησιμότητα των προτεινόμενων μεθόδων για τη λήψη αποφάσεων στον οικονομικό τομέα.

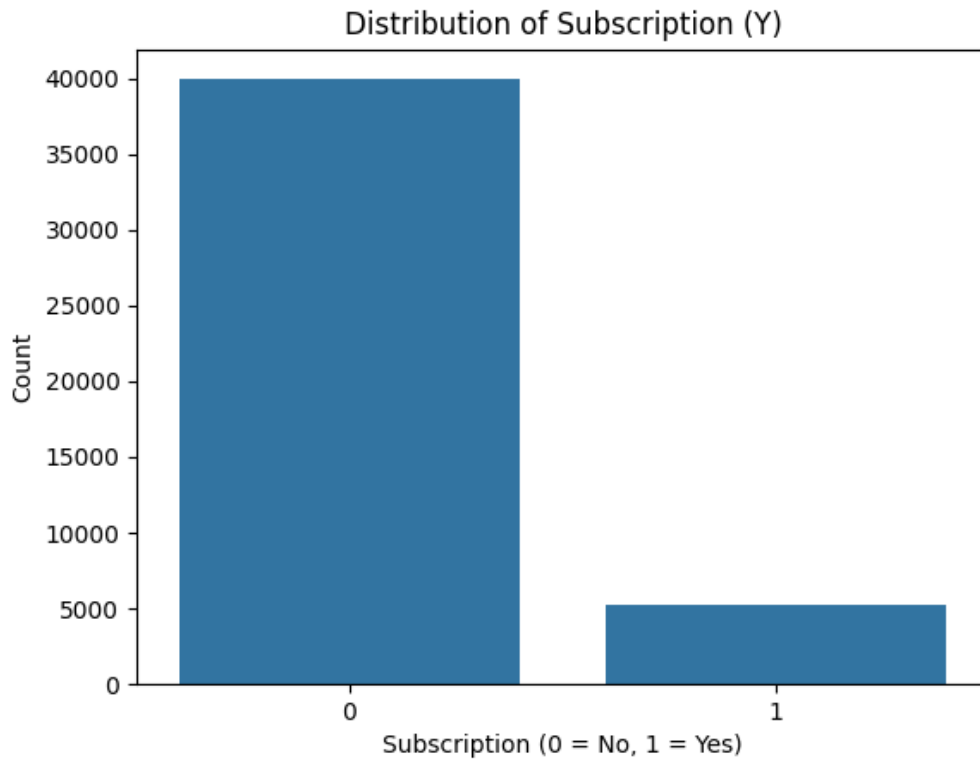
6.3.1. Bank Marketing Dataset

Γραφήματα

Η κατανόηση της φύσης των δεδομένων αποτελεί ένα από τα πιο σημαντικά βήματα σε οποιαδήποτε ανάλυση. Τα γραφήματα που παρουσιάζονται επιτρέπουν την εξερεύνηση της δομής και των ιδιαιτεροτήτων του συνόλου δεδομένων, αποκαλύπτοντας κρίσιμες πληροφορίες για τη συμπεριφορά των χαρακτηριστικών και τη σχέση τους με τη μεταβλητή στόχο. Μέσα από αυτά, εντοπίζονται προβλήματα, όπως η ανισορροπία των κατηγοριών ή η ύπαρξη πλεονάζουσας πληροφορίας μεταξύ των χαρακτηριστικών, τα οποία επηρεάζουν σημαντικά την απόδοση των μοντέλων. Στη συνέχεια, αναλύονται δύο βασικά γραφήματα: η κατανομή της μεταβλητής στόχου, που υπογραμμίζει το πρόβλημα της ανισορροπίας, και ο πίνακας συσχέτισης, που αποκαλύπτει τις σχέσεις μεταξύ των αριθμητικών χαρακτηριστικών.

1) Διάγραμμα Κατανομής Στόχου

Το γράφημα παρουσιάζει την κατανομή της μεταβλητής στόχου (Y) στο σύνολο δεδομένων.



Εικόνα 4: Κατανομή της μεταβλητής στόχου Subscription (Y) του συνόλου δεδομένων Bank Marketing.

Όπως παρατηρείται, αν και έχει ήδη προαναφερθεί στην ενότητα 1 (Δεδομένα), η μεταβλητή στόχος παίρνει δύο τιμές:

- 0: Αντιπροσωπεύει τους πελάτες που δεν εγγράφηκαν στην υπηρεσία.
- 1: Αντιπροσωπεύει τους πελάτες που εγγράφηκαν στην υπηρεσία.

Ο κάθετος άξονας (Count) δείχνει τον αριθμό των δειγμάτων σε κάθε κατηγορία, ενώ ο οριζόντιος άξονας (Subscription) περιλαμβάνει τις δύο κατηγορίες της μεταβλητής στόχου. Η κατανομή είναι εμφανώς ανισόροπη, με την κατηγορία 0 να υπερτερεί σημαντικά έναντι της κατηγορίας 1.

Η κατηγορία 0 συγκεντρώνει πάνω από 40.000 παρατηρήσεις, αντιπροσωπεύοντας περίπου το 89% του συνόλου των δεδομένων, ενώ η κατηγορία 1 περιλαμβάνει περίπου 5.000 παρατηρήσεις, καλύπτοντας το υπόλοιπο 11%. Αυτή η έντονη ανισοροπία δημιουργεί προκλήσεις για τα μοντέλα μηχανικής μάθησης. Τα μοντέλα τείνουν να προσαρμόζονται καλύτερα στην κυρίαρχη κατηγορία 0, καθώς αυτή αντιπροσωπεύεται καλύτερα στα δεδομένα, γεγονός που οδηγεί σε χαμηλή απόδοση για την υποεκπροσωπούμενη κατηγορία 1. Αν το πρόβλημα της ανισοροπίας δεν αντιμετωπιστεί κατάλληλα, το

μοντέλο μπορεί να εμφανίσει υψηλή συνολική ακρίβεια, αλλά θα αποτύχει να αναγνωρίσει σωστά περιπτώσεις της κατηγορίας 1.

Αυτή η κατανομή υποδεικνύει ότι η κατηγορία 0 είναι πιο συχνή και κυριαρχεί στο σύνολο δεδομένων, κάτι που μπορεί να αντανακλά επιχειρησιακές πραγματικότητες, όπως ότι οι περισσότεροι πελάτες δεν ενδιαφέρονται να εγγραφούν στην προσφερόμενη υπηρεσία. Αντίθετα, η κατηγορία 1, αν και λιγότερο συχνή, έχει μεγάλη επιχειρησιακή σημασία, καθώς αφορά τους πελάτες που είναι πιθανότερο να εγγραφούν. Η σωστή αναγνώριση αυτών των πελατών είναι κρίσιμη για τη βελτίωση της απόδοσης των καμπανιών μάρκετινγκ και την αύξηση της επιτυχίας της υπηρεσίας.

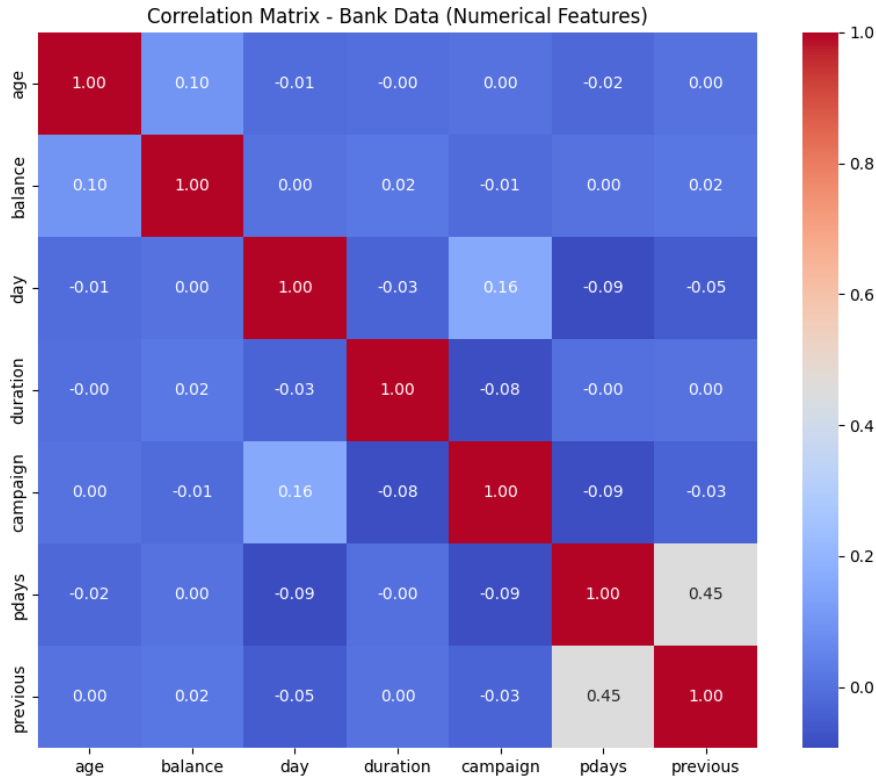
Η ανισορροπία αυτή αναδεικνύει την ανάγκη για τη χρήση τεχνικών εξισορρόπησης δεδομένων, όπως το SMOTE (Synthetic Minority Over-sampling Technique) και το undersampling. Το SMOTE μπορεί να δημιουργήσει συνθετικά δείγματα για την κατηγορία 1, αυξάνοντας το μέγεθός της, ενώ το undersampling μειώνει τις παρατηρήσεις της κατηγορίας 0, εξισορροπώντας έτσι τις δύο κατηγορίες [6]. Τέτοιες τεχνικές είναι απαραίτητες για να διασφαλιστεί ότι τα μοντέλα μηχανικής μάθησης θα δώσουν έμφαση στην κατηγορία 1, αποφεύγοντας την προκατάληψη υπέρ της κατηγορίας 0.

Η υπεροχή της κατηγορίας 0 μπορεί να οδηγήσει σε μοντέλα που επικεντρώνονται κυρίως στην κυρίαρχη κατηγορία, προβλέποντας με μεγάλη ακρίβεια τους πελάτες που δεν εγγράφονται, αλλά με χαμηλή αποτελεσματικότητα στον εντοπισμό των πελατών που εγγράφονται. Αυτό έχει ως αποτέλεσμα υψηλό recall για την κατηγορία 0 και χαμηλό recall για την κατηγορία 1, κάτι που είναι ανεπιθύμητο για πολλές εφαρμογές, όπως αυτή που εξετάζεται. Η κατανόηση της κατανομής της μεταβλητής στόχου είναι επομένως καθοριστική, καθώς επηρεάζει άμεσα τη στρατηγική που θα ακολουθηθεί για την επεξεργασία των δεδομένων και την εκπαίδευση των μοντέλων.

Συνολικά, το γράφημα υπογραμμίζει ένα σοβαρό πρόβλημα ανισορροπίας που πρέπει να αντιμετωπιστεί κατάλληλα για τη βελτίωση της απόδοσης των μοντέλων. Η σωστή επεξεργασία των δεδομένων και η χρήση εξειδικευμένων τεχνικών εξισορρόπησης είναι απαραίτητες για την επίτευξη πιο ακριβών προβλέψεων, ιδιαίτερα για την κατηγορία 1, που αποτελεί το κύριο ενδιαφέρον αυτής της ανάλυσης. Η κατανόηση της ανισορροπίας αυτής είναι θεμελιώδης για την εφαρμογή επιτυχών στρατηγικών μάρκετινγκ και τη μεγιστοποίηση της απόδοσης της υπηρεσίας.

2) Heatmap Συσχέτισης των αριθμητικών χαρακτηριστικών

Το γράφημα απεικονίζει τον πίνακα συσχέτισης (correlation matrix) μεταξύ των αριθμητικών χαρακτηριστικών του συνόλου δεδομένων, επιτρέποντας την ανάλυση της γραμμικής σχέσης μεταξύ αυτών των χαρακτηριστικών.



Εικόνα 5: Πίνακας συσχέτισης μεταξύ των αριθμητικών μεταβλητών του συνόλου δεδομένων Bank Marketing.

Κάθε κελί στον πίνακα συσχέτισης περιέχει την τιμή συσχέτισης Pearson, η οποία κυμαίνεται από -1 έως 1. Μια τιμή κοντά στο 1 υποδεικνύει ισχυρή θετική συσχέτιση, ενώ μια τιμή κοντά στο -1 υποδεικνύει ισχυρή αρνητική συσχέτιση. Μια τιμή κοντά στο 0 υποδεικνύει ότι τα χαρακτηριστικά δεν έχουν γραμμική συσχέτιση μεταξύ τους. Η κλίμακα χρωμάτων στο πλάι παρέχει μια οπτική αναπαράσταση αυτών των τιμών: οι αποχρώσεις του κόκκινου αντιπροσωπεύουν ισχυρή θετική συσχέτιση, ενώ οι αποχρώσεις του μπλε υποδεικνύουν αρνητική ή ανύπαρκτη συσχέτιση.

Η διαγώνιος του πίνακα περιέχει πάντα την τιμή 1, καθώς κάθε χαρακτηριστικό είναι τέλεια συσχετισμένο με τον εαυτό του. Παρατηρώντας τα υπόλοιπα στοιχεία του πίνακα, προκύπτουν σημαντικά

συμπεράσματα για τις σχέσεις μεταξύ των αριθμητικών χαρακτηριστικών. Για παράδειγμα, οι συσχετίσεις είναι γενικά αδύναμες ή ανύπαρκτες μεταξύ των περισσότερων χαρακτηριστικών, γεγονός που υποδηλώνει ότι δεν υπάρχει έντονη γραμμική εξάρτηση μεταξύ τους. Αυτή η έλλειψη ισχυρών συσχετίσεων μπορεί να σημαίνει ότι τα χαρακτηριστικά παρέχουν συμπληρωματική πληροφορία στο μοντέλο και δεν είναι ιδιαίτερα πλεονάζοντα.

Ένα από τα χαρακτηριστικά που ξεχωρίζουν είναι το `duration`, το οποίο φαίνεται να έχει τη μεγαλύτερη συσχέτιση με άλλα χαρακτηριστικά, αν και οι τιμές παραμένουν σχετικά χαμηλές. Αυτό είναι αναμενόμενο, καθώς η διάρκεια μιας τηλεφωνικής επικοινωνίας μπορεί να σχετίζεται με την πιθανότητα επιτυχίας μιας καμπάνιας. Παράλληλα, το `rdays` εμφανίζει μέτρια θετική συσχέτιση με το `previous` (τιμή συσχέτισης περίπου 0.45). Αυτό υποδεικνύει ότι ο αριθμός ημερών από την προηγούμενη επικοινωνία σχετίζεται με τον αριθμό των προηγούμενων επαφών, γεγονός που μπορεί να είναι χρήσιμο για την κατανόηση της συμπεριφοράς των πελατών.

Άλλα χαρακτηριστικά, όπως το `age`, το `balance` και το `campaign`, εμφανίζουν σχεδόν μηδενικές συσχετίσεις με τα υπόλοιπα, υποδηλώνοντας ότι δεν υπάρχει έντονη γραμμική σχέση. Αυτή η ανεξαρτησία είναι χρήσιμη σε μοντέλα όπως το Random Forest, τα οποία δεν απαιτούν χαρακτηριστικά με ισχυρή συσχέτιση για να αποδώσουν αποτελεσματικά.

Η απουσία υψηλών συσχετίσεων μεταξύ των περισσότερων χαρακτηριστικών είναι ένα θετικό στοιχείο, καθώς μειώνει τον κίνδυνο πολυδιγραμμικότητας (*multicollinearity*) στο μοντέλο. Η πολυδιγραμμικότητα μπορεί να δημιουργήσει προκλήσεις, ειδικά σε γραμμικά μοντέλα, ωστόσο, σε μοντέλα όπως το Random Forest ή το Gradient Boosting, δεν αποτελεί μεγάλο πρόβλημα. Παρ' όλα αυτά, η αδύναμη συσχέτιση μπορεί να υποδεικνύει ότι ορισμένα χαρακτηριστικά ενδέχεται να μην έχουν ισχυρή προβλεπτική δύναμη και μπορεί να είναι χρήσιμη μια ανάλυση χαρακτηριστικών για τον εντοπισμό αυτών που συμβάλλουν περισσότερο στη μεταβλητή στόχο.

Συνοψίζοντας, ο πίνακας συσχέτισης παρέχει σημαντικές πληροφορίες για τη γραμμική σχέση μεταξύ των αριθμητικών χαρακτηριστικών. Η ανάλυση δείχνει ότι οι περισσότερες συσχετίσεις είναι αδύναμες, γεγονός που μειώνει τον κίνδυνο πλεονάζουσας πληροφορίας μεταξύ των χαρακτηριστικών. Ωστόσο, χαρακτηριστικά όπως το `duration` και το ζεύγος `rdays-previous` παρουσιάζουν ορισμένες αξιοσημείωτες συσχετίσεις, οι οποίες μπορούν να αξιοποιηθούν περαιτέρω. Η πληροφορία αυτή μπορεί να καθοδηγήσει τη στρατηγική του μοντέλου, καθώς και τις προσεγγίσεις επιλογής και προεπεξεργασίας χαρακτηριστικών.

Μοντέλα Μηχανικής Μάθησης

A. Αρχικά, χρησιμοποιείται ο Random Forest Classifier για την αξιολόγηση της απόδοσης μέσω διαφορετικών τεχνικών, όπως το baseline μοντέλο, το SMOTE και το Undersampling. Ο Random Forest εφαρμόζεται με τις προκαθορισμένες παραμέτρους του Scikit-Learn, όπως `n_estimators=100`, `max_depth=None`, και `max_features='sqrt'`, χωρίς περαιτέρω παραμετροποίηση. Εξετάζεται πώς οι μέθοδοι εξισορρόπησης επηρεάζουν βασικά metrics, όπως το precision, το recall, το F1-score, η accuracy και το ROC-AUC, για τις δύο κατηγορίες της μεταβλητής στόχου. Η σύγκριση των αποτελεσμάτων αυτών επιτρέπει την εξαγωγή πολύτιμων συμπερασμάτων σχετικά με την ικανότητα του Random Forest να εντοπίζει σωστά τους πελάτες που είναι πιθανότερο να εγγραφούν, ενώ παράλληλα αναδεικνύονται οι προκλήσεις που δημιουργούνται από την ανισορροπία των δεδομένων και τις διαφορετικές στρατηγικές που εφαρμόστηκαν για την αντιμετώπισή της.

1) Baseline Model

Η αρχική αξιολόγηση του Random Forest έγινε στα πρωτογενή δεδομένα, χωρίς την εφαρμογή τεχνικών εξισορρόπησης. Τα αποτελέσματα αντικατοπτρίζουν την απόδοση του μοντέλου όταν τα δεδομένα εκπαίδευσης έχουν την αρχική τους ανισορροπία.

Class	Precision	Recall	F1-Score	Support
0	0.92	0.97	0.95	7952
1	0.67	0.40	0.50	1091
Accuracy			0.90	9043
Macro Avg	0.80	0.69	0.72	9043
Weighted Avg	0.89	0.90	0.89	9043

Table 1: Αποτελέσματα ταξινόμησης: Precision, Recall και F1-Score ανά κατηγορία για το baseline μοντέλο της Random Forest, καθώς και συνολικές επιδόσεις (Accuracy, Macro και Weighted Averages) στο σύνολο δεδομένων Bank Marketing.

Accuracy (Baseline): 0.9040

ROC-AUC (Baseline): 0.9239

Accuracy (Ακρίβεια): Η ακρίβεια μετρά το ποσοστό των σωστά προβλεφθέντων δειγμάτων (0 και 1) στο σύνολο των προβλέψεων. Η τιμή 90.4% φαίνεται εντυπωσιακή, αλλά πρέπει να εξεταστεί στο πλαίσιο της ανισορροπίας:

- Η κλάση 0 είναι πολύ πιο συχνή στα δεδομένα (οι περισσότεροι πελάτες δεν εγγράφονται).
- Το μοντέλο μαθαίνει να προβλέπει κυρίως την κλάση 0 για να πετύχει υψηλή ακρίβεια.
- Αυτό επιβεβαιώνεται και από την ανάλυση των παρακάτω metrics.

Precision (Ακρίβεια): Για την κατηγορία 0 (πελάτες που δεν εγγράφηκαν), η ακρίβεια είναι 92%, δηλαδή από όλες τις προβλέψεις που έγιναν για την κατηγορία αυτή, το 92% ήταν σωστές. Αυτό είναι αναμενόμενο, καθώς η κατηγορία 0 είναι η επικρατέστερη και το μοντέλο προσαρμόζεται εύκολα στις πολλές παρατηρήσεις της. Για την κατηγορία 1 (πελάτες που εγγράφηκαν), η ακρίβεια είναι 67%, υποδηλώνοντας ότι μόνο τα δύο τρίτα των προβλέψεων για την κατηγορία 1 είναι αληθείς. Η σχετικά χαμηλή ακρίβεια εδώ αντανακλά την ανισορροπία δεδομένων, καθώς το μοντέλο κάνει λιγότερες προβλέψεις για την κατηγορία 1.

Recall (Ανάκληση): Για την κατηγορία 0, το recall είναι 97%, δείχνοντας ότι το μοντέλο σχεδόν δεν χάνει καμία περίπτωση πελάτη που δεν εγγράφηκε. Αντίθετα, για την κατηγορία 1, η ανάκληση είναι μόνο 40%, πράγμα που σημαίνει ότι το μοντέλο εντοπίζει λιγότερους από τους μισούς πελάτες που πραγματικά εγγράφηκαν. Αυτό είναι το αποτέλεσμα της ανισορροπίας, καθώς το μοντέλο τείνει να μην «βλέπει» την κατηγορία 1.

F1-Score: Το F1-Score είναι ο αρμονικός μέσος του precision και του recall, παρέχοντας μια συνολική εκτίμηση της απόδοσης του μοντέλου για κάθε κατηγορία. Για την κατηγορία 0, το F1-Score είναι 95%, δείχνοντας εξαιρετική συνολική απόδοση. Ωστόσο, για την κατηγορία 1, το F1-Score είναι μόλις 50%, το οποίο υπογραμμίζει την περιορισμένη ικανότητα του μοντέλου να ισορροπήσει την ακρίβεια και την ανάκληση για την υποεκπροσωπούμενη κατηγορία.

Macro Average και Weighted Average: Η macro average υπολογίζει τον μέσο όρο των metrics μεταξύ των δύο κλάσεων, χωρίς να λαμβάνει υπόψη την ανισορροπία των δεδομένων. Το macro average F1-Score είναι 72%, το οποίο αντανακλά τη διαφορά στην απόδοση μεταξύ των δύο κατηγοριών. Αντίθετα, το weighted average υπολογίζει τον μέσο όρο, λαμβάνοντας υπόψη τη συχνότητα κάθε κλάσης. Το weighted F1-Score είναι 89%, γεγονός που επηρεάζεται κυρίως από την καλή απόδοση της κατηγορίας 0.

ROC-AUC (Receiver Operating Characteristic - Area Under Curve): Το ROC-AUC για το Baseline Model είναι 92.4%, μια τιμή που υποδηλώνει ισχυρή διακριτική ικανότητα του μοντέλου. Δηλαδή, το μοντέλο μπορεί σε μεγάλο βαθμό να ξεχωρίσει πελάτες που θα εγγραφούν από εκείνους που δεν θα εγγραφούν. Αυτό αποδεικνύει ότι, παρότι υπάρχει ανισορροπία μεταξύ των δύο κατηγοριών, το μοντέλο καταφέρνει να διατηρήσει μια καλή συνολική απόδοση ως προς τη διακριτική ισχύ, πιθανόν επειδή η Random Forest επωφελείται από το ensemble learning και είναι εγγενώς ανθεκτική σε τέτοιες προκλήσεις.

Συνολικά, το Baseline Model αποδίδει εξαιρετικά καλά για την κατηγορία 0, ενώ παρουσιάζει σημαντικές ελλείψεις στην κατηγορία 1. Η υψηλή ακρίβεια και το F1-Score για την πλειοψηφική κατηγορία ενδέχεται να αποκρύπτουν προβλήματα που σχετίζονται με την ανισορροπία των δεδομένων, όπως το χαμηλό recall για την κατηγορία 1 (μόλις 40%).

2) SMOTE

Η τεχνική SMOTE εφαρμόστηκε για να αντιμετωπιστεί η ανισορροπία στις κλάσεις του συνόλου δεδομένων, που προκαλεί την υπεροχή της κυρίαρχης κατηγορίας (0 - πελάτες που δεν εγγράφηκαν) έναντι της υποεκπροσωπούμενης κατηγορίας (1 - πελάτες που εγγράφηκαν). Η προσέγγιση αυτή δημιουργεί συνθετικά δείγματα για την υποεκπροσωπούμενη κατηγορία με βάση τα υπάρχοντα δεδομένα, ενισχύοντας έτσι τη δυνατότητα του μοντέλου να μάθει για την κατηγορία 1. Παρακάτω εξετάζονται αναλυτικά τα αποτελέσματα που προέκυψαν:

Class	Precision	Recall	F1-Score	Support
0	0.95	0.93	0.94	7952
1	0.55	0.63	0.59	1091
Accuracy			0.89	9043
Macro Avg	0.75	0.78	0.76	9043
Weighted Avg	0.90	0.89	0.90	9043

Table 2: Αποτελέσματα ταξινόμησης: Precision, Recall και F1-Score ανά κατηγορία για το μοντέλο της Random Forest, μετά την εφαρμογή της τεχνικής SMOTE, καθώς και συνολικές επιδόσεις (Accuracy, Macro και Weighted Averages) στο σύνολο δεδομένων Bank Marketing.

Accuracy (SMOTE): 0.8920

ROC-AUC (SMOTE): 0.9241

Η ακρίβεια (accuracy) του μοντέλου με SMOTE ανέρχεται στο 89.2%, υποδεικνύοντας καλή συνολική απόδοση. Ωστόσο, όπως ήδη προαναφέρθηκε, η ακρίβεια από μόνη της δεν είναι πάντα κατάλληλη για την αξιολόγηση ενός μοντέλου σε ανισόρροπα δεδομένα, καθώς μπορεί να επηρεαστεί από την κυρίαρχη κλάση. Η εφαρμογή του SMOTE μειώνει την έμφαση στην ακρίβεια και μετατοπίζει την προσοχή σε metrics όπως το recall και το F1-score, που αποτυπώνουν καλύτερα την απόδοση για την κατηγορία 1.

Αναλύοντας τα metrics ανά κλάση, για την κλάση 0 (πελάτες που δεν εγγράφηκαν), η ακρίβεια (precision) είναι **95%**, υποδηλώνοντας ότι το 95% των προβλέψεων της κλάσης αυτής είναι σωστές. Η ανάκληση (recall) είναι 93%, δείχνοντας ότι το μοντέλο εντοπίζει τη συντριπτική πλειονότητα των περιπτώσεων που ανήκουν στην κλάση 0. Ο συνδυασμός των δύο αυτών metrics αποδίδεται με το F1-score, που είναι 94%, και υποδηλώνει άριστη συνολική απόδοση για την κυρίαρχη κατηγορία.

Για την υποεκπροσωπούμενη κατηγορία 1 (πελάτες που εγγράφηκαν), το precision μειώνεται στο 55%, πράγμα που σημαίνει ότι μόνο το 55% των προβλέψεων της κλάσης 1 είναι σωστές. Αυτό είναι συνέπεια της χρήσης συνθετικών δεδομένων, που αυξάνουν την πιθανότητα ψευδών θετικών (false positives). Ωστόσο, η ανάκληση (recall) για την κλάση 1 αυξάνεται σημαντικά στο 63%, υποδηλώνοντας ότι το μοντέλο εντοπίζει πλέον σχεδόν τα δύο τρίτα των πραγματικών περιπτώσεων της κλάσης αυτής. Το F1-score για την κλάση 1 ανέρχεται στο 59%, μια σημαντική βελτίωση σε σχέση με το baseline μοντέλο (όπου ήταν μόλις 50%). Αυτή η αύξηση δείχνει την ικανότητα του SMOTE να βελτιώνει τη συνολική απόδοση για την υποεκπροσωπούμενη κατηγορία, χωρίς να «θυσιάζει» εντελώς την ακρίβεια.

Το ROC-AUC για το μοντέλο με SMOTE ανέρχεται στο 92.41%, υποδηλώνοντας ότι το μοντέλο διατηρεί εξαιρετική ικανότητα διάκρισης μεταξύ των πελατών που εγγράφονται και αυτών που δεν εγγράφονται. Η τιμή αυτή είναι ισοδύναμη με εκείνη του baseline μοντέλου, γεγονός που δείχνει ότι η εφαρμογή του SMOTE δεν επηρεάζει αρνητικά τη συνολική διακριτική ισχύ του μοντέλου. Το αποτέλεσμα αυτό επιβεβαιώνει ότι η στρατηγική εξισορρόπησης διατηρεί την ποιότητα ταξινόμησης στο σύνολο, παρά τις παρεμβάσεις στη δομή των δεδομένων.

Σε γενικές γραμμές, η εφαρμογή του SMOTE οδήγησε σε σημαντική αύξηση του recall και του F1-score για την κατηγορία 1, βελτιώνοντας την ικανότητα του μοντέλου να εντοπίζει πελάτες που είναι πιθανότερο να εγγραφούν. Αν και υπήρξε μια μικρή μείωση στο precision για την ίδια κατηγορία, αυτή η απώλεια θεωρείται αποδεκτή, καθώς το βασικό μέλημα της ανάλυσης είναι να μην χάνονται πραγματικές περιπτώσεις της κατηγορίας 1. Το SMOTE, επομένως, αποδεικνύεται μια αποτελεσματική στρατηγική για την αντιμετώπιση της ανισορροπίας δεδομένων.

3) Undersampling

Άλλη μία τεχνική που εφαρμόστηκε για την ανισορροπία των δεδομένων είναι η Undersampling, η οποία μειώνει το μέγεθος της υπερεκπροσωπούμενης κατηγορίας (0 - πελάτες που δεν εγγράφηκαν) ώστε να ταιριάζει με το μέγεθος της υποεκπροσωπούμενης κατηγορίας (1 - πελάτες που εγγράφηκαν). Με τη μέθοδο αυτή, το μοντέλο εκπαιδεύεται σε ένα πιο ισορροπημένο υποσύνολο δεδομένων, γεγονός που μπορεί να βελτιώσει την απόδοση για την κατηγορία 1, αλλά πιθανώς να μειώσει τη συνολική ακρίβεια. Τα αποτελέσματα που προέκυψαν αναλύονται παρακάτω:

Class	Precision	Recall	F1-Score	Support
0	0.98	0.82	0.90	7952
1	0.41	0.89	0.56	1091
Accuracy			0.83	9043
Macro Avg	0.70	0.86	0.73	9043
Weighted Avg	0.91	0.83	0.85	9043

Table 3: Αποτελέσματα ταξινόμησης: Precision, Recall και F1-Score ανά κατηγορία για το μοντέλο της Random Forest, μετά την εφαρμογή της τεχνικής Undersampling, καθώς και συνολικές επιδόσεις (Accuracy, Macro και Weighted Averages) στο σύνολο δεδομένων Bank Marketing.

Accuracy (Undersampling): 0.8310

ROC-AUC (Undersampling): 0.9240

Η ακρίβεια (accuracy) του μοντέλου με undersampling είναι 83.1%, σημαντικά χαμηλότερη από το baseline μοντέλο και τη μέθοδο SMOTE. Η μείωση αυτή είναι αναμενόμενη, καθώς η μέθοδος undersampling μειώνει τον όγκο δεδομένων της κατηγορίας 0, με αποτέλεσμα το μοντέλο να μην έχει πρόσβαση σε όλες τις πληροφορίες για την κυρίαρχη κλάση.

Αναλύοντας τα metrics ανά κλάση, για την κλάση 0 (πελάτες που δεν εγγράφηκαν), το precision είναι 98%, πράγμα που σημαίνει ότι σχεδόν όλες οι προβλέψεις για αυτή την κατηγορία είναι σωστές. Το recall, ωστόσο, μειώνεται στο 82%, γεγονός που δείχνει ότι το μοντέλο αποτυγχάνει να αναγνωρίσει όλες τις περιπτώσεις της κατηγορίας 0. Αυτό οφείλεται στη μείωση του πλήθους των δεδομένων της κλάσης 0 μέσω undersampling, που περιορίζει την ικανότητα του μοντέλου να μαθαίνει πλήρως τα χαρακτηριστικά

αυτής της κατηγορίας. Παρόλα αυτά, το F1-score για την κατηγορία 0 ανέρχεται στο 90%, δείχνοντας καλή συνολική απόδοση για την κυρίαρχη κλάση, αν και ελαφρώς μειωμένη σε σχέση με τις άλλες μεθόδους.

Για την υποεκπροσωπούμενη κατηγορία 1 (πελάτες που εγγράφηκαν), το precision είναι 41%, το οποίο υποδηλώνει ότι μόνο το 41% των προβλέψεων της κατηγορίας 1 είναι σωστές. Αν και το precision είναι σχετικά χαμηλό, το recall αυξάνεται στο 89%, υποδηλώνοντας ότι το μοντέλο εντοπίζει σχεδόν όλες τις περιπτώσεις πελατών που πραγματικά εγγράφηκαν. Το F1-score για την κατηγορία 1 είναι 56%, το οποίο είναι βελτιωμένο σε σχέση με το baseline μοντέλο αλλά χαμηλότερο από το F1-score με SMOTE. Αυτό δείχνει ότι η μέθοδος undersampling είναι αποτελεσματική στο να εστιάζει στην κατηγορία 1, αλλά με κόστος την απόδοση για την κατηγορία 0.

Το ROC-AUC του μοντέλου με undersampling ανέρχεται σε 92.40%, παραμένοντας στα ίδια υψηλά επίπεδα με το baseline και τη μέθοδο SMOTE. Η τιμή αυτή υποδεικνύει ότι το μοντέλο διατηρεί εξαιρετική συνολική διακριτική ικανότητα, ακόμα και όταν εκπαιδεύεται σε μειωμένο υποσύνολο της κατηγορίας 0. Αυτό καταδεικνύει ότι, παρά τη μείωση του όγκου της κυρίαρχης κατηγορίας, το μοντέλο εξακολουθεί να ξεχωρίζει αποτελεσματικά μεταξύ πελατών που θα εγγραφούν και αυτών που δεν θα εγγραφούν, χωρίς να επηρεάζεται η συνολική ικανότητα ταξινόμησης.

Σε γενικές γραμμές, η μέθοδος undersampling είναι ιδιαίτερα αποτελεσματική για την υποεκπροσωπούμενη κατηγορία 1. Το υψηλό recall της κατηγορίας 1 (89%) είναι κρίσιμο σε εφαρμογές όπου η αναγνώριση πελατών που θα εγγραφούν έχει μεγαλύτερη σημασία από την ακρίβεια των προβλέψεων. Ωστόσο, η μείωση του precision για την ίδια κατηγορία (41%) υποδηλώνει ότι το μοντέλο κάνει περισσότερες ψευδείς θετικές προβλέψεις. Επίσης, το undersampling μειώνει την απόδοση για την κατηγορία 0, καθώς το μοντέλο έχει λιγότερα δεδομένα για να μάθει πλήρως τα χαρακτηριστικά της.

Η μέθοδος undersampling προσφέρει μια ισχυρή επιλογή για προβλήματα με έντονη ανισορροπία, ειδικά όταν η υποεκπροσωπούμενη κατηγορία είναι το κύριο ενδιαφέρον. Ωστόσο, η χρήση της μπορεί να οδηγήσει σε απώλεια πληροφορίας για την κυρίαρχη κατηγορία, γεγονός που απαιτεί προσεκτική εξισορρόπηση μεταξύ recall και precision μέσω παραμετροποίησης.

B. Ενώ τα αποτελέσματα ανέδειξαν τη δυνατότητα του Random Forest να διαχειρίζεται ανισόρροπα δεδομένα και να επιτυγχάνει υψηλά επίπεδα ακρίβειας, παρατηρήθηκε ότι οι προκαθορισμένες παράμετροι που χρησιμοποιήθηκαν ενδεχομένως να μην αξιοποιούν πλήρως τις δυνατότητες του αλγορίθμου.

Επιπλέον, η χρήση μόνο ενός μοντέλου μηχανικής μάθησης περιορίζει τη δυνατότητα σύγκρισης διαφορετικών προσεγγίσεων και μεθόδων, γεγονός που μπορεί να επηρεάσει την εξαγωγή γενικότερων συμπερασμάτων για το σύνολο δεδομένων. Για τον λόγο αυτό, κρίθηκε αναγκαία η διεύρυνση της ανάλυσης με τη χρήση εναλλακτικών αλγορίθμων (όπως Logistic Regression και Gradient Boosting), καθώς και η εφαρμογή τεχνικών βελτιστοποίησης υπερπαραμέτρων, όπως το Grid Search, για να εντοπιστούν οι καλύτερες δυνατές ρυθμίσεις για τον Random Forest.

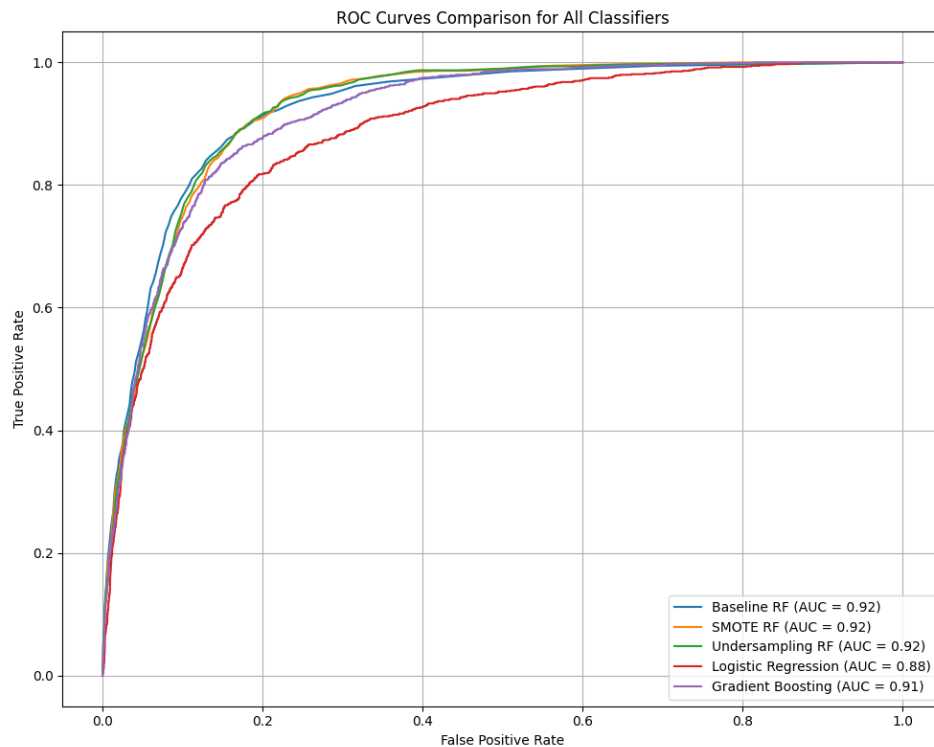
Στη συνέχεια, παρουσιάζονται πιο σύνθετες προσεγγίσεις, όπως η συστηματική σύγκριση μοντέλων, η δημιουργία ROC καμπυλών, και η αξιολόγηση της επίδρασης της παραμετροποίησης στα τελικά αποτελέσματα.

Γραφήματα

1) Διαγράμμα ROC Καμπυλών

Το παρακάτω διάγραμμα απεικονίζει τη σύγκριση των καμπυλών ROC για πέντε ταξινομητές. Συγκεκριμένα, περιλαμβάνονται τα εξής μοντέλα: Baseline Random Forest (χωρίς εξισορρόπηση), SMOTE Random Forest, Undersampling Random Forest, Logistic Regression (με SMOTE) και Gradient Boosting (με SMOTE). Η καμπύλη ROC (Receiver Operating Characteristic) αποτυπώνει τη σχέση μεταξύ του True Positive Rate και του False Positive Rate για κάθε μοντέλο, ενώ η επιφάνεια κάτω από την καμπύλη (AUC

- Area Under Curve) αποδίδει συνοπτικά την ικανότητα του εκάστοτε μοντέλου να διαχωρίζει σωστά τις δύο κατηγορίες της μεταβλητής στόχου.



Εικόνα 6: Σύγκριση ROC καμπυλών για τα μοντέλα ταξινόμησης υπό διαφορετικές στρατηγικές διαχείρισης ανισορροπίας δεδομένων: Baseline, SMOTE και Undersampling – Bank Marketing.

Αρχικά, το baseline μοντέλο του Random Forest, χωρίς καμία παραμετροποίηση, πέτυχε ROC-AUC ίσο με 0.92. Αυτό υποδεικνύει την εξαιρετική ικανότητα του μοντέλου να διαχωρίζει σωστά τις δύο κατηγορίες, ακόμα και όταν τα δεδομένα παραμένουν ανισόρροπα. Το γεγονός ότι ένα τόσο ισχυρό αποτέλεσμα επιτυγχάνεται χωρίς προσαρμογή υπερπαραμέτρων, οφείλεται στη φυσική ανθεκτικότητα του Random Forest σε προβλήματα ανισορροπίας, λόγω του τρόπου με τον οποίο χρησιμοποιεί δέντρα απόφασης σε συνδυασμό με δειγματοληψία bagging. Η προσθήκη τεχνικών εξισορρόπησης, όπως το SMOTE και το Undersampling, δεν οδήγησε σε αύξηση του ROC-AUC, που παρέμεινε σταθερό στο 0.92. Ωστόσο, είναι πιθανό ότι αυτές οι μέθοδοι να βελτίωσαν άλλα metrics, όπως το recall για την κατηγορία μειονότητας, το οποίο δεν αντικατοπτρίζεται απευθείας στο ROC-AUC.

Η χρήση του SMOTE, συγκεκριμένα, επικεντρώνεται στη δημιουργία συνθετικών δεδομένων για την κατηγορία μειονότητας, προσφέροντας στο μοντέλο περισσότερα δείγματα από αυτή την κατηγορία για να εκπαιδευτεί. Παρόλο που η συνολική ικανότητα διαχωρισμού των κατηγοριών δεν αυξήθηκε, η ενίσχυση του recall για την κατηγορία 1 είναι ιδιαίτερα σημαντική σε εφαρμογές όπου η σωστή

αναγνώριση της μειονότητας έχει μεγαλύτερη αξία από την ακρίβεια. Για παράδειγμα, σε περιπτώσεις όπου είναι κρίσιμο να εντοπιστούν όλοι οι πελάτες που ενδέχεται να εγγραφούν σε μια υπηρεσία, το SMOTE μπορεί να είναι προτιμητέο, παρά την πιθανή αύξηση των ψευδών θετικών (false positives).

Αντίστοιχα, το Undersampling, το οποίο μειώνει τα δεδομένα της κυρίαρχης κατηγορίας, παρέχει ισορροπημένα δεδομένα εκπαίδευσης, διευκολύνοντας το μοντέλο να αναγνωρίσει πιο αποτελεσματικά τις περιπτώσεις της μειονότητας. Ωστόσο, αυτή η μέθοδος συνοδεύεται από το κόστος της απώλειας πληροφοριών από την κατηγορία πλειονότητας, κάτι που μπορεί να επηρεάσει τη γενική ακρίβεια του μοντέλου. Το γεγονός ότι το ROC-AUC παρέμεινε αμετάβλητο (0.92) υποδηλώνει ότι το Random Forest είναι σε θέση να διαχειριστεί αυτή τη μείωση του όγκου δεδομένων χωρίς σημαντική απώλεια διακριτικής ικανότητας.

Στη συνέχεια, η ενσωμάτωση της Logistic Regression ως γραμμικού μοντέλου έδωσε ROC-AUC ίσο με 0.88, ελαφρώς χαμηλότερο από το Random Forest. Αυτό είναι αναμενόμενο, καθώς η Logistic Regression βασίζεται σε γραμμικές σχέσεις μεταξύ των χαρακτηριστικών και της μεταβλητής στόχου. Παρότι αποτελεί μια αξιόπιστη βάση σύγκρισης, τα αποτελέσματα δείχνουν ότι η πολυπλοκότητα και η μη γραμμική φύση του Random Forest το καθιστούν πιο αποτελεσματικό για το συγκεκριμένο dataset. Ωστόσο, η Logistic Regression είναι ιδιαίτερα χρήσιμη για την ερμηνευσιμότητα, παρέχοντας μια ξεκάθαρη εικόνα της σχέσης μεταξύ των χαρακτηριστικών και του στόχου.

Το Gradient Boosting παρουσίασε εξαιρετική απόδοση, με ROC-AUC ίσο με 0.91, πλησιάζοντας πολύ κοντά στο Random Forest. Η ικανότητά του να μαθαίνει από τα λάθη προηγούμενων βημάτων, μέσω της ενίσχυσης (boosting), το καθιστά ιδιαίτερα αποδοτικό για δεδομένα με σύνθετες σχέσεις και ανισορροπία. Παρά το γεγονός ότι η εκπαίδευσή του είναι πιο απαιτητική από άποψη χρόνου και υπολογιστικής ισχύος, το Gradient Boosting είναι συχνά προτιμώμενο σε εφαρμογές όπου η υψηλή ακρίβεια είναι ζωτικής σημασίας.

Βελτιστοποίηση Random Forest με τη χρήση Grid Search

Ο παρακάτω πίνακας παρουσιάζει τα αποτελέσματα ταξινόμησης του μοντέλου Random Forest, έπειτα από βελτιστοποίηση των υπερπαραμέτρων μέσω της μεθόδου Grid Search. Η εκπαίδευση του μοντέλου πραγματοποιήθηκε σε εξισορροπημένο σύνολο δεδομένων, καθώς εφαρμόστηκε η τεχνική SMOTE στο training set, ώστε να ενισχυθεί η παρουσία της μειοψηφικής κλάσης (κατηγορία 1 – πελάτες που εγγράφηκαν). Η αξιολόγηση του μοντέλου, ωστόσο, έγινε στο αρχικό test set που διατηρεί την

πραγματική ανισορροπία των κατηγοριών, γεγονός που επιτρέπει ρεαλιστική εκτίμηση της απόδοσης σε πραγματικές συνθήκες.

Class	Precision	Recall	F1-Score	Support
0	0.96	0.91	0.93	7952
1	0.51	0.70	0.59	1091
Accuracy			0.88	9043
Macro Avg	0.73	0.81	0.76	9043
Weighted Avg	0.90	0.88	0.89	9043

Table 4: Αποτελέσματα ταξινόμησης (Precision, Recall, F1-Score) για τον βελτιστοποιημένο Random Forest με Grid Search και εφαρμογή SMOTE στο σύνολο δεδομένων Bank Marketing.

Accuracy (Tuned RF): 0.8834

ROC-AUC (Tuned RF): 0.9224

Η διαδικασία βελτιστοποίησης πραγματοποιήθηκε με χρήση του Grid Search σε συνδυασμό με διασταυρούμενη επικύρωση 3-πλής αναδίπλωσης (3-fold cross-validation). Το πλέγμα τιμών που εξετάστηκε περιλάμβανε παραμέτρους για τον αριθμό των δέντρων (`n_estimators`), το μέγιστο βάθος κάθε δέντρου (`max_depth`) και τον αριθμό χαρακτηριστικών που εξετάζονται σε κάθε split (`max_features`). Οι βέλτιστες τιμές που εντοπίστηκαν ήταν οι εξής: `n_estimators=200`, `max_depth=20` και `max_features='sqrt'`.

Αναλύοντας τα αποτελέσματα του ταξινομητή, παρατηρείται αξιοσημείωτη βελτίωση στο recall της κατηγορίας 1 (πελάτες που εγγράφηκαν), το οποίο ανέρχεται στο 70%, έναντι 63% του μοντέλου με SMOTE χωρίς βελτιστοποίηση. Αυτή η αύξηση δηλώνει ότι το μοντέλο κατόρθωσε να εντοπίσει περισσότερους πελάτες που πραγματικά ανήκουν στη θετική κατηγορία, γεγονός σημαντικό για εφαρμογές όπου η πρόβλεψη της μειοψηφίας έχει μεγαλύτερη βαρύτητα. Το precision για την κατηγορία 1, ωστόσο, παρουσιάζει ελαφρά μείωση σε σχέση με προηγούμενα μοντέλα, διαμορφούμενο στο 0.51, κάτι που σημαίνει ότι ένα ποσοστό των προβλέψεων της κατηγορίας 1 αντιστοιχεί σε false positives.

Η σταθερότητα του F1-score στην κατηγορία 1, το οποίο παραμένει στο 0.59, επιβεβαιώνει ότι η συνολική ισορροπία μεταξύ precision και recall για αυτή την κατηγορία διατηρείται, παρά την αύξηση του recall. Το μοντέλο επιτυγχάνει συνολική ακρίβεια 88.3%, ενώ και οι μέσοι όροι (macro και weighted) κινούνται σε υψηλά επίπεδα, αντανακλώντας τη γενική σταθερότητα των προβλέψεων.

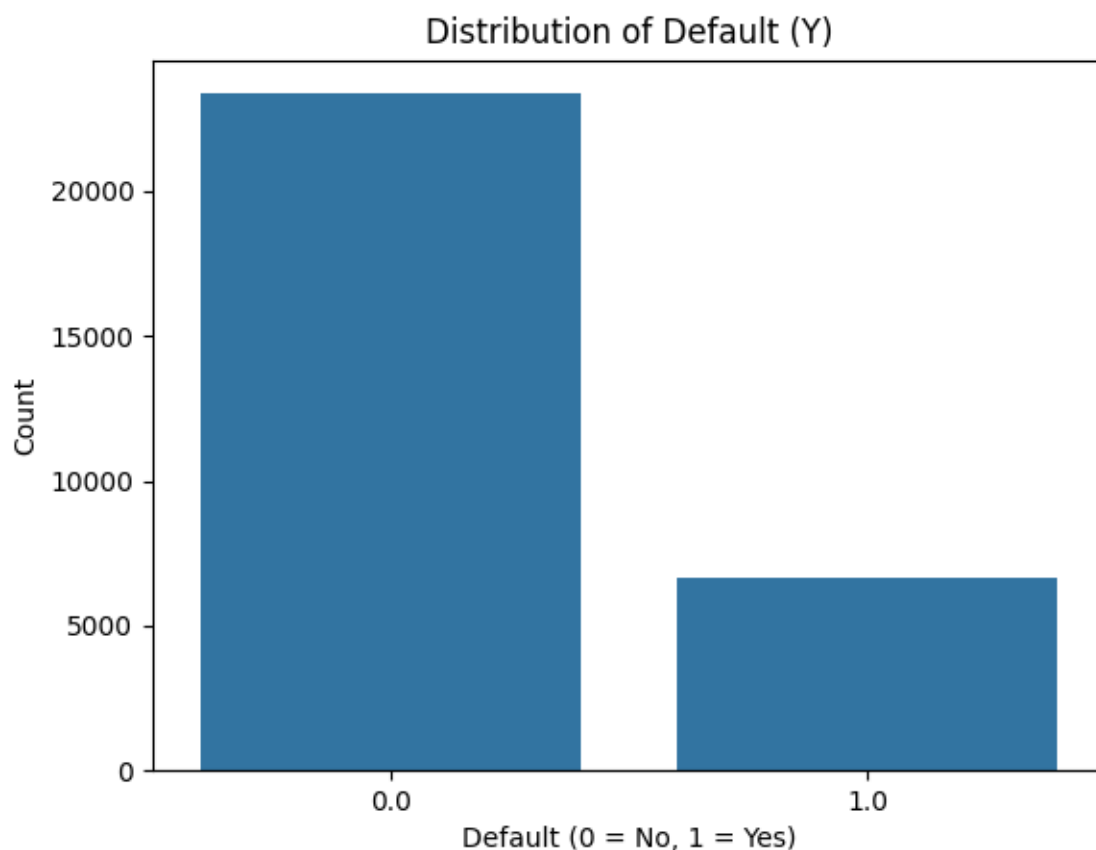
Η τιμή του ROC-AUC (0.9224) παραμένει ίδια με εκείνη του baseline μοντέλου, γεγονός που υποδηλώνει ότι η διακριτική ικανότητα του Random Forest δεν επηρεάστηκε ουσιαστικά από την παραμετροποίηση. Ωστόσο, η αλλαγή στη συμπεριφορά του μοντέλου απέναντι στην κατηγορία 1, και ειδικά η αύξηση του recall, καθιστά τη βελτιστοποίηση μέσω Grid Search σημαντική, καθώς επιτυγχάνεται καλύτερη κάλυψη της μειοψηφίας χωρίς απώλεια της συνολικής απόδοσης.

6.3.2. Credit Card Clients Dataset

Γραφήματα

1) Διάγραμμα Κατανομής Στόχου

Το γράφημα παρουσιάζει την κατανομή της μεταβλητής στόχου "Default (Y)" στο σύνολο δεδομένων. Η μεταβλητή αυτή υποδηλώνει εάν ένας πελάτης έχει αποτύχει να εκπληρώσει τις οικονομικές του υποχρεώσεις (τιμή 1) ή όχι (τιμή 0). Ο κάθετος άξονας αποτυπώνει τη συχνότητα εμφάνισης κάθε κατηγορίας, ενώ ο οριζόντιος άξονας απεικονίζει τις δύο τιμές της μεταβλητής στόχου.



Εικόνα 7: Κατανομή συχνότητων της μεταβλητής στόχου Default (Y) του συνόλου δεδομένων Credit Card Clients.

Όπως καταδεικνύεται από το διάγραμμα, υπάρχει σαφής ανισορροπία στις δύο κατηγορίες. Η κατηγορία 0, που αντιπροσωπεύει τους πελάτες που δεν έχουν αθετήσει τις υποχρεώσεις τους, είναι σημαντικά μεγαλύτερη από την κατηγορία 1, που περιλαμβάνει τους πελάτες που έχουν αποτύχει να ανταποκριθούν στις οικονομικές απαιτήσεις. Η κατανομή αυτή υποδεικνύει ότι η πλειοψηφία των παρατηρήσεων στο dataset ανήκει στην κατηγορία 0, ενώ η κατηγορία 1 είναι υποεκπροσωπούμενη.

Η εν λόγω ανισορροπία έχει σημαντικές επιπτώσεις για την ανάλυση των δεδομένων και την εφαρμογή μοντέλων μηχανικής μάθησης. Συγκεκριμένα, τα μοντέλα που εκπαιδεύονται σε τέτοιου είδους δεδομένα τείνουν να αποδίδουν καλύτερα στην πλειονότητα των παρατηρήσεων, δηλαδή στην κατηγορία 0, και συχνά αποτυγχάνουν να αναγνωρίσουν σωστά περιπτώσεις που ανήκουν στην κατηγορία 1. Αυτό μπορεί να οδηγήσει σε φαινομενικά υψηλή ακρίβεια του μοντέλου, αλλά με σοβαρές ελλείψεις στην πρόβλεψη της κατηγορίας που έχει το μεγαλύτερο επιχειρησιακό ενδιαφέρον.

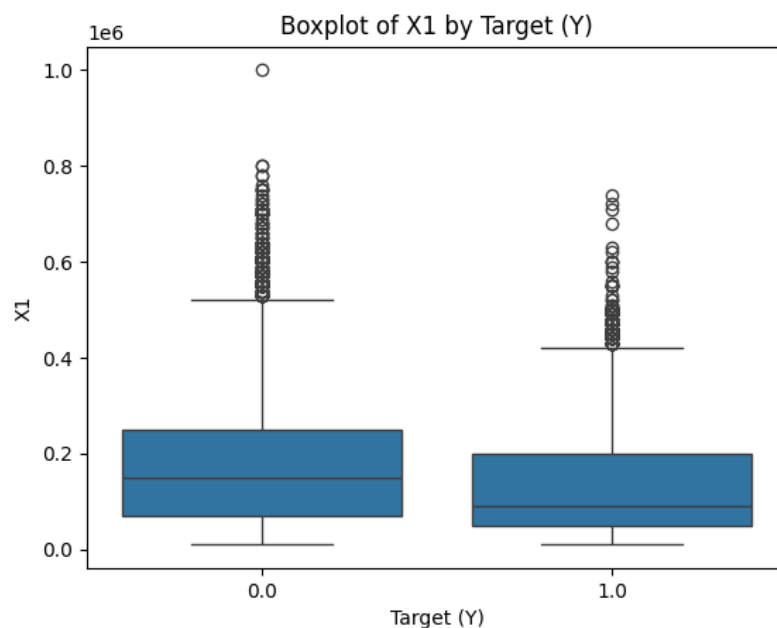
Η παρουσία αυτή της ασυμμετρίας καθιστά απαραίτητη την υιοθέτηση στρατηγικών για τη διαχείριση της ανισορροπίας των δεδομένων. Εργαλεία όπως το SMOTE και το Random UnderSampling χρησιμοποιούνται ευρέως για την εξισορρόπηση των κατηγοριών. Το SMOTE συνθέτει νέες παρατηρήσεις για την κατηγορία 1, αυξάνοντας το μέγεθος της υποεκπροσωπούμενης κατηγορίας, ενώ το undersampling μειώνει τον αριθμό των παρατηρήσεων της κατηγορίας 0, προσεγγίζοντας τη συχνότητα της κατηγορίας 1. Παράλληλα, η κατάλληλη επιλογή μετρικών αξιολόγησης, όπως το ROC-AUC και το F1 score, είναι καίριας σημασίας για την αντικειμενική αξιολόγηση της απόδοσης του μοντέλου, λαμβάνοντας υπόψη την ευαισθησία της κατηγορίας 1.

Από επιχειρησιακής άποψης, η κατηγορία 1, αν και σπανιότερη, έχει ιδιαίτερη σημασία, καθώς αφορά τους πελάτες που ενδέχεται να προκαλέσουν οικονομικές ζημιές λόγω αθέτησης υποχρεώσεων. Η σωστή αναγνώριση αυτών των πελατών μπορεί να βοηθήσει στη λήψη κατάλληλων μέτρων διαχείρισης κινδύνου και στη βελτίωση των οικονομικών αποτελεσμάτων.

Συμπερασματικά, η συγκεκριμένη οπτικοποίηση παρέχει μια σαφή εικόνα της διαφοράς στην κατανομή και αποτελεί βασικό βήμα για την κατανόηση των δεδομένων πριν από την εφαρμογή πιο σύνθετων μεθόδων ανάλυσης

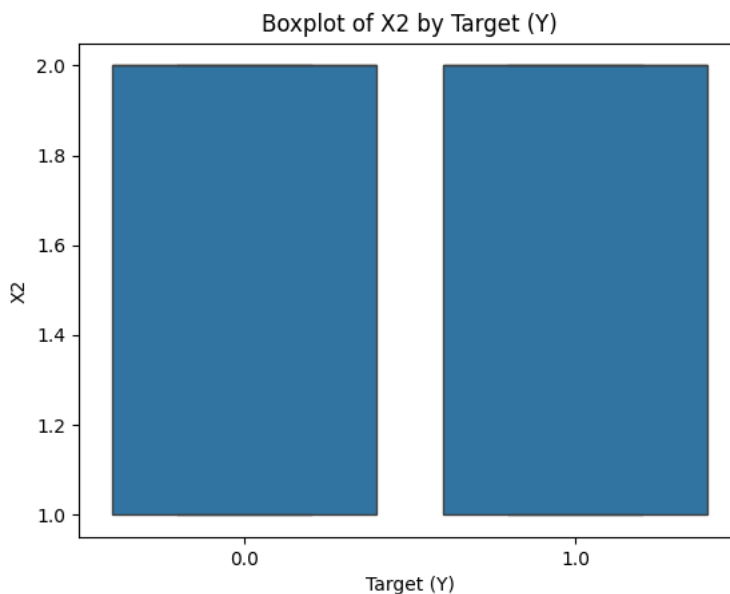
2) Διάγραμματα (Boxplots) Χαρακτηριστικών (X1 έως X23)

Ξεκινώντας με τη μεταβλητή LIMIT_BAL (X1), παρατηρείται ότι και για τις δύο κατηγορίες του στόχου (0 και 1) υπάρχουν αρκετές ακραίες τιμές. Η διάμεσος είναι ελαφρώς χαμηλότερη για την κατηγορία 1, ενώ η κατανομή παρουσιάζει μεγάλη διασπορά και στις δύο περιπτώσεις. Αυτό υποδηλώνει ότι το πιστωτικό όριο ενδέχεται να επηρεάζει την πιθανότητα αθέτησης, αλλά η επίδραση δεν είναι απόλυτα καθοριστική.



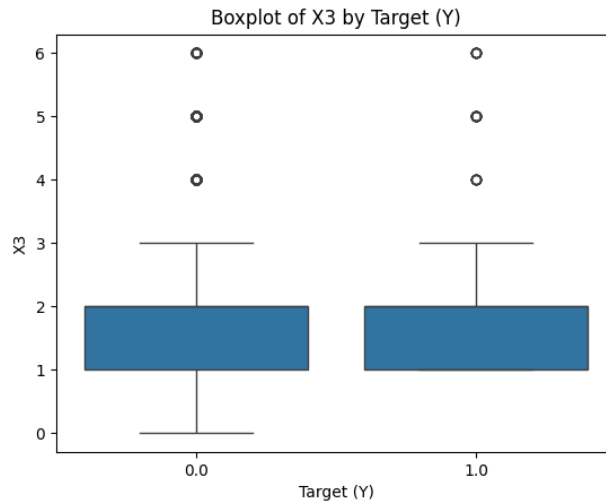
Εικόνα 8: Βoxplot της μεταβλητής X1 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.

Στη μεταβλητή SEX (X2), δεν παρατηρούνται ακραίες τιμές, ενώ η κατανομή είναι εξαιρετικά περιορισμένη. Αυτό υποδηλώνει ότι πρόκειται για δυαδικά κατηγορικά δεδομένα (1=Άνδρας, 2=Γυναίκα), με ελάχιστη διαφοροποίηση μεταξύ των κατηγοριών του στόχου, γεγονός που δείχνει ότι το φύλο πιθανώς δεν αποτελεί ισχυρό παράγοντα πρόβλεψης.



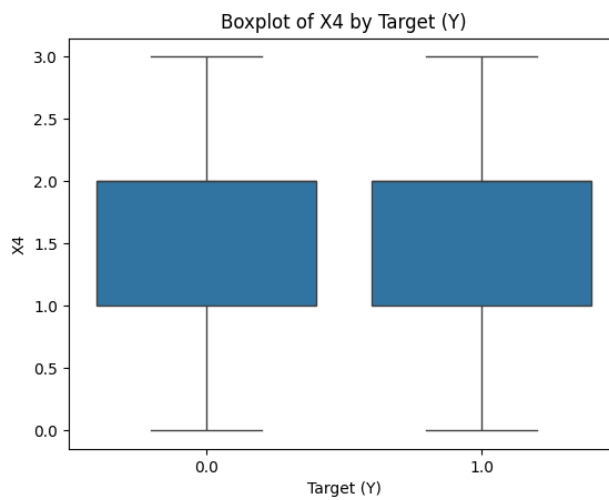
Εικόνα 9: Βoxplot της μεταβλητής X2 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.

Όσον αφορά τη μεταβλητή EDUCATION (X3), οι ακραίες τιμές εμφανίζονται και στις δύο κατηγορίες, αλλά η κατανομή είναι πιο συμμετρική σε σύγκριση με την LIMIT_BAL. Η διάμεσος παραμένει σχεδόν σταθερή και για τις δύο κατηγορίες, υποδηλώνοντας ότι το επίπεδο εκπαίδευσης δεν διαχωρίζει σαφώς τους δανειολήπτες που αθετούν ή όχι.



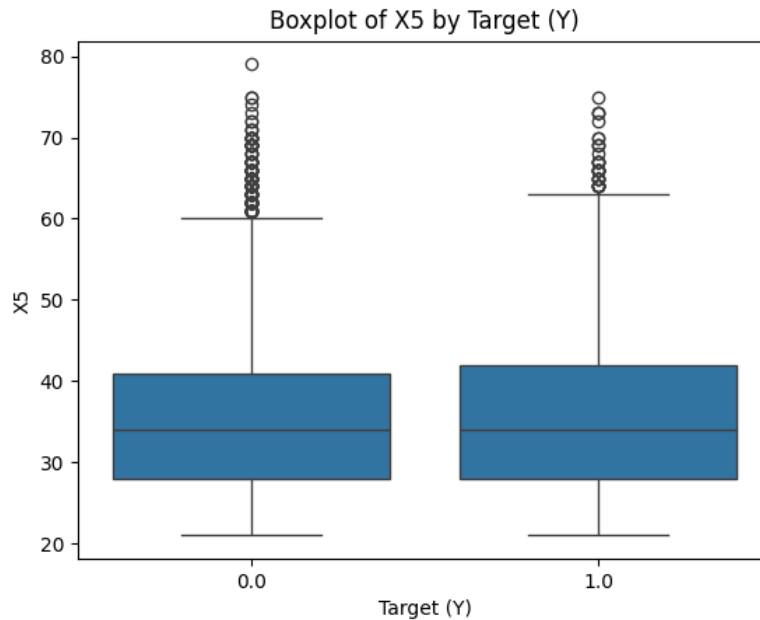
Εικόνα 10: Boxplot της μεταβλητής X3 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.

Η μεταβλητή MARRIAGE (X4) παρουσιάζει σχεδόν πανομοιότυπη κατανομή για τις δύο κατηγορίες του στόχου, χωρίς εμφανείς ακραίες τιμές. Αυτό υποδηλώνει ότι η οικογενειακή κατάσταση πιθανόν να μην έχει ισχυρή διακριτική δύναμη ως προς την αθέτηση πληρωμών.



Εικόνα 11: Boxplot της μεταβλητής X4 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.

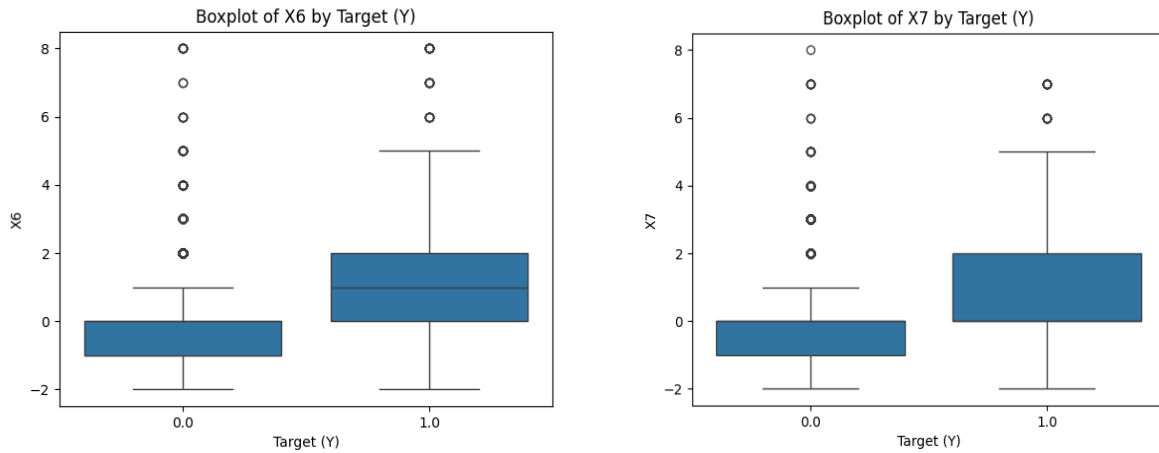
Για τη μεταβλητή AGE (X5), παρατηρούνται αρκετές ακραίες τιμές και στις δύο κατηγορίες, με τη διάμεσο να παραμένει σχεδόν αμετάβλητη. Ωστόσο, η διασπορά φαίνεται ελαφρώς μεγαλύτερη για την κατηγορία 0, κάτι που μπορεί να υποδηλώνει ότι οι νεότεροι δανειολήπτες παρουσιάζουν μεγαλύτερη πιθανότητα αθέτησης.



Εικόνα 12: Βοξπλοτ της μεταβλητής X5 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.

Οι μεταβλητές PAY_0 (X6) και PAY_2 (X7), που αντιπροσωπεύουν την κατάσταση καθυστερημένων πληρωμών στους τελευταίους μήνες, παρουσιάζουν σημαντική διαφοροποίηση στη διάμεσο μεταξύ των δύο κατηγοριών του στόχου. Η κατηγορία 1 (αθέτηση) εμφανίζει υψηλότερες τιμές, κάτι που υποδηλώνει ότι οι προηγούμενες καθυστερήσεις πληρωμών αποτελούν ισχυρό δείκτη πρόβλεψης αθέτησης.

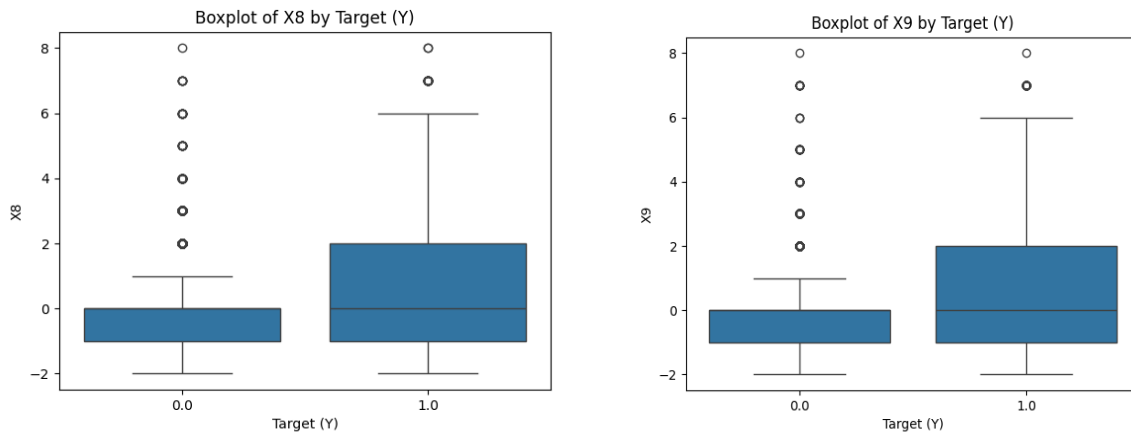
Εικόνα 13: Βoxplot της μεταβλητής X6 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.



Εικόνα 14: Βoxplot της μεταβλητής X7 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.

Ανάλογη τάση παρατηρείται και στις μεταβλητές PAY_3 (X8) και PAY_4 (X9), όπου η κατηγορία 1 εμφανίζει υψηλότερες διαμέσους σε σύγκριση με την κατηγορία 0, υποδεικνύοντας ότι οι συνεχείς καθυστερήσεις πληρωμών ενδέχεται να σχετίζονται θετικά με την πιθανότητα αθέτησης.

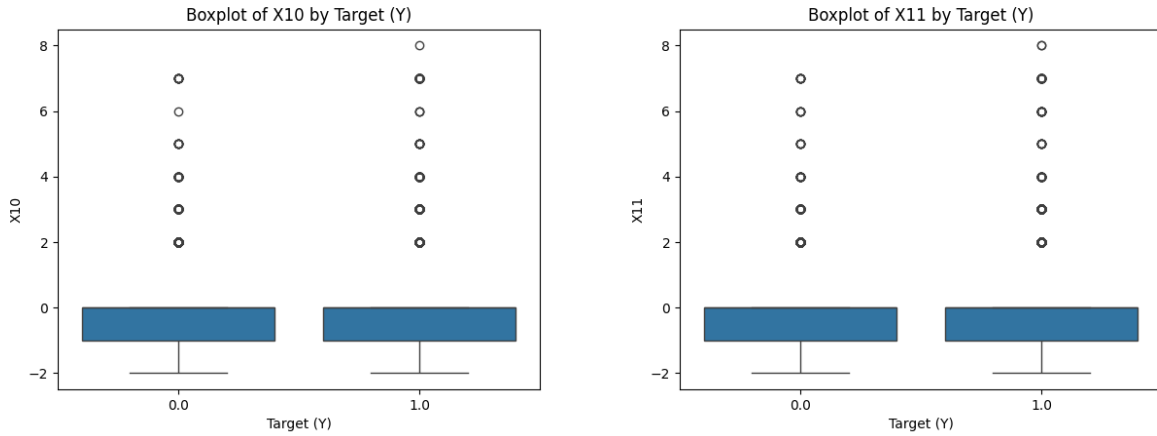
Εικόνα 15: Βoxplot της μεταβλητής X8 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.



Εικόνα 16: Βoxplot της μεταβλητής X9 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.

Η κατανομή των μεταβλητών PAY_5 (X10) και PAY_6 (X11) είναι σχεδόν ταυτόσημη για τις δύο κατηγορίες, με αρκετές ακραίες τιμές αλλά χωρίς σαφή διαφοροποίηση στη διάμεσο. Αυτό μπορεί να υποδηλώνει ότι οι πληρωμές σε μεταγενέστερους μήνες έχουν λιγότερη επιρροή στην έκβαση, πιθανώς επειδή οι προηγούμενες καθυστερήσεις έχουν ήδη ισχυρό αντίκτυπο.

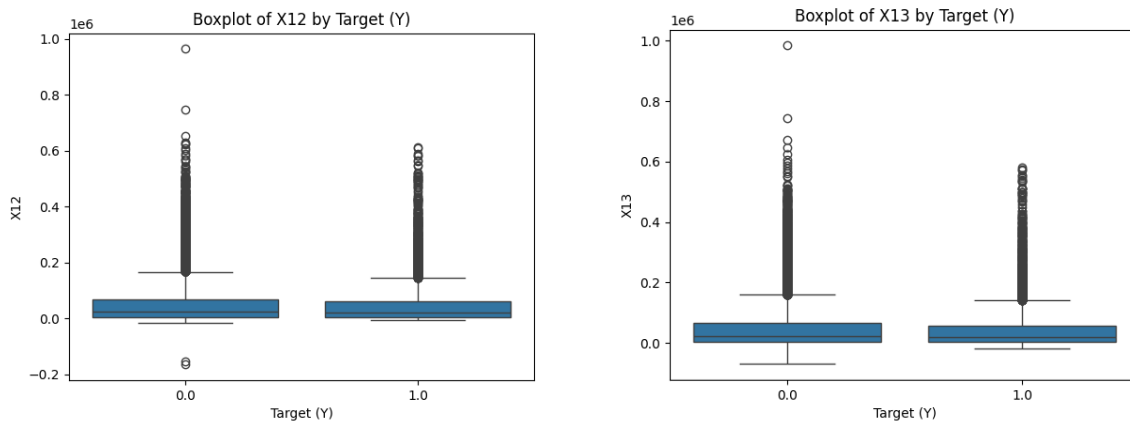
Εικόνα 17: Βοξplot της μεταβλητής X10 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.



Εικόνα 18: Βοξplot της μεταβλητής X11 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.

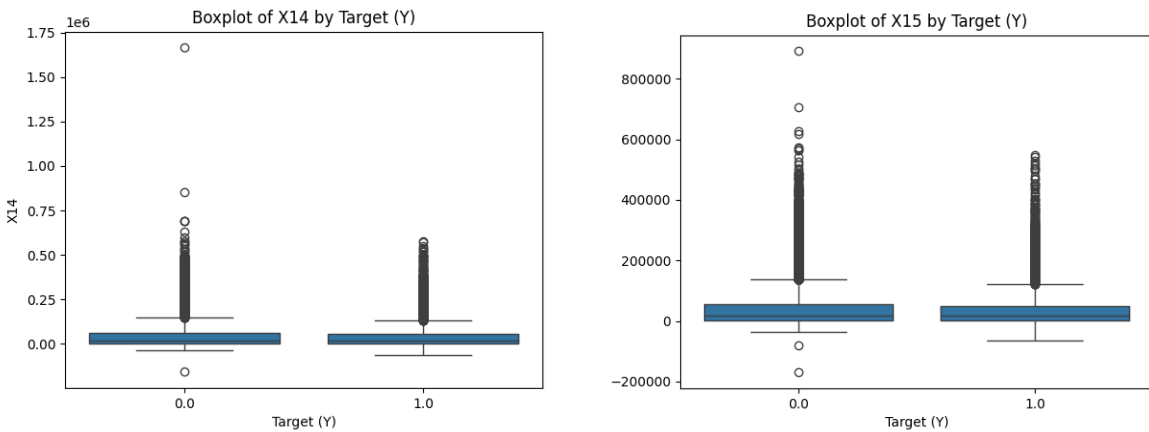
Στις μεταβλητές BILL_AMT1 έως BILL_AMT6 (X12 έως X17), που αντιπροσωπεύουν τα ποσά των οφειλών για τους τελευταίους έξι μήνες, παρατηρούνται μεγάλες διασπορές και πολλές ακραίες τιμές. Αν και οι διαφορές στις διαμέσους μεταξύ των κατηγοριών δεν είναι πάντα εμφανείς, η παρουσία ακραίων τιμών μπορεί να επηρεάζει τη σταθερότητα των μοντέλων και να απαιτεί περαιτέρω επεξεργασία, όπως κανονικοποίηση ή αφαίρεση των ακραίων τιμών.

Εικόνα 19: Βοξplot της μεταβλητής X12 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.



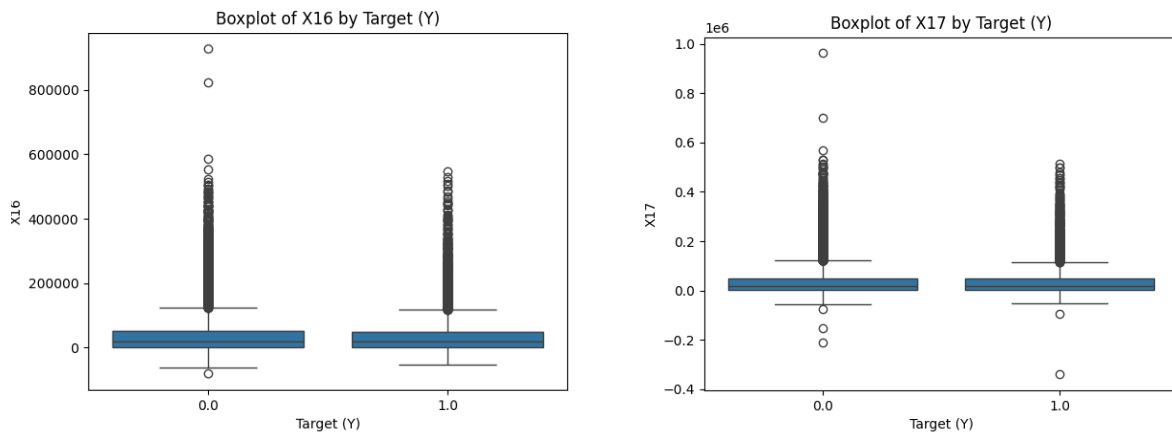
Εικόνα 20: Βοξplot της μεταβλητής X13 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.

Εικόνα 21: Βοχplot της μεταβλητής X14 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.



Εικόνα 22: Βοχplot της μεταβλητής X15 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.

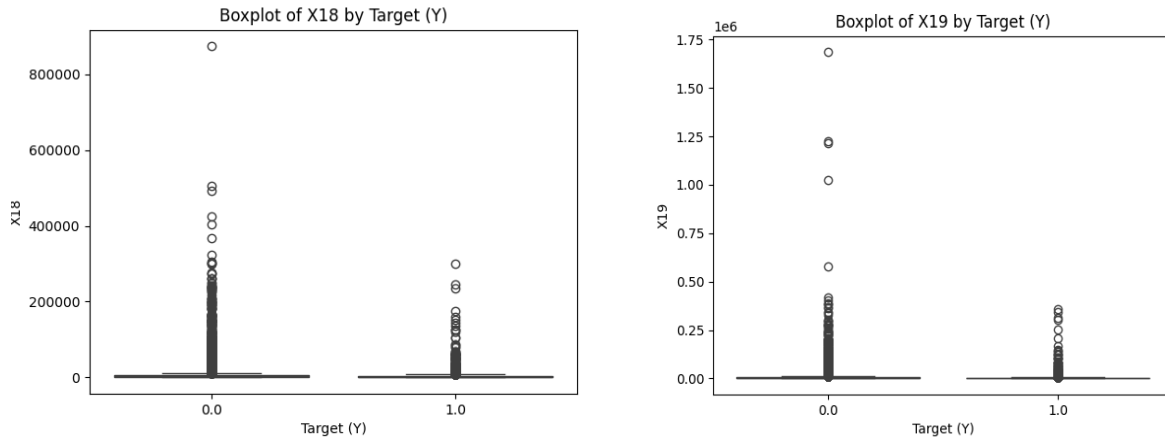
Εικόνα 23: Βοχplot της μεταβλητής X16 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.



Εικόνα 24: Βοχplot της μεταβλητής X17 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.

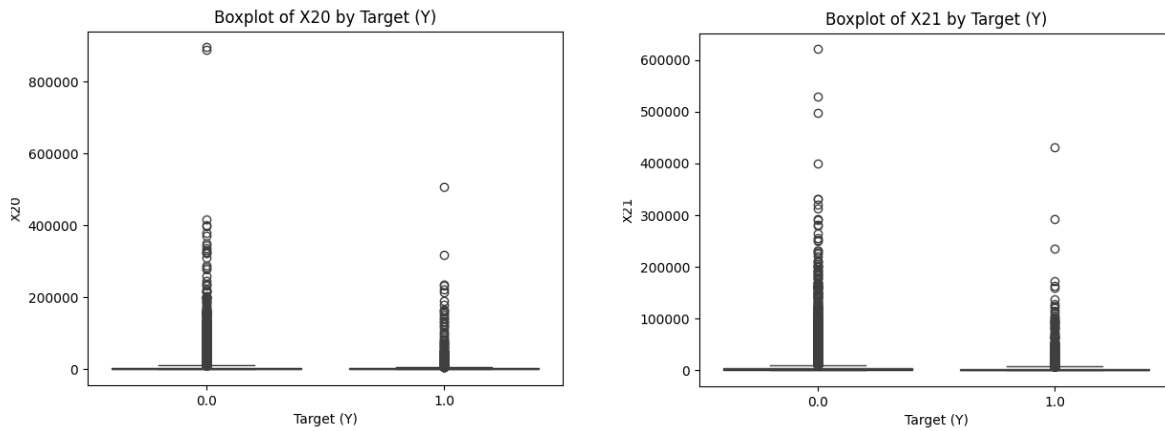
Οι μεταβλητές PAY_AMT1 έως PAY_AMT6 (X18 έως X23), που αφορούν τα ποσά πληρωμών στους τελευταίους έξι μήνες, εμφανίζουν επίσης σημαντικές ακραίες τιμές, με την κατηγορία 0 να παρουσιάζει μεγαλύτερες πληρωμές σε σύγκριση με την κατηγορία 1. Αυτή η τάση ενδέχεται να υποδηλώνει αρνητική συσχέτιση αυτών των μεταβλητών με την αθέτηση, καθώς οι μεγαλύτερες πληρωμές υποδηλώνουν υπεύθυνη οικονομική συμπεριφορά.

Εικόνα 25: Βoxplot της μεταβλητής X18 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.



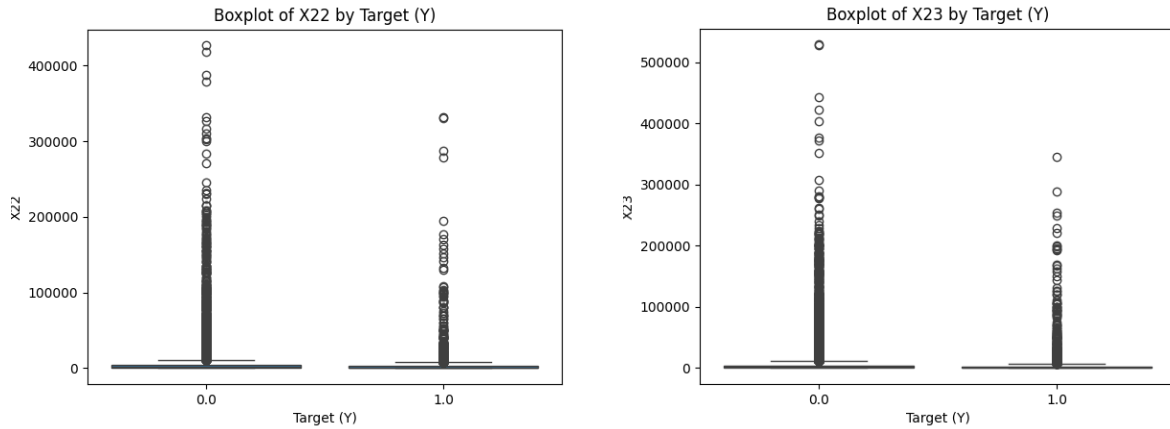
Εικόνα 26: Βoxplot της μεταβλητής X19 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.

Εικόνα 27: Βoxplot της μεταβλητής X20 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.



Εικόνα 28: Βoxplot της μεταβλητής X21 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.

Εικόνα 29: Βoxplot της μεταβλητής X22 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.

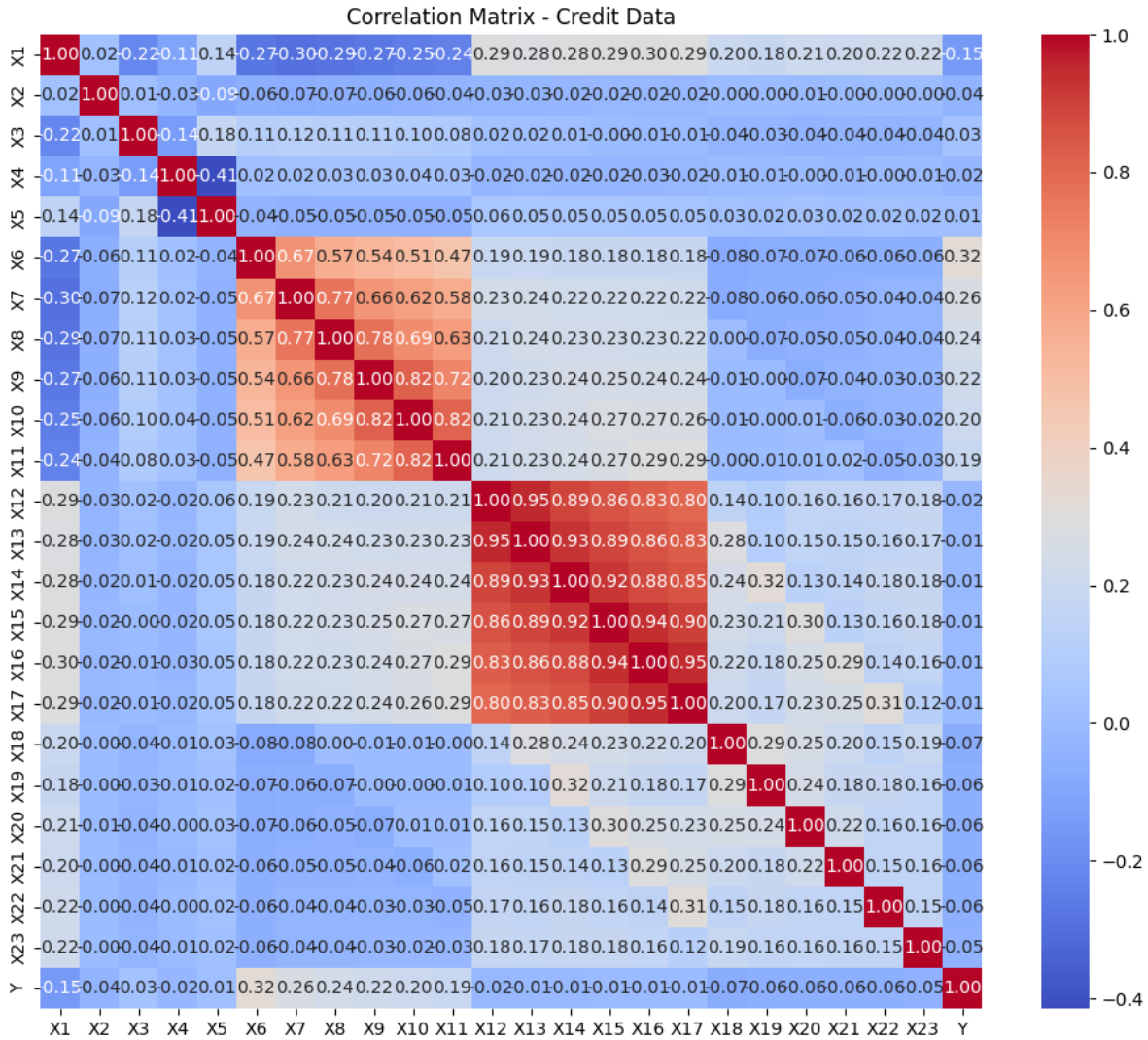


Εικόνα 30: Βoxplot της μεταβλητής X23 κατά τις δύο κατηγορίες της εξαρτημένης μεταβλητής στόχου Y.

Συνολικά, Οι μεταβλητές PAY_0 έως PAY_4 (X6 έως X9) φαίνεται να παρουσιάζουν πιο εμφανή διαφοροποίηση ως προς τον στόχο, καθιστώντας τις ιδιαίτερα χρήσιμες για μοντελοποίηση. Από την άλλη πλευρά, οι μεταβλητές SEX (X2), MARRIAGE (X4), PAY_5 (X10) και PAY_6 (X11) παρουσιάζουν μικρή διαφοροποίηση, γεγονός που μπορεί να υποδηλώνει χαμηλότερη σημασία για την πρόβλεψη αθέτησης. Η παρουσία ακραίων τιμών σε πολλές μεταβλητές, ιδίως στα ποσά οφειλών και πληρωμών, υποδηλώνει την ανάγκη για προσεκτική επεξεργασία, όπως κανονικοποίηση ή αφαίρεση ακραίων τιμών, πριν από την ανάπτυξη οποιουδήποτε μοντέλου πρόβλεψης.

3) Πίνακας Συσχέτισης Χαρακτηριστικών με την Αθέτηση Πληρωμών

Η ανάλυση του heatmap συσχέτισης για το dataset των πιστωτικών καρτών αποκαλύπτει σημαντικές σχέσεις μεταξύ των μεταβλητών και του στόχου, ο οποίος είναι η αθέτηση πληρωμών (Y).



Εικόνα 31: Πίνακας συσχέτισης των αριθμητικών μεταβλητών του συνόλου δεδομένων Credit Card Clients.

Αρχίζοντας με τη μεταβλητή στόχο (Y), παρατηρείται ότι η ισχυρότερη θετική συσχέτιση εμφανίζεται με τη μεταβλητή PAY_0 (X6), που αντιπροσωπεύει την κατάσταση καθυστερημένων πληρωμών του πιο πρόσφατου μήνα. Η συσχέτιση φτάνει το 0.32, υποδηλώνοντας ότι οι καθυστερήσεις πληρωμών αποτελούν έναν από τους σημαντικότερους παράγοντες που σχετίζονται με την πιθανότητα αθέτησης. Παρόμοια τάση παρατηρείται και για τις μεταβλητές PAY_2 έως PAY_5 (X7 έως X10), με συσχετίσεις που κυμαίνονται από 0.26 έως 0.24, επιβεβαιώνοντας ότι οι καθυστερήσεις πληρωμών στους προηγούμενους μήνες επηρεάζουν επίσης την έκβαση.

Αντίθετα, οι μεταβλητές που σχετίζονται με τις πληρωμές, συγκεκριμένα οι PAY_AMT1 έως PAY_AMT6 (X18 έως X23), παρουσιάζουν αρνητικές συσχετίσεις με την αθέτηση, με τιμές γύρω στο -0.06. Αυτό

σημαίνει ότι όσο περισσότερα πληρώνει κάποιος, τόσο μειώνεται η πιθανότητα να αθετήσει την πληρωμή του.

Η ανάλυση αποκαλύπτει επίσης πολύ ισχυρές συσχετίσεις μεταξύ ορισμένων μεταβλητών. Οι μεταβλητές που σχετίζονται με τα ποσά οφειλών, δηλαδή οι BILL_AMT1 έως BILL_AMT6 (X12 έως X17), εμφανίζουν πολύ υψηλή θετική συσχέτιση μεταξύ τους, με τιμές που κυμαίνονται από 0.89 έως 0.95. Αυτή η συσχέτιση είναι λογική, καθώς τα οφειλόμενα ποσά σε συνεχόμενους μήνες είναι συχνά παρόμοια. Ωστόσο, αυτή η υψηλή συσχέτιση μπορεί να οδηγήσει σε πολυπαραγοντικότητα (multicollinearity), γεγονός που μπορεί να επηρεάσει αρνητικά την απόδοση κάποιων μοντέλων. Ενδεχομένως να χρειαστεί να συνδυαστούν ή να μειωθούν αυτές οι μεταβλητές για να περιοριστεί αυτό το φαινόμενο.

Επιπλέον, το πιστωτικό όριο (LIMIT_BAL - X1) παρουσιάζει αρνητική συσχέτιση με την αθέτηση πληρωμών (Y), με τιμή -0.15. Αυτό υποδηλώνει ότι οι πελάτες με υψηλότερο πιστωτικό όριο είναι λιγότερο πιθανό να αθετήσουν τις πληρωμές τους. Παράλληλα, το LIMIT_BAL παρουσιάζει επίσης αρνητική συσχέτιση με τις μεταβλητές καθυστερημένων πληρωμών (PAY_0 έως PAY_6), με τιμές που κυμαίνονται από -0.27 έως -0.19, επιβεβαιώνοντας ότι οι πελάτες με μεγαλύτερο πιστωτικό όριο τείνουν να έχουν λιγότερες καθυστερήσεις στις πληρωμές τους.

Συνολικά, η ανάλυση του heatmap αποκαλύπτει ότι οι καθυστερήσεις πληρωμών (PAY_0 έως PAY_5) αποτελούν τους σημαντικότερους παράγοντες που σχετίζονται με την αθέτηση, ενώ οι μεγαλύτερες πληρωμές (PAY_AMT1 έως PAY_AMT6) συνδέονται με μειωμένη πιθανότητα αθέτησης. Οι ισχυρές συσχετίσεις μεταξύ των μεταβλητών οφειλών (BILL_AMT1 έως BILL_AMT6) υποδεικνύουν την ανάγκη για προσεκτική διαχείριση, ώστε να αποφευχθεί η πολυπαραγοντικότητα στα μοντέλα πρόβλεψης. Το πιστωτικό όριο (LIMIT_BAL) επίσης διαδραματίζει σημαντικό ρόλο, καθώς συσχετίζεται αρνητικά με την αθέτηση, υποδηλώνοντας ότι οι πελάτες με μεγαλύτερα όρια είναι πιο αξιόπιστοι.

Μοντέλα Μηχανικής Μάθησης

A. Μετά την ανάλυση των γραφημάτων και την εξερεύνηση της κατανομής των χαρακτηριστικών, αναδείχθηκαν σημαντικές παρατηρήσεις σχετικά με την παρουσία ακραίων τιμών (outliers) σε συγκεκριμένες μεταβλητές, όπως τα ποσά χρέωσης και πληρωμής (BILL_AMT και PAY_AMT). Οι ακραίες αυτές τιμές ενδέχεται να επηρεάσουν αρνητικά την απόδοση των μοντέλων μηχανικής μάθησης, προκαλώντας μεροληψία στις προβλέψεις και υποβαθμίζοντας τη γενίκευση των αποτελεσμάτων. Λαμβάνοντας υπόψη την πιθανή επίδραση των outliers, κρίθηκε σκόπιμη η εφαρμογή λογαριθμικού μετασχηματισμού σε επιλεγμένα χαρακτηριστικά (X1, X12-X23), με στόχο τη μείωση της ασυμμετρίας

των δεδομένων και την ενίσχυση της σταθερότητας των μοντέλων. Στη συνέχεια, ακολούθησε η κανονικοποίηση των δεδομένων και η διαχείριση της ανισορροπίας των κατηγοριών μέσω τεχνικών όπως το SMOTE και το undersampling.

Η ανάλυση που ακολουθεί επικεντρώνεται στην αξιολόγηση της απόδοσης αυτών των τεχνικών εξισορρόπησης σε σύγκριση με ένα baseline μοντέλο. Ο Random Forest Classifier εφαρμόζεται σε όλες τις περιπτώσεις, ακολουθώντας την ίδια μεθοδολογική προσέγγιση που υιοθετήθηκε και για το «Bank Marketing Dataset», προκειμένου να διατηρηθεί η συνοχή και η συγκρισιμότητα των αποτελεσμάτων. Μέσα από αυτή τη διαδικασία, επιδιώκεται η αξιολόγηση της αποτελεσματικότητας των διαφορετικών στρατηγικών εξισορρόπησης και η αποτίμηση της συμβολής τους στη βελτίωση της πρόβλεψης αθέτησης πληρωμών.

1) Baseline Model

Η αρχική αξιολόγηση του dataset «Default of Credit Card Clients» πραγματοποιήθηκε μέσω του baseline μοντέλου, το οποίο εκπαιδεύτηκε χωρίς την εφαρμογή τεχνικών εξισορρόπησης. Τα αποτελέσματα του μοντέλου προσφέρουν μια καλή αφετηρία για την κατανόηση της δυναμικής του προβλήματος και της συμπεριφοράς των δεδομένων.

Class	Precision	Recall	F1-Score	Support
0	0.84	0.94	0.89	4687
1	0.64	0.37	0.47	1313
Accuracy			0.82	6000
Macro Avg	0.74	0.66	0.68	6000
Weighted Avg	0.80	0.82	0.80	6000

Table 5: Αποτελέσματα ταξινόμησης του baseline μοντέλου Random Forest χωρίς εξισορρόπηση δεδομένων για το σύνολο δεδομένων Credit Card Clients.

Accuracy (Baseline): 0.8166666666666667

ROC-AUC (Baseline): 0.7548158272195898

Η συνολική ακρίβεια του baseline μοντέλου ανέρχεται στο 81,67%, υποδεικνύοντας ότι περισσότερο από το 80% των προβλέψεων του μοντέλου είναι σωστές. Αν και η τιμή αυτή μπορεί να φαίνεται υψηλή, είναι

σημαντικό να ερμηνευτεί προσεκτικά λόγω της ανισορροπίας των κατηγοριών στο dataset. Συγκεκριμένα, η πλειονότητα των πελατών δεν αθετεί τις πληρωμές τους, γεγονός που μπορεί να οδηγήσει το μοντέλο να προσαρμοστεί υπερβολικά σε αυτή την κυρίαρχη κατηγορία, παραβλέποντας την υποεκπροσωπούμενη κατηγορία των αθετήσεων.

Αναλύοντας τις μετρικές απόδοσης ανά κατηγορία:

- Για την κατηγορία 0 (πελάτες που δεν αθέτησαν τις πληρωμές τους), η ακρίβεια (precision) φτάνει το 84%. Αυτό σημαίνει ότι όταν το μοντέλο προβλέπει ότι ένας πελάτης δεν θα αθετήσει, έχει δίκιο στο 84% των περιπτώσεων. Παράλληλα, η ανάκληση (recall) για την κατηγορία αυτή είναι ιδιαίτερα υψηλή, στο 94%, υποδεικνύοντας ότι το μοντέλο καταφέρνει να εντοπίσει σχεδόν όλους τους πελάτες που πληρώνουν κανονικά. Το F1-score, ο αρμονικός μέσος της ακρίβειας και της ανάκλησης, ανέρχεται στο 0,89, αποδεικνύοντας την εξαιρετική συνολική απόδοση για αυτή την κατηγορία.
- Αντίθετα, για την κατηγορία 1 (πελάτες που αθέτησαν τις πληρωμές τους), η εικόνα είναι λιγότερο ενθαρρυντική. Η ακρίβεια για αυτή την κατηγορία είναι 64%, υποδεικνύοντας ότι μόνο τα δύο τρίτα των προβλέψεων αθέτησης είναι σωστές. Ακόμη πιο ανησυχητική είναι η ανάκληση, η οποία ανέρχεται μόλις στο 37%. Αυτό σημαίνει ότι το μοντέλο αποτυγχάνει να εντοπίσει την πλειονότητα των πελατών που πραγματικά αθετούν τις πληρωμές τους. Το F1-score για την κατηγορία 1 είναι 0,47, υποδηλώνοντας τη δυσκολία του μοντέλου να εξισορροπήσει την ακρίβεια και την ανάκληση για την υποεκπροσωπούμενη κατηγορία.

Η τιμή του ROC-AUC (Receiver Operating Characteristic - Area Under Curve) ανέρχεται σε 0,75. Η μέτρηση αυτή εκφράζει την ικανότητα του μοντέλου να διαχωρίζει τις δύο κατηγορίες και δείχνει ότι το baseline μοντέλο έχει σχετικά καλή, αλλά όχι εξαιρετική, διακριτική ικανότητα. Με άλλα λόγια, μπορεί να ξεχωρίσει σε ικανοποιητικό βαθμό τους πελάτες που αθετούν από αυτούς που δεν αθετούν, ωστόσο εξακολουθεί να απαιτείται βελτίωση ώστε να επιτευχθεί πιο αξιόπιστη πρόβλεψη.

Συνοψίζοντας, το baseline μοντέλο δείχνει καλή απόδοση στην πρόβλεψη των πελατών που πληρώνουν κανονικά, αλλά αποτυγχάνει να εντοπίσει επαρκώς τους πελάτες που αθετούν τις πληρωμές τους. Η υψηλή συνολική ακρίβεια είναι «παραπλανητική», καθώς δεν αντικατοπτρίζει τη χαμηλή αποτελεσματικότητα του μοντέλου στην ανίχνευση της κατηγορίας που παρουσιάζει το μεγαλύτερο επιχειρησιακό ενδιαφέρον. Αυτό το αποτέλεσμα υπογραμμίζει την ανάγκη για εφαρμογή τεχνικών

εξισορρόπησης των δεδομένων, όπως το SMOTE και το undersampling, για τη βελτίωση της απόδοσης στην πρόβλεψη αθετήσεων.

2) SMOTE

Η εφαρμογή της τεχνικής SMOTE στο dataset στοχεύει στην αντιμετώπιση του προβλήματος ανισορροπίας των κατηγοριών, δημιουργώντας συνθετικά δείγματα για την υποεκπροσωπούμενη κατηγορία (πελάτες που αθετούν).

Class	Precision	Recall	F1-Score	Support
0	0.86	0.90	0.88	4687
1	0.56	0.47	0.51	1313
Accuracy			0.80	6000
Macro Avg	0.71	0.68	0.69	6000
Weighted Avg	0.79	0.80	0.80	6000

Table 6: Αποτελέσματα ταξινόμησης του μοντέλου Random Forest, μετά την εφαρμογή της τεχνικής SMOTE, για το σύνολο δεδομένων Credit Card Clients.

Accuracy (Baseline): 0.8026666666666666

ROC-AUC (Baseline): 0.7532212626163242

Το αποτέλεσμα αυτής της διαδικασίας είναι η βελτίωση του recall για την κατηγορία 1 από 37% σε 47%, γεγονός που δείχνει ότι το μοντέλο καταφέρνει να εντοπίσει περισσότερους πελάτες που αθετούν. Ωστόσο, αυτή η βελτίωση συνοδεύεται από μια μείωση στην ακρίβεια της ίδιας κατηγορίας, η οποία από 64% μειώνεται σε 56%. Το F1-score αυξάνεται από 0,47 σε 0,51, υποδηλώνοντας μια πιο ισορροπημένη απόδοση μεταξύ ακρίβειας και ανάκλησης.

Η συνολική ακρίβεια του μοντέλου μειώνεται ελαφρώς στο 80,27%, ενώ η τιμή του ROC-AUC ανέρχεται σε 0,7532, δηλαδή παραμένει σχεδόν στα ίδια επίπεδα με το baseline μοντέλο (0,7548). Αυτό δείχνει ότι, αν και το SMOTE συμβάλλει στην καλύτερη ανίχνευση της μειονοτικής κατηγορίας, δεν βελτιώνει ουσιαστικά τη συνολική διακριτική ικανότητα του μοντέλου. Παρόλα αυτά, η ενίσχυση της ικανότητας εντοπισμού των αθετήσεων καθιστά τη χρήση της τεχνικής σημαντική σε περιπτώσεις όπου προέχει η κάλυψη της υποεκπροσωπούμενης κατηγορίας.

3) Undersampling

Η μέθοδος του undersampling εφαρμόζεται για να μειώσει τον αριθμό των δειγμάτων της κυρίαρχης κατηγορίας (πελάτες που δεν αθετούν), ώστε να επιτευχθεί ισορροπία με την υποεκπροσωπούμενη κατηγορία.

Class	Precision	Recall	F1-Score	Support
0	0.88	0.77	0.82	4687
1	0.43	0.63	0.51	1313
Accuracy			0.74	6000
Macro Avg	0.66	0.70	0.67	6000
Weighted Avg	0.78	0.74	0.75	6000

Table 7: Αποτελέσματα ταξινόμησης του μοντέλου Random Forest, μετά την εφαρμογή της μεθόδου Undersampling, για το σύνολο δεδομένων Credit Card Clients.

Accuracy (Baseline): 0.7378333333333333

ROC-AUC (Baseline): 0.7615632745431408

Το αποτέλεσμα αυτής της προσέγγισης είναι η σημαντική βελτίωση της ανάκλησης για την κατηγορία 1, η οποία φτάνει το 63%, η υψηλότερη από όλες τις μεθόδους. Αυτό σημαίνει ότι το μοντέλο καταφέρνει να εντοπίσει την πλειονότητα των πελατών που αθετούν, καθιστώντας το ιδιαίτερα χρήσιμο σε εφαρμογές όπου η ανίχνευση των αθετήσεων είναι κρίσιμη.

Ωστόσο, αυτή η βελτίωση στο recall συνοδεύεται από μείωση στην ακρίβεια της κατηγορίας 1, η οποία φτάνει μόλις το 43%. Το F1-score παραμένει σταθερό στο 0,51, παρόμοιο με αυτό του μοντέλου με SMOTE, γεγονός που υποδηλώνει παρόμοια ισορροπία μεταξύ ακρίβειας και ανάκλης. Η συνολική ακρίβεια του μοντέλου μειώνεται στο 73,78%, η χαμηλότερη από όλες τις μεθόδους, κάτι αναμενόμενο λόγω της μείωσης των δειγμάτων της κατηγορίας 0.

Παρά τη μείωση στην ακρίβεια, το ROC-AUC ανέρχεται σε 0,7616, δηλαδή ελαφρώς υψηλότερα σε σχέση με το baseline (0,7548) και το SMOTE (0,7532). Αυτό δείχνει ότι το μοντέλο με undersampling επιτυγχάνει τη σχετικά καλύτερη διακριτική ικανότητα ανάμεσα στις μεθόδους που εφαρμόστηκαν, χωρίς όμως η διαφορά να είναι ιδιαίτερα μεγάλη. Συνολικά, η μέθοδος του undersampling προσφέρει ουσιαστική

βελτίωση στην ανίχνευση των αθετήσεων, αλλά με κόστος την αύξηση των ψευδώς θετικών και τη μείωση της συνολικής ακρίβειας.

B. Μετά την αρχική αξιολόγηση του baseline μοντέλου και την εφαρμογή τεχνικών εξισορρόπησης, όπως το SMOTE και το undersampling, η ανάλυση επεκτάθηκε με την εισαγωγή επιπλέον ταξινομητών και τη βελτιστοποίηση παραμέτρων. Η αναβαθμισμένη αυτή μεθοδολογική προσέγγιση στοχεύει στη διερεύνηση της δυνατότητας περαιτέρω βελτίωσης της απόδοσης των μοντέλων όσον αφορά τη διάκριση μεταξύ πελατών που αθετούν τις πληρωμές τους και εκείνων που είναι συνεπείς.

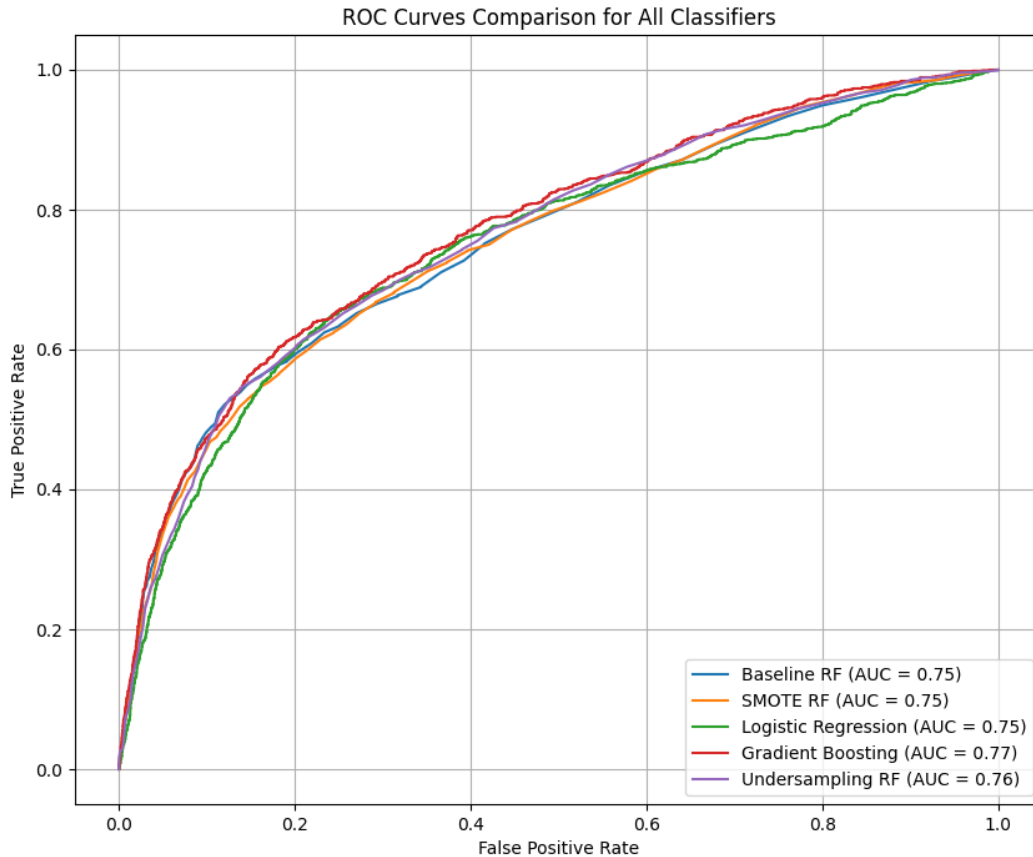
Η ανάλυση επικεντρώθηκε στη χρήση του Random Forest Classifier ως σημείο αναφοράς, ενώ εξετάστηκαν και άλλοι αλγόριθμοι όπως η Λογιστική Παλινδρόμηση (Logistic Regression) και η Ενίσχυση Κλίσης (Gradient Boosting). Τα δεδομένα εξισορροπήθηκαν μέσω SMOTE πριν την εκπαίδευση των μοντέλων, προκειμένου να αντιμετωπιστεί η έντονη ανισορροπία των κατηγοριών. Επιπλέον, πραγματοποιήθηκε βελτιστοποίηση παραμέτρων (Grid Search) στον Random Forest, προκειμένου να ενισχυθεί η ακρίβεια και η γενίκευση του μοντέλου.

Ο παρακάτω πίνακας συνοψίζει τις τιμές του ROC-AUC για κάθε μέθοδο και μοντέλο που εφαρμόστηκε, παρέχοντας μια συγκριτική εικόνα της διακριτικής ικανότητας των διαφορετικών αλγορίθμων και τεχνικών εξισορρόπησης.

Model	ROC-AUC
Baseline Random Forest	0.75
SMOTE Random Forest	0.75
Undersampling Random Forest	0.76
Logistic Regression (SMOTE)	0.75
Gradient Boosting (SMOTE)	0.77
Best Random Forest Parameters	0.76

Table 8: Τιμές ROC-AUC για τα εξεταζόμενα μοντέλα (Random Forest, Logistic Regression, Gradient Boosting) με και χωρίς εφαρμογή τεχνικών εξισορρόπησης δεδομένων για το σύνολο δεδομένων Credit Card Clients.

Για τη σαφέστερη κατανόηση της απόδοσης των μοντέλων, δημιουργήθηκε το παρακάτω γράφημα ROC Curves Comparison for All Classifiers, το οποίο απεικονίζει τις καμπύλες χαρακτηριστικών λειτουργίας του δέκτη (ROC) για κάθε μία από τις εξεταζόμενες μεθόδους.



Εικόνα 32: Σύγκριση ROC καμπυλών για διαφορετικούς ταξινομητές και στρατηγικές διαχείρισης ανισορροπίας: Random Forest (Baseline και SMOTE), Logistic Regression, Gradient Boosting και Undersampling RF.

Το baseline μοντέλο Random Forest πέτυχε ROC-AUC 0.75, σημειώνοντας σαφή βελτίωση σε σύγκριση με το αρχικό μοντέλο της πρώτης ανάλυσης (ROC-AUC 0.6551). Αυτή η αύξηση της απόδοσης οφείλεται κυρίως στον τροποποιημένο τρόπο υπολογισμού του ROC-AUC μέσω των πιθανοτήτων πρόβλεψης αντί για άμεσες ταξινομήσεις. Η αξιολόγηση με προβλεπόμενες πιθανότητες παρέχει μια πιο ευαίσθητη και ακριβή μέτρηση της διακριτικής ικανότητας του μοντέλου, αναδεικνύοντας καλύτερα την απόδοσή του στη διάκριση μεταξύ πελατών που αθετούν και εκείνων που πληρώνουν κανονικά.

Η εφαρμογή του SMOTE, οδήγησε σε ROC-AUC 0.75, ισοδύναμο με το baseline μοντέλο. Παρότι η βελτίωση του ROC-AUC δεν ήταν αξιοσημείωτη, υπήρξε αύξηση του recall για την κατηγορία των αθετήσεων, γεγονός που υποδεικνύει βελτίωση στην ανίχνευση πελατών υψηλού κινδύνου. Ωστόσο, αυτή η βελτίωση συνοδεύτηκε από μείωση της ακρίβειας, καθώς το μοντέλο παράγαγε περισσότερα ψευδώς θετικά αποτελέσματα.

Αντίστοιχα, η μέθοδος του undersampling, επέφερε μικρή αύξηση στο ROC-AUC (0.76). Αυτή η προσέγγιση βελτίωσε σημαντικά την ανάκληση για την κατηγορία 1, εντοπίζοντας περισσότερους πελάτες που αθετούν, αλλά με κόστος τη μείωση της συνολικής ακρίβειας λόγω της απώλειας πληροφοριών από την κυρίαρχη κατηγορία.

Η Λογιστική Παλινδρόμηση (Logistic Regression) με εφαρμογή SMOTE πέτυχε ROC-AUC 0.75, επιδεικνύοντας αντίστοιχη απόδοση με το baseline Random Forest. Αυτό το αποτέλεσμα επιβεβαιώνει ότι, αν και η λογιστική παλινδρόμηση αποτελεί έναν απλό γραμμικό αλγόριθμο, μπορεί να προσφέρει ανταγωνιστική απόδοση όταν τα δεδομένα είναι καλά εξισορροπημένα και προεπεξεργασμένα.

Από την άλλη πλευρά, ο αλγόριθμος Gradient Boosting (SMOTE) παρουσίασε την καλύτερη απόδοση με ROC-AUC 0.77. Η βελτίωση αυτή καταδεικνύει τη δυναμική των μοντέλων ενίσχυσης κλίσης στην αντιμετώπιση πολύπλοκων προβλημάτων ταξινόμησης με ανισόρροπα δεδομένα. Η ικανότητα του Gradient Boosting να συνδυάζει πολλαπλά ασθενή μοντέλα για τη δημιουργία ενός ισχυρού ταξινομητή αποδεικνύεται κρίσιμη για την ακριβή ανίχνευση πελατών υψηλού κινδύνου.

Η βελτιστοποίηση παραμέτρων του Random Forest με χρήση Grid Search οδήγησε σε περαιτέρω αύξηση της απόδοσης του μοντέλου, με το ROC-AUC να ανέρχεται στο 0.76. Οι βέλτιστες παράμετροι που προέκυψαν ήταν οι εξής: `max_depth = 20`, `max_features = 'sqrt'`, και `n_estimators = 200`. Η επιλογή μεγαλύτερου βάθους δέντρων και περισσότερων εκτιμητών βελτίωσε την ικανότητα του μοντέλου να συλλαμβάνει πιο πολύπλοκα μοτίβα στα δεδομένα, διατηρώντας παράλληλα την ισορροπία ανάμεσα στην ακρίβεια και την γενίκευση.

Συμπερασματικά, η εφαρμογή τεχνικών εξισορρόπησης όπως το SMOTE και το undersampling αποδείχθηκε καθοριστική για τη βελτίωση της ικανότητας των μοντέλων να εντοπίζουν αθετήσεις πληρωμών. Αν και το SMOTE προσέφερε σημαντική βελτίωση στην ανίχνευση πελατών υψηλού κινδύνου, η μέθοδος του undersampling παρουσίασε την υψηλότερη διακριτική ικανότητα με ROC-AUC 0.76, αποδεικνύοντας ότι η στρατηγική μείωσης της κυρίαρχης κατηγορίας μπορεί να είναι εξαιρετικά αποτελεσματική σε ορισμένα σενάρια. Ο αλγόριθμος Gradient Boosting αναδείχθηκε ως ο πιο

αποδοτικός ταξινομητής, με ROC-AUC 0.77, γεγονός που υπογραμμίζει τη δύναμή του στην αντιμετώπιση πολύπλοκων προβλημάτων ταξινόμησης με ανισόρροπα δεδομένα. Η βελτιστοποίηση παραμέτρων του Random Forest οδήγησε σε επιδόσεις που προσεγγίζουν αυτές του Gradient Boosting, επιβεβαιώνοντας τη σημασία της κατάλληλης προσαρμογής υπερπαραμέτρων. Τα αποτελέσματα της ανάλυσης υπογραμμίζουν τη σημασία της διερεύνησης διαφορετικών τεχνικών εξισορρόπησης και της προσαρμογής παραμέτρων για τη βελτίωση της απόδοσης των μοντέλων μηχανικής μάθησης. Η επιλογή της κατάλληλης στρατηγικής εξαρτάται από τη φύση του προβλήματος και τις επιχειρησιακές απαιτήσεις, ειδικά όταν πρόκειται για την πρόβλεψη αθετήσεων πληρωμών, όπου η ανίχνευση πελατών υψηλού κινδύνου αποτελεί κρίσιμο παράγοντα επιτυχίας.

7. Συμπεράσματα

Η παρούσα διπλωματική εργασία αναλύει την εφαρμογή τεχνικών εξόρυξης γνώσης και μηχανικής μάθησης σε οικονομικά δεδομένα, με επίκεντρο τον τραπεζικό τομέα και τη διαχείριση πιστωτικού κινδύνου καθώς και τις στρατηγικές μάρκετινγκ. Μέσα από την επεξεργασία δύο πραγματικών συνόλων δεδομένων (Bank Marketing και Default of Credit Card Clients), εφαρμόζονται και συγκρίνονται διαφορετικές τεχνικές ανάλυσης, προεπεξεργασίας και μοντελοποίησης. Το παρόν κεφάλαιο συνοψίζει τα βασικά ευρήματα, αποτιμά τη συνεισφορά της μελέτης, αναδεικνύει τους περιορισμούς και προτείνει κατευθύνσεις για μελλοντική έρευνα.

Αρχικά, αναδεικνύεται η σημασία των τεχνικών προεπεξεργασίας δεδομένων για την επιτυχή εφαρμογή αλγορίθμων μηχανικής μάθησης στα οικονομικά σύνολα δεδομένων. Η μετατροπή των κατηγορικών μεταβλητών σε αριθμητική μορφή μέσω του One-Hot Encoding και η κανονικοποίηση των αριθμητικών χαρακτηριστικών μέσω StandardScaler αποτελούν βασικά στάδια για τη διασφάλιση της ομοιομορφίας και της σωστής λειτουργίας των μοντέλων. Ιδιαίτερη έμφαση δίνεται στη διαχείριση της ανισορροπίας των δεδομένων, πρόβλημα σύνηθες στις εφαρμογές πρόβλεψης σπάνιων γεγονότων, όπως η αθέτηση πληρωμών ή η ανταπόκριση σε καμπάνιες μάρκετινγκ. Οι τεχνικές SMOTE και undersampling βελτιώνουν την απόδοση των ταξινομητών, επιτρέποντας ακριβέστερη εκμάθηση της μειοψηφικής κλάσης.

Η εφαρμογή διαφορετικών αλγορίθμων μηχανικής μάθησης οδηγεί σε χρήσιμα συμπεράσματα σχετικά με την καταλληλότητα κάθε μεθόδου στα υπό εξέταση προβλήματα. Ο αλγόριθμος Random Forest παρουσιάζει υψηλή ακρίβεια και στα δύο σύνολα δεδομένων. Το πλεονέκτημά του εντοπίζεται στην ικανότητά του να χειρίζεται δεδομένα με πολλαπλές μεταβλητές και να εντοπίζει περίπλοκες συσχετίσεις, ενώ παρέχει και ερμηνεύσιμες πληροφορίες μέσω των δεικτών «feature importance». Παράλληλα, το

Gradient Boosting εμφανίζεται ιδιαίτερα αποδοτικό, με μεγαλύτερη ευαισθησία στη διαχείριση πολύπλοκων αλληλεπιδράσεων και ασθενών συσχετίσεων, επιτυγχάνοντας συχνά τις υψηλότερες τιμές ROC-AUC.

Η Logistic Regression, αν και απλούστερη μέθοδος, αποδίδει ικανοποιητικά αποτελέσματα, ιδίως μετά την εφαρμογή τεχνικών προεπεξεργασίας και εξισορρόπησης των κλάσεων. Η ευκολία στην ερμηνεία των παραμέτρων της καθιστά τη Logistic Regression ιδιαίτερα χρήσιμη σε περιβάλλοντα όπου η κατανόηση της λογικής του μοντέλου αποτελεί προτεραιότητα, όπως συμβαίνει συχνά στον τραπεζικό κλάδο λόγω των κανονιστικών απαιτήσεων διαφάνειας.

Σε επίπεδο πρακτικής συνεισφοράς, η εργασία αποδεικνύει τη δυναμική της μηχανικής μάθησης στην αποτελεσματικότερη στόχευση πελατών και στη βελτίωση των πιστωτικών διαδικασιών. Στην περίπτωση του Bank Marketing Dataset, η χρήση ταξινομητών υψηλής ακρίβειας επιτρέπει στις τράπεζες να επικεντρώνουν τις καμπάνιες μάρκετινγκ σε πελάτες με υψηλή πιθανότητα αποδοχής των προσφορών, μειώνοντας το κόστος άσκοπων προσπαθειών και βελτιώνοντας την αποδοτικότητα των ενεργειών. Αντίστοιχα, στο Default of Credit Card Clients Dataset, η επιτυχής πρόβλεψη αθετήσεων πληρωμών συνιστά εργαλείο υψηλής στρατηγικής σημασίας για τη διαχείριση πιστωτικού κινδύνου, τον περιορισμό επισφαλειών και την έγκαιρη λήψη προληπτικών μέτρων.

Τα αποτελέσματα αναδεικνύουν τον ουσιαστικό ρόλο των τεχνικών διαχείρισης ανισορροπίας δεδομένων στη βελτίωση της απόδοσης των μοντέλων. Ιδιαίτερα το SMOTE, δημιουργώντας συνθετικά δείγματα για τη μειοψηφική κλάση, συμβάλλει στην εξισορρόπηση των ταξινομητών χωρίς να διαστρεβλώνει την πληροφορία. Αυτό οδηγεί σε αύξηση των δεικτών AUC και F1-score, οι οποίοι αποτυπώνουν με μεγαλύτερη ακρίβεια την προγνωστική ικανότητα των μοντέλων, ειδικά σε προβλήματα όπου τα ψευδώς θετικά ή αρνητικά αποτελέσματα επιφέρουν σημαντικές επιχειρησιακές συνέπειες.

Παρά τα πλεονεκτήματα, η εργασία αναγνωρίζει και συγκεκριμένους περιορισμούς. Η εκτεταμένη χρήση τεχνικών oversampling ενδέχεται να οδηγήσει σε υπερπροσαρμογή, ειδικά όταν δημιουργούνται συνθετικά δείγματα που δεν απεικονίζουν πλήρως την πραγματικότητα. Επιπλέον, η εκπαίδευση μοντέλων όπως το Gradient Boosting απαιτεί αυξημένο υπολογιστικό κόστος, γεγονός που καθιστά την υλοποίηση απαιτητική για πολύ μεγάλα σύνολα δεδομένων σε πραγματικές επιχειρησιακές υποδομές.

Επιπρόσθετα, η εργασία βασίζεται σε δημόσια διαθέσιμα σύνολα δεδομένων τα οποία, αν και επαρκή για ερευνητικούς σκοπούς, δεν αποτυπώνουν πλήρως την πολυπλοκότητα των πραγματικών τραπεζικών δεδομένων. Στην πράξη, τα δεδομένα που διαθέτουν οι τράπεζες περιλαμβάνουν μεγαλύτερη χρονική

δυναμική, σύνθετες συμπεριφορικές μεταβλητές και εξωτερικούς οικονομικούς δείκτες, οι οποίοι θα μπορούσαν να εμπλουτίσουν περαιτέρω την ανάλυση.

Η εργασία επιβεβαιώνει επίσης τη σημασία της οπτικοποίησης δεδομένων ως εργαλείου κατανόησης και ερμηνείας των αποτελεσμάτων. Μέσω των heatmaps, των καμπυλών ROC και των classification reports, καθίσταται δυνατή η βαθύτερη κατανόηση της συμπεριφοράς των μοντέλων και η ουσιαστικότερη επικοινωνία των ευρημάτων μεταξύ τεχνικών και διοικητικών στελεχών.

Με βάση τα ευρήματα προκύπτουν και οι εξής προτάσεις για περαιτέρω έρευνα:

- Εφαρμογή των τεχνικών αυτών σε πραγματικά τραπεζικά δεδομένα που εμφανίζουν μεγαλύτερη πολυπλοκότητα.
- Διερεύνηση σύγχρονων αλγορίθμων βαθιάς μάθησης (Deep Learning), όπως τα Recurrent Neural Networks (RNN) και τα Graph Neural Networks (GNN), τα οποία αξιοποιούν χρονοσειρές και σχέσεις μεταξύ συναλλαγών.
- Ανάπτυξη υβριδικών μοντέλων που συνδυάζουν αλγοριθμικά αποτελέσματα με κανόνες επιχειρηματικής λογικής, ενισχύοντας τη διαφάνεια και τη συμμόρφωση.
- Ενσωμάτωση δεικτών fairness και ελέγχου bias, ώστε να αποφεύγεται η αλγοριθμική μεροληψία σε ευαίσθητες εφαρμογές.

Συνολικά, η εργασία αναδεικνύει τις δυνατότητες και τις προκλήσεις της εφαρμογής της εξόρυξης γνώσης και της μηχανικής μάθησης στον τραπεζικό τομέα, παρέχοντας πολύτιμες κατευθύνσεις για βελτίωση της λήψης αποφάσεων, ενίσχυση της αποδοτικότητας και καλύτερη εξυπηρέτηση των πελατών.

Βιβλιογραφία

- [1] J. Han, M. Kamber, and J. Pei, *Data Mining: Concepts and Techniques*, 3rd ed., Morgan Kaufmann, 2011.
- [2] I. H. Witten, E. Frank, M. A. Hall, and C. J. Pal, *Data Mining: Practical Machine Learning Tools and Techniques*, 4th ed., Morgan Kaufmann, 2016.
- [3] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5-32, 2001.
- [4] J. Friedman, "Greedy function approximation: a gradient boosting machine," *Annals of Statistics*, vol. 29, no. 5, pp. 1189–1232, 2001.
- [5] D. W. Hosmer, S. Lemeshow, and R. X. Sturdivant, *Applied Logistic Regression*, 3rd ed., Wiley, 2013.
- [6] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "SMOTE: Synthetic Minority Over-sampling Technique," *Journal of Artificial Intelligence Research*, vol. 16, pp. 321–357, 2002.
- [7] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*, MIT Press, 2016.
- [8] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [9] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, 2016, pp. 785–794.
- [10] F. Doshi-Velez and B. Kim, "Towards a rigorous science of interpretable machine learning," *arXiv preprint arXiv:1702.08608*, 2017.
- [11] M. Hardt, E. Price, and N. Srebro, "Equality of opportunity in supervised learning," in *Proc. 30th Advances in Neural Information Processing Systems (NeurIPS)*, 2016.
- [12] European Union, "General Data Protection Regulation (GDPR)," *Regulation (EU) 2016/679*, *Official Journal of the European Union*, 2016.
- [13] S. Barocas, M. Hardt, and A. Narayanan, *Fairness and Machine Learning*, fairmlbook.org, 2021.

[14] A. T. Nguyen, J. W. Tucker, and J. N. Zhang, "Machine Learning in Credit Risk Modeling," SSRN Electronic Journal, 2021.

[15] UCI Machine Learning Repository, "Bank Marketing Dataset," [Online]. Available: <https://archive.ics.uci.edu/dataset/222/bank+marketing>.

[16] UCI Machine Learning Repository, "Default of Credit Card Clients Dataset," [Online]. Available: <https://archive.ics.uci.edu/dataset/350/default+of+credit+card+clients>.

[17] Javelin Strategy & Research, *2020 Identity Fraud Study: Genesis of the Identity Fraud Crisis*, Pleasanton, CA, 2020. [Online]. Available: <https://www.javelinstrategy.com>