



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

Διπλωματική Εργασία

Μαθηματικές Μέθοδοι Ανίχνευσης Συμβάντων και
Εντοπισμού μη Ομαλών Χαρακτηριστικών σε Χρονολογικές
Σειρές

ΓΕΩΡΓΙΟΣ ΜΑΡΔΕΛΗΣ

$$\frac{\partial}{\partial a} \ln f_{a, \sigma^2}(\xi_1) = \frac{(\xi_1 - a)}{\sigma^2} f_{a, \sigma^2}(\xi_1) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left\{-\frac{(\xi_1 - a)^2}{2\sigma^2}\right\}$$
$$\int_{R_n} T(x) \cdot \frac{\partial}{\partial \theta} f(x, \theta) dx = M\left(T(\xi) \cdot \frac{\partial}{\partial \theta} \ln l(\xi, \theta)\right)$$

Υπεβλήθη για την εκπλήρωση μέρους των απαιτήσεων για την απόκτηση του



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

Διπλώματος Πολιτικού Μηχανικού

ΒΟΛΟΣ 2023

© 2023 Γεώργιος Μαρδέλης

Η έγκριση της διπλωματικής εργασίας από το Τμήμα Πολιτικών Μηχανικών της Πολυτεχνικής Σχολής του Πανεπιστημίου Θεσσαλίας δεν υποδηλώνει αποδοχή των απόψεων του συγγραφέα (Ν. 5343/32 αρ. 202 παρ. 2).

Εγκρίθηκε από τα Μέλη της Τριμελούς Εξεταστικής Επιτροπής:

Πρώτος Εξεταστής [Δρ. Αθανάσιος Φράγκου](#)
(Επιβλέπων) [Τμήμα Πολιτικών Μηχανικών, Πανεπιστήμιο Θεσσαλίας](#)

Δεύτερος Εξεταστής [Αθανάσιος Θεοφιλάτος](#)
[Επίκουρος Καθηγητής, Τμήμα Πολιτικών Μηχανικών](#)

Τρίτος Εξεταστής [Δρ. Μάριος Σπηλιωτόπουλος](#)
[Τμήμα Πολιτικών Μηχανικών, Πανεπιστήμιο Θεσσαλίας](#)



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

“Ευχαριστίες”

Πρώτα απ’ όλα, θέλω να ευχαριστήσω τον επιβλέποντα της διπλωματικής εργασίας μου, Δρ. κ. Αθανάσιο Φράγκου για την πολύτιμη βοήθεια και καθοδήγησή του κατά τη διάρκεια της εργασίας μου ειδικότερα κάτω υπό τις συγκεκριμένες και αντίξοες συνθήκες της περιόδου εκπόνησης της διπλωματικής, εκτίοντας ταυτοχρόνως την στρατιωτική μου θητεία στην παραμεθόριο περιοχή της χώρας . Είμαι ευγνώμων στον κ. Φράγκου που από την πρώτη στιγμή στάθηκε αρωγός σε όλη αυτή την ύψιστη προσωπική μου προσπάθεια να επέλθει ένα άριστο αποτέλεσμα επικοινωνώντας μαζί μου ανά πάσα στιγμή χωρίς καμία κωλυσιεργία και πάντοτε με προθυμία για την εύρεση λύσης σε οποιοδήποτε πρόβλημα . Θα ήθελα να εκφράσω επίσης, τις ευχαριστίες μου, στα υπόλοιπα μέλη της εξεταστικής επιτροπής της διπλωματικής εργασίας μου, Καθηγητές κκ. Μάριο Σπηλιωτόπουλο και Αθανάσιο Θεοφιλάτο για την προσεκτική ανάγνωση της εργασίας μου και για τις πολύτιμες υποδείξεις τους. Ευχαριστώ όλους τους φίλους(ες) μου για την ηθική και ψυχολογική υποστήριξή τους. Πάνω απ’ όλα, είμαι ευγνώμων στους γονείς μου, Δημήτριο και Ελένη Μαρδέλη για την ολόψυχη αγάπη και υποστήριξή τους όλα αυτά τα χρόνια. Τέλος, σας ευχαριστώ όλους για την κατανόησή και την στήριξη σας, ιδιαίτερα κατά τη διάρκεια των τελευταίων δύσκολων μηνών της προσπάθειάς μου. Άλλωστε η ζωή αποκτά μεγαλύτερο νόημα όταν υπάρχουν προκλήσεις...

Βόλος, 24/06/2023

Γεώργιος Δ. Μαρδέλης



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

Περίληψη

Η παρούσα εργασία είχε ως κεντρικό αντικείμενο πραγμάτευσης την ανίχνευση ανωμαλιών σε δεδομένα χρονοσειρών. Για την μελέτη του εν λόγω ερευνητικού αντικειμένου επιλέχθηκε η μέθοδος της βιβλιογραφικής ανασκόπησης ως καταλληλότερη, καθώς προσφέρει την δυνατότητα της σφαιρικής προσέγγισης του. Σε πρώτη φάση μελετήθηκαν τα δυναμικά συστήματα και οι βασικές ιδιότητες των χρονοσειρών ως προαπαιτούμενες βάσεις για την πληρέστερη κατανόηση των ανωμαλιών σε χρονοσειρές. Στην συνέχεια έγινε μια εισαγωγή στην έννοια των ανωμαλιών και στην συνήθη τυπολογία τους, η οποία πλαισιώθηκε από τους τρόπους και τις βασικές τεχνικές ανίχνευσής τους, όπως τα Νευρωνικά δίκτυα, τα Μπαγесиανά δίκτυα κ.ά. Πέραν της θεωρητικής πραγμάτευσης του αντικειμένου, στην εργασία περιλαμβάνεται και η μελέτη δύο ενδεικτικών περιπτώσεων-εφαρμογών τεχνικών ανίχνευσης ανωμαλιών με την χρήση μεθόδων βαθιάς μάθησης (Deep Learning). Η πρώτη περίπτωση είναι μια εφαρμοσμένη μελέτη με την μέθοδο των Συνελικτικών Νευρωνικών Δικτύων (CNN) σε δεδομένα χρονοσειρών που προκύπτουν από κυκλοφοριακές ροές. Η δεύτερη περίπτωση αφορά τον τομέα της πολιτικής μηχανικής και συγκεκριμένα τον εντοπισμό ανωμαλιών σε χρονοσειρές δομικών δεδομένων από αισθητήρες σε γέφυρα. Στην δεύτερη περίπτωση χρησιμοποιείται πάλι η μέθοδος CNN, με αυξημένο βαθμό πολυπλοκότητας και με μια διαφορετική προσέγγιση στην προεπεξεργασία των δεδομένων. Το συμπέρασμα της παρούσας εργασίας είναι πως οι τεχνικές βαθιάς μάθησης προσφέρουν ισχυρές δυνατότητες εντοπισμού και κατηγοριοποίησης των ανωμαλιών σε χρονοσειρές.

Περιεχόμενα

I. Κατάλογος Διαγραμμάτων	7
II. Κατάλογος σχημάτων.....	Σφάλμα! Δεν έχει οριστεί σελιδοδείκτης.
III. Κατάλογος πινάκων	8
Κεφάλαιο 1 ^ο : Εισαγωγή.....	8



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

Κεφάλαιο 2 ^ο : Δυναμικά συστήματα	10
2.1 Ορισμός.....	10
2.2 Παραδείγματα δυναμικών συστημάτων	12
2.3 Ροή	14
2.4 Εξέλιξη.....	16
2.5 Διακριτά δυναμικά συστήματα	16
2.5.1 Διαγράμματα και ροές	16
2.5.2 Τροχιές.....	17
2.5.3 Διαγράμματα φάσεων	17
Κεφάλαιο 3 ^ο : Χρονοσειρές	18
3.1 Ορισμός και παραδείγματα χρονοσειρών	18
3.1.1 Παραδείγματα.....	19
3.2 Χαρακτηριστικά χρονοσειρών.....	23
3.2.1 Συνιστώσες των χρονοσειρών.....	24
3.2.2 Διαγράμματα χρονοσειρών.....	26
3.3 Μοντέλα χρονοσειρών.....	27
3.3.1 Προσθετικό μοντέλο.....	27
3.3.2 Πολλαπλασιαστικό μοντέλο	28
3.4 Στασιμότητα.....	29
3.4.1 Τρόποι ελέγχου της στασιμότητας.....	30
3.4.2 Τρόποι απαλοιφής τάσης	31
Κεφάλαιο 4 ^ο : Ανίχνευση ανωμαλιών.....	34
4.1 Προβλήματα ανίχνευσης.....	37
4.2 Φύση των δεδομένων εισόδου	38
4.3 Τύπος ανωμαλίας	39
4.3.1 Σημειακές ανωμαλίες.....	39



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

4.3.2	Ανωμαλίες πλαισίου	40
4.3.3	Συλλογικές ανωμαλίες	42
Κεφάλαιο 5ο : Μέθοδοι ανίχνευσης ανωμαλιών		43
5.1	Τεχνικές ανίχνευσης ανωμαλιών με βάση την ταξινόμηση.....	43
5.1.1	Νευρωνικά δίκτυα.....	44
5.1.2	Δίκτυα Bayesian.....	47
5.1.3	Μηχανές διανυσμάτων υποστήριξης	48
5.1.4	Βασισμένες σε κανόνες.....	49
5.1.5	Υπολογιστική πολυπλοκότητα.....	50
5.2	Τεχνικές ανίχνευσης ανωμαλιών με βάση τον πλησιέστερο γείτονα	51
5.2.1	Χρήση της απόστασης από τον <i>k</i> th πλησιέστερο γείτονα	52
5.2.2	Χρήση της σχετικής πυκνότητας	53
5.2.3	Υπολογιστική πολυπλοκότητα.....	55
5.3	Τεχνικές ανίχνευσης ανωμαλιών με βάση την ομαδοποίηση	56
5.3.1	Υπολογιστική πολυπλοκότητα.....	58
5.4	Στατιστικές τεχνικές ανίχνευσης ανωμαλιών.....	59
5.4.1	Παραμετρικές τεχνικές	60
5.4.2	Μη παραμετρικές τεχνικές.....	62
5.4.3	Υπολογιστική πολυπλοκότητα.....	63
5.5	Χειρισμός ανωμαλιών του πλαισίου	64
Κεφάλαιο 6ο : Εφαρμογές ανίχνευσης ανωμαλιών		67
6.1	Ανίχνευση ανωμαλιών στην κυκλοφοριακή ροή.....	67
6.1.1	Δεδομένα και μεθοδολογία	68
6.1.2	Αποτελέσματα μοντέλων	72
6.2	Ανίχνευση ανωμαλιών σε κατασκευές γεφυρών με χρήση συνελκτικού νευρωνικού δικτύου (CNN).....	77



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

6.2.1 Επεξεργασία των δεδομένων και μεθοδολογία.....	78
6.2.2 Αποτελέσματα του μοντέλου CNN	85
Κεφάλαιο 7 ^ο : Συμπεράσματα.....	89
Βιβλιογραφία	92

I. Κατάλογος Διαγραμμάτων

Διάγραμμα 1: Γραφική αναπαράσταση του μοντέλου ανάπτυξης του πληθυσμού.....	14
Διάγραμμα 2: Εξαγωγή της τρέχουσας τάσης στο μέλλον για την πρόβλεψη των πωλήσεων.....	20
Διάγραμμα 3: Ακολουθία κλεισίματος του δείκτη μετοχών S&P	21
Διάγραμμα 4: Σύγκριση μεταξύ των πραγματικών τιμών και των τιμών πρόβλεψης στη θερμοκρασία στο Nanjing.....	22
Διάγραμμα 5: Παραδείγματα χρονοσειρών RR για τις κατηγορίες ασθενών.....	23
Διάγραμμα 6: Διάγραμμα χρονοσειράς με διαφορετικές συνιστώσες.....	27
Διάγραμμα 7: Προσθετικό μοντέλο με $M(t) + S(t) + E(t)$	28
Διάγραμμα 8: Παράδειγμα για το πολλαπλασιαστικό μοντέλο.....	29
Διάγραμμα 9 : Διάγραμμα αυτοσυσχέτισης (Acf).....	31
Διάγραμμα 10: Διάγραμμα της τιμής του χρυσού μετά την εξάλειψη της τάσης με τη μέθοδο της διαφοράς.....	32
Διάγραμμα 11: Διάγραμμα της τιμής του χρυσού μετά την εξάλειψη της τάσης με τη μέθοδο των ελαχίστων τετραγώνων.....	33
Διάγραμμα 12: Διάγραμμα των επιβατών με τη μέθοδο των κινητών μέσων όρων ..	34
Διάγραμμα 13: Ένα απλό παράδειγμα ανωμαλιών σε ένα δισδιάστατο σύνολο δεδομένων	36
Διάγραμμα 14. Συναφής ανωμαλία ανωμαλία t_2 σε χρονοσειρά θερμοκρασίας. Η θερμοκρασία τη στιγμή t_1 είναι η ίδια με εκείνη τη στιγμή t_2 , αλλά εμφανίζεται σε διαφορετικό πλαίσιο και ως εκ τούτου δεν θεωρείται ως ανωμαλία.	41
Διάγραμμα 15. Συλλογική ανωμαλία που αντιστοιχεί σε μια πρόωρη κολπική συστολή στην έξοδο ενός ανθρώπινου ηλεκτροκαρδιογραφήματος	42
Διάγραμμα 16: Χρήση ταξινόμησης για ανίχνευση ανωμαλιών	44
Διάγραμμα 17: Ιστόγραμμα της κατανομής της απώλειας ανομοιότητας στα δεδομένα εκπαίδευσης	73
Διάγραμμα 18: Συγκεκριμένα δεδομένα που ανιχνεύονται ως ανώμαλα με βάση το κατώτατο όριο (κατώφλι).....	73
Διάγραμμα 19 : Συγκεκριμένα δεδομένα που ανιχνεύονται ως ανώμαλα βάσει κατωφλίου με την μέθοδο CNN.....	74



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

Διάγραμμα 20 : Συγκεκριμένα δεδομένα που ανιχνεύονται ως ανώμαλα βάσει κατώφλιου με την μέθοδο BiLSTM.....	75
Διάγραμμα 21 : % της ανίχνευσης ανωμαλιών ανά τεταρτημόριο ντιλο που χρησιμοποιείται ως κατώφλι με την μέθοδο CNN (αριστερά) και % ανίχνευσης ανωμαλιών ανά τεταρτημόριο ντιλο που χρησιμοποιείται ως κατώφλι με την μέθοδο BiLSTM (δεξιά).	76
Διάγραμμα 22: Συγκριτικό διάγραμμα πριν και μετά την επέκταση των δεδομένων. (α) Επέκταση δεδομένων για ακραία δεδομένα (β) Επέκταση δεδομένων για δεδομένα παρέκκλισης.....	82
Διάγραμμα 23 : Καμπύλη απωλειών εκπαίδευ εκπαίδευσης και επικύρωσης	86
Διάγραμμα 24: Καμπύλη ακρίβειας εκπαίδευσης και επικύρωσης.....	87

II. Κατάλογος σχημάτων

Σχήμα 1: Φασικό πορτραίτο της $f(\theta) = \theta + 03\sin(3\theta)$	18
Σχήμα 2: Απλουστευμένη άποψη ενός τεχνητού νευρωνικού δικτύου τροφοδότησης	46
Σχήμα 3: Πλεονέκτημα της τεχνολογίας με βάση την τοπική πυκνότητα.....	54
Σχήμα 4 : Η παρακολουθούμενη γέφυρα και η θέση του επιταχυνσιόμετρου στο σώμα της γέφυρας και στον πύργο	79
Σχήμα 5 : Σχηματική αναπαράσταση της προτεινόμενης αρχιτεκτονικής CNN.....	85

III. Κατάλογος πινάκων

Πίνακας 1. Ορισμένα παραδείγματα τεχνικών ανίχνευσης ανωμαλιών με βάση την ταξινόμηση με χρήση νευρωνικών δικτύων	46
Πίνακας 2: Περιγραφή κάθε τύπου του μοτίβου των δεδομένων.....	80
Πίνακας 3: Αποτελέσματα της πρόβλεψης του συνόλου δοκιμής.....	88

Κεφάλαιο 1^ο: Εισαγωγή



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

Η ανίχνευση ανωμαλιών είναι ένα σημαντικό πρόβλημα που έχει ερευνηθεί σε διάφορους ερευνητικούς τομείς και τομείς εφαρμογών. Πολλές τεχνικές ανίχνευσης ανωμαλιών έχουν αναπτυχθεί ειδικά για ορισμένους τομείς εφαρμογών, ενώ άλλες είναι πιο γενικές. Η παρούσα έρευνα προσπαθεί να παράσχει μια δομημένη και περιεκτική επισκόπηση της έρευνας για την ανίχνευση ανωμαλιών. Έχουμε ομαδοποιήσει τις υπάρχουσες τεχνικές σε διάφορες κατηγορίες με βάση την υποκείμενη προσέγγιση που υιοθετεί κάθε τεχνική. Για κάθε κατηγορία προσδιορίσαμε τις βασικές παραδοχές που χρησιμοποιούνται από τις τεχνικές για τη διάκριση μεταξύ κανονικής και ανώμαλης συμπεριφοράς. Κατά την εφαρμογή μιας δεδομένης τεχνικής σε ένα συγκεκριμένο τομέα, αυτές οι παραδοχές μπορούν να χρησιμοποιηθούν ως κατευθυντήριες γραμμές για την αξιολόγηση της αποτελεσματικότητας της στον συγκεκριμένο τομέα. Για κάθε κατηγορία, παρέχουμε μια βασική τεχνική ανίχνευσης ανωμαλιών και στη συνέχεια παρουσιάζονται διάφορες παραλλαγές της. Αυτό το πρότυπο παρέχει μια ευκολότερη και συνοπτική κατανόηση των τεχνικών που ανήκουν σε κάθε κατηγορία. Περαιτέρω, για κάθε κατηγορία, προσδιορίζουμε τα πλεονεκτήματα και τα μειονεκτήματα των τεχνικών της συγκεκριμένης κατηγορίας. Παρέχουμε επίσης μια συζήτηση σχετικά με την υπολογιστική πολυπλοκότητα των τεχνικών, δεδομένου ότι αποτελεί σημαντικό ζήτημα σε πραγματικούς τομείς εφαρμογών.

Πιο συγκεκριμένα, ως προς τη δομή της εργασίας, το παρόν πρώτο κεφάλαιο είναι εισαγωγικό και προσφέρει τις βασικές πληροφορίες για την εργασία. Στο δεύτερο κεφάλαιο παρουσιάζονται τα δυναμικά συστήματα, ο ορισμός τους, διάφορα παραδείγματα συστημάτων, οι έννοιες της ροής και της εξέλιξης, καθώς και τα διακριτά δυναμικά συστήματα. Το τρίτο κεφάλαιο αναφέρεται στις χρονοσειρές και αναλύονται ο ορισμός και παραδείγματα χρονοσειρών, τα χαρακτηριστικά τους, τα διάφορα μοντέλα που έχουν αναπτυχθεί, καθώς και η έννοια της στασιμότητας. Το τέταρτο κεφάλαιο αναφέρεται στην ανίχνευση ανωμαλιών. Παρουσιάζονται τα διάφορα προβλήματα ανίχνευσης, η φύση των δεδομένων εισόδου και οι διάφοροι τύποι ανωμαλιών. Στο πέμπτο κεφάλαιο, παρατίθενται διάφοροι μέθοδοι ανίχνευσης ανωμαλιών. Αναλύονται οι τεχνικές που βασίζονται στην ταξινόμηση, στον πλησιέστερο γείτονα και στην ομαδοποίηση, στατιστικές τεχνικές και ο χειρισμός των



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

ανωμαλιών πλαισίου. Στο έκτο κεφάλαιο, παρουσιάζονται εφαρμογές ανίχνευσης ανωμαλιών στον τομέα της κυκλοφοριακής ροής και στα δομικά δεδομένα. Τέλος, στο έβδομο κεφάλαιο παρατίθεται τα συνολικά συμπεράσματα της εργασίας.

Κεφάλαιο 2^ο: Δυναμικά συστήματα

2.1 Ορισμός

Η δυναμική είναι πρωτίστως η μελέτη μιας χρονικά εξελικτικής διαδικασίας και το αντίστοιχο σύστημα εξισώσεων είναι γνωστό ως δυναμικό σύστημα. Γενικά, ένα σύστημα n διαφορικών εξισώσεων πρώτης τάξης στο χώρο R^n ονομάζεται δυναμικό σύστημα διάστασης n , το οποίο καθορίζει τη χρονική συμπεριφορά της εξελικτικής διαδικασίας. Οι εξελικτικές διαδικασίες μπορεί να είναι ντετερμινιστικές ή μη ντετερμινιστικές, να είναι πεπερασμένης ή μη πεπερασμένης διάστασης και να διαθέτουν τις ιδιότητες της διαφοροποιησιμότητας. Μια διαδικασία καλείται **ντετερμινιστική** εάν ολόκληρη η μελλοντική της πορεία και ολόκληρο το παρελθόν της καθορίζονται με μοναδικό τρόπο από την κατάστασή της στον παρόντα χρόνο. Διαφορετικά, η διαδικασία καλείται **μη ντετερμινιστική**. Ωστόσο, η διαδικασία μπορεί να είναι **ημι-ντετερμινιστική** (καθορισμένη, αλλά όχι με μοναδικό τρόπο). Στην κλασική μηχανική η κίνηση ενός συστήματος του οποίου το μέλλον και το παρελθόν καθορίζονται μοναδικά από τις αρχικές θέσεις και τις αρχικές ταχύτητες είναι ένα παράδειγμα ντετερμινιστικού δυναμικού συστήματος. Η εξελικτική διαδικασία μπορεί να περιγράφει, ως

- (i) Μια διαδικασία συνεχούς χρόνου
- (ii) Μια διαδικασία διακριτού χρόνου

Η διαδικασία συνεχούς χρόνου αναπαρίσταται με διαφορικές εξισώσεις, ενώ η διαδικασία διακριτού χρόνου με εξισώσεις διαφορών (ή χάρτες). Τα δυναμικά συστήματα συνεχούς χρόνου μπορούν να περιγραφούν μαθηματικά ως εξής. Έστω: $x = x(t) \in R, t \in I \subseteq R$ το διάνυσμα που αντιπροσωπεύει τη δυναμική ενός συνεχούς συστήματος (σύστημα συνεχούς χρόνου). Η μαθηματική αναπαράσταση του συστήματος μπορεί να γραφεί ως εξής:

$$\frac{dx}{dt} = \dot{x} = f(x, t)$$



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

(1.1)

όπου $f(x, t)$ είναι μια επαρκώς ομαλή συνάρτηση που ορίζεται σε κάποιο υποσύνολο $U \subset \mathbb{R}^n \times \mathbb{R}$.

Η μεταβλητή t ερμηνεύεται συνήθως ως χρόνος και η συνάρτηση $f(x, t)$ είναι γενικά μη γραμμική. Το χρονικό διάστημα μπορεί να είναι πεπερασμένο, ημιπεπερασμένο ή άπειρο. Από την άλλη πλευρά, το διακριτό σύστημα σχετίζεται με έναν διακριτό χάρτη (που δίνεται μόνο σε ισαπέχοντα χρονικά σημεία), έτσι ώστε από ένα σημείο x_0 , μπορεί κανείς να λάβει ένα σημείο x_1 το οποίο με τη σειρά του απεικονίζει σε

x_2 , και ούτω καθεξής. Με άλλα λόγια, $x_{n+1} = g(x_n) = g(g(x_{n-1}))$, κ.λπ. Αυτό γράφεται επίσης στη μορφή $x_{n+1} = g(x_n) = g^2(x_{n-1})$.

Εάν η δεξιά πλευρά της Εξίσωσης (1.1) είναι ρητά ανεξάρτητη από το χρόνο, τότε το σύστημα ονομάζεται **αυτόνομο**. Οι τροχιές ενός τέτοιου συστήματος δεν αλλάζουν στο χρόνο. Από την άλλη πλευρά, εάν η δεξιά πλευρά της Εξίσωσης (1.1) έχει ρητή εξάρτηση από το χρόνο, τότε το σύστημα ονομάζεται **μη αυτόνομο**. Ένα n -διάστατο μη αυτόνομο σύστημα μπορεί να μετατραπεί σε αυτόνομο μορφή εισάγοντας μια νέα εξαρτημένη μεταβλητή x_{n+1} , όπου $x_{n+1} = t$. Γενικά, η εύρεση της λύσης της Εξίσωσης (1.1) είναι δύσκολη ή μερικές φορές αδύνατη όταν η συνάρτηση $f(x, t)$ είναι μη γραμμική, εκτός από ορισμένες ειδικές περιπτώσεις. Παραδείγματα αυτόνομων και μη αυτόνομων συστημάτων δίνονται παρακάτω:

— **Αυτόνομα συστήματα:**

- $\ddot{x} + a\dot{x} + b = 0, a, b > 0$: Πρόκειται για έναν αποσβεσμένο γραμμικό αρμονικό ταλαντωτή. Οι παράμετροι a και b είναι, αντίστοιχα, η ισχύς της απόσβεσης και η ισχύς της γραμμικής δύναμης επαναφοράς.
- $\ddot{x} + \omega^2 \sin x = 0, \omega = \sqrt{g/L}$, όπου g είναι η επιτάχυνση της βαρύτητας, L το μήκος του νήματος. Πρόκειται για έναν απλό μη γραμμικό ταλαντωτή χωρίς απόσβεση (εκκρεμές).

— **Μη αυτόνομα συστήματα:**



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

- $\ddot{x} + a\dot{x} + b = f \cos \omega t$, $a, b > 0$. Αυτό είναι ένα παράδειγμα γραμμικού ταλαντωτή με εξωτερική χρονοεξαρτώμενη δύναμη f και ω είναι το πλάτος και η συχνότητα της κινητήριας δύναμης.
- $\ddot{x} + a\dot{x} + \omega_0^2 x + bx^3 = f \sin \omega t$. Πρόκειται για έναν μη γραμμικό ταλαντωτή Duffing με κυβική δύναμη επαναφοράς. a είναι η ισχύς της απόσβεσης, ω_0 είναι η ιδιοσυχνότητα και b είναι η ισχύς της μη γραμμικής δύναμης επαναφοράς.

2.2 Παραδείγματα δυναμικών συστημάτων

Το πιο συνηθισμένο παράδειγμα δυναμικού συστήματος είναι τα Νευτώνεια συστήματα που διέπονται από τον νόμο της κίνησης. Ο νόμος αυτός ορίζει ότι η επιτάχυνση ενός σώματος καθορίζεται από τη δύναμη ανά μονάδα μάζας. Η δύναμη μπορεί να είναι συνάρτηση της ταχύτητας (\dot{x}) και της θέσης (x) και έτσι τα Νευτώνεια συστήματα παίρνουν τη μορφή:

$$m\ddot{x} = F(x, \dot{x}) \text{ όπου } (m=\text{μάζα και } F=\text{δύναμη}) \quad (1.2)$$

Η εξίσωση (1.2) μπορεί να γραφεί ως ένα σύστημα δύο διαφορικών εξισώσεων πρώτης τάξης ως εξής:

$$\dot{x} = y \text{ και } \dot{y} = F(x, y) \quad (1.3)$$

Το σύστημα (1.3) μπορεί να θεωρηθεί ως ένα δυναμικό σύστημα δύο διαστάσεων στο επίπεδο xy και η δυναμική είναι ένα σύνολο τροχιών που δίνουν τη χρονική εξέλιξη της κίνησης.

Οι πληθυσμοί δύο ανταγωνιστικών ειδών (πληθυσμοί θηρευτών και θηραμάτων) θα μπορούσαν να μοντελοποιηθούν μαθηματικά. Το μοντέλο πληθυσμού θηρευτή-θηράματος διατυπώθηκε για πρώτη φορά από τον *Alfred J. Lotka* (1880-1949) το έτος 1910 και αργότερα από τον *Vito Volterra* (1860-1940) το έτος 1926. Το μοντέλο αυτό έχει μείνει γνωστό ως μοντέλο αρπακτικού-θηράματος Lotka-Volterra. Σε αυτό το μοντέλο ο πληθυσμός της αλεπούς θηρεύει τον πληθυσμό των κουνελιών. Η πληθυσμιακή πυκνότητα του κουνελιού επηρεάζει την πληθυσμιακή πυκνότητα της

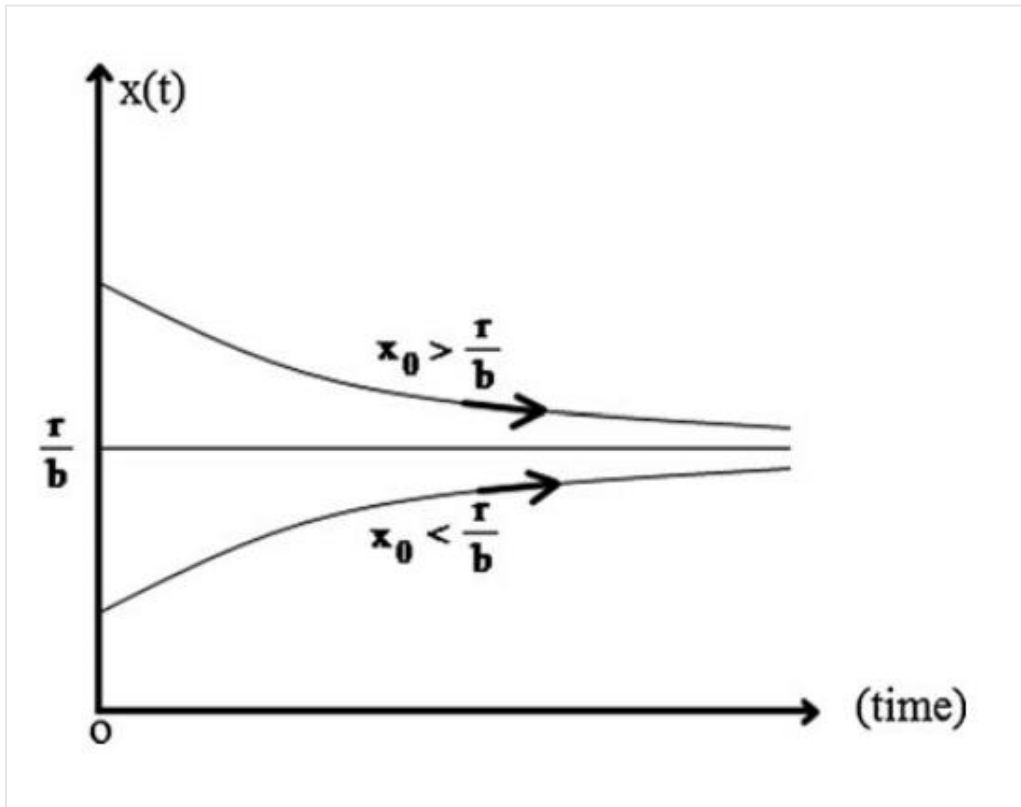


ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

αλεπούς, καθώς η τελευταία βασίζεται στην πρώτη για τροφή. Εάν η πυκνότητα του κουνελιού είναι υψηλή, ο πληθυσμός της αλεπούς μειώνεται, ενώ όταν ο πληθυσμός της αλεπούς αυξάνεται, ο πληθυσμός των κουνελιών μειώνεται. Όταν μειώνεται ο πληθυσμός των κουνελιών, μειώνεται και ο πληθυσμός των αλεπούδων. Όταν ο πληθυσμός των αλεπούδων πέφτει, τα κουνέλια μπορούν να πολλαπλασιαστούν ξανά και ούτω καθεξής. Η αύξηση ή η μείωση δύο πληθυσμών θα μπορούσε να αναλυθεί χρησιμοποιώντας τις αρχές του δυναμικού συστήματος. Οι δυναμικές εξισώσεις για το μοντέλο θηρευτή - θηράματος δίνονται ως εξής:

$$\left. \begin{aligned} \dot{x} &= ax + bxy \\ \dot{y} &= -\gamma y + \delta xy \end{aligned} \right\} \quad (1.4)$$

όπου x δηλώνει την πληθυσμιακή πυκνότητα του θηράματος και y , την πληθυσμιακή πυκνότητα του θηρευτή. Η παράμετρος a αντιπροσωπεύει τον ρυθμό ανάπτυξης του θηράματος απουσία αλληλεπίδρασης με τα αρπακτικά, ενώ η παράμετρος γ αντιπροσωπεύει τον ρυθμό θανάτου των αρπακτικών απουσία αλληλεπίδρασης με το θήραμα και β , d είναι οι παράμετροι αλληλεπίδρασης και είναι όλες σταθερές (για το απλό μοντέλο). Χρησιμοποιώντας τη δυναμική αρχή μπορεί κανείς να λάβει μια αναγκαία συνθήκη για τη συνύπαρξη των δύο ειδών. Σε αυτό το μοντέλο η επιβίωση των θηρευτών εξαρτάται εξ ολοκλήρου από τον πληθυσμό του θηράματος. Εάν αρχικά $x = 0$, τότε $y = -\gamma x$, δηλαδή $y(t) = y(0)e^{-\gamma t}$ και $y(t) \rightarrow 0$ όταν $t \rightarrow \infty$ (...).



Διάγραμμα 1: Γραφική αναπαράσταση του μοντέλου ανάπτυξης του πληθυσμού

Πηγή: (Layek, 2016)

2.3 Ροή

Η χρονικά εξελικτική διαδικασία μπορεί να περιγραφεί ως ροή ενός διανυσματικού πεδίου. Γενικά, η ροή χρησιμοποιείται συχνά εντός του πλαισίου των δυναμικών συστημάτων στο σύνολό τους και όχι για την εξέλιξη ενός συστήματος σε ένα συγκεκριμένο σημείο. Η λύση $x(t)$ ενός συστήματος $\dot{x} = f(x)$ το οποίο ικανοποιεί τη συνθήκη $x(t_0) = x_0$ δίνει το παρελθόν $t < t_0$ και το μέλλον $t > t_0$ εξελίξεις του συστήματος. Μαθηματικά, η ροή ορίζεται από τη σχέση $\varphi_t(x) : U \rightarrow R^n$ όπου $\varphi_t(x) = \varphi(t, x)$ είναι μια ομαλή διανυσματική συνάρτηση των $x \in U \subset R^n$ και $t \in I \subset R$ που ικανοποιεί την εξίσωση:

$$\frac{d}{dt} \varphi_t(x) = f(\varphi_t(x))$$

(1.5)

Ροές στην R:



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

Έστω ότι ένα ρευστό ρέει κατά μήκος της πραγματικής γραμμής με τοπική ταχύτητα $f(x)$, σε ένα μονοδιάστατο αυτόνομο σύστημα που αναπαρίσταται από το $x = f(x)$, $x \in R$. Αυτό το φανταστικό ρευστό ονομάζεται ρευστό φάσης και η πραγματική γραμμή ονομάζεται γραμμή φάσης. Για την επίλυση του συστήματος $x = f(x)$ ξεκινώντας από μια αυθαίρετη αρχική θέση x_0 , τοποθετούμε ένα φανταστικό σωματίδιο, που ονομάζεται σημείο φάσης, στο x_0 και παρακολουθούμε πώς κινείται μαζί με τη ροή στη γραμμή φάσης σε μεταβαλλόμενο χρόνο t . Καθώς περνάει ο χρόνος, το σημείο φάσης (x, t) στο μονοδιάστατο σύστημα $x = f(x)$ με $x(0) = x_0$ κινείται κατά μήκος του άξονα x σύμφωνα με κάποια συνάρτηση $\varphi(t, x_0)$. Η συνάρτηση $\varphi(t, x_0)$ ονομάζεται τροχιά για δεδομένη αρχική κατάσταση x_0 και το σύνολο $\{\varphi(t, x_0) | t \in I \subset R\}$ είναι η τροχιά του $x_0 \in R$. Το σύνολο όλων των ποιοτικών τροχιών του συστήματος ονομάζεται **πορτρέτο φάσης**.

Ροές στο R^2 :

Έστω ένα δισδιάστατο σύστημα που αντιπροσωπεύεται από τις ακόλουθες εξισώσεις $x = f(x, y)$ και $y = g(x, y)$ με $(x, y) \in R^2$. Ένα φανταστικό σωματίδιο ρευστού ρέει στο επίπεδο R^2 , γνωστό ως επίπεδο φάσης του συστήματος. Η διαδοχή των καταστάσεων που δίνονται παραμετρικά από $x = x(t)$ και $y = y(t)$ διαγράφουν μια καμπύλη μέσω κάποιου αρχικού σημείου $P(x(t_0), y(t_0))$ ονομάζεται **διαδρομή φάσης**. Το σύνολο $\{\varphi(t, x_0) | t \in I \subset R\}$ είναι η τροχιά του x στο R^2 . Υπάρχει άπειρος αριθμός τροχιών που θα γέμιζαν το επίπεδο φάσης. Αλλά η ποιοτική συμπεριφορά μπορεί να προσδιοριστεί με την απεικόνιση μερικών τροχιών με διαφορετικές αρχικές συνθήκες. Το πορτραίτο φάσης δείχνει πώς μεταβάλλεται η ποιοτική συμπεριφορά ενός συστήματος καθώς τα x και y μεταβάλλονται με το χρόνο t . Μια τροχιά ονομάζεται περιοδική αν $x(t + p) = x(t)$ για κάποιο $p > 0$, για όλα τα t . Ο μικρότερος ακέραιος p για τον οποίο ικανοποιείται η σχέση ονομάζεται πρώτη περίοδος της τροχιάς. Οι ροές στο R δεν μπορούν να έχουν ταλαντωτική ή κλειστή διαδρομή.

Επιπλέον, υπάρχουν και ροές σε χώρους μεγαλύτερων διαστάσεων R^n , με αυτόνομα συστήματα και n διαφορικές εξισώσεις. Ωστόσο, στα πλαίσια της παρούσας εργασίας δεν θα αναφερθούμε περαιτέρω.



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

2.4 Εξέλιξη

Έστω ένα σύστημα $\dot{x} = f(x)$, $x \in R^N$ με αρχικές συνθήκες $x(t_0) = x_0$. Έστω $E \subset R^N$ ένα ανοικτό σύνολο και $f \in C^{-1}(E)$. Για $x_0 \in E$, έστω $\varphi(t, x_0)$ μια λύση του παραπάνω συστήματος στο μέγιστο διάστημα ύπαρξης $I(x_0) \subset R$. Η απεικόνιση $\varphi_t: R^N \rightarrow R^N$ που ορίζεται από τη σχέση $\varphi_t(x_0) = \varphi(t, x_0)$ είναι γνωστή ως τελεστής εξέλιξης του συστήματος. Η γραμμική ροή για το σύστημα $\dot{x} = Ax$ με $x(t_0) = x$, ορίζεται από τη σχέση $\varphi_t: R^N \rightarrow R^N$ και $\varphi_t = e^{-At}$, ο εκθετικός πίνακας.

Γενικά, ένα δυναμικό σύστημα μπορεί να θεωρηθεί ως ομάδα μη γραμμικών/ γραμμικών τελεστών που εξελίσσονται ως $\{\varphi_t(x), t \in R, x \in R^N\}$.

2.5 Διακριτά δυναμικά συστήματα

Μια εξελικτική διαδικασία μπορεί επίσης να εκφραστεί μαθηματικά ως διακριτά συστήματα στο χρόνο. Τα διακριτά συστήματα περιγράφονται από διαγράμματα (εξισώσεις διαφορών). Η σύνθεση του διαγράμματος παράγει τη δυναμική ή τη ροή ενός διακριτού συστήματος. Πρόκειται για μια ακολουθία επαναλήψεων, όπως για μια δεδομένη συνάρτηση $f: E \rightarrow E \subseteq R^n$ με αρχικό σημείο x_0 , η ακολουθία των επαναλήψεων μπορεί να παραχθεί ως εξής

$$x_0, f(x_0), f(f(x_0)), f(f(f(x_0))), \dots \tag{1.6}$$

Η ακολουθία αυτή μπορεί να είναι πεπερασμένη ή άπειρη. Είναι ενδιαφέρον το πώς συμπεριφέρεται αυτή η ακολουθία μετά από μερικές επαναλήψεις. Τα διακριτά διαγράμματα καλύπτουν πολύ μεγαλύτερο εύρος δυναμικής από τα συνεχή συστήματα. Η έννοια της ροής που δημιουργείται από ένα διακριτό σύστημα, η μαθηματική αναπαράστασή της, οι συνθέσεις διαγραμμάτων, οι τροχιές, τα πορτραίτα φάσης, τα σταθερά σημεία, τα περιοδικά σημεία, οι περιοδικοί κύκλοι, οι σταθερότητες, τα υπερβολικά, μη υπερβολικά σταθερά σημεία με ορισμένα σημαντικά θεωρήματα και παραδείγματα θα παρουσιαστούν παρακάτω.

2.5.1 Διαγράμματα και ροές

Γενικά, ένα διάγραμμα είναι μια συνάρτηση $f: E \rightarrow E$. Η κατάσταση x_{n+1} στο $(n+1)$ -οστό στάδιο εκφράζεται ως προς το προηγούμενο στάδιο x_n μέσω της



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

σχέσης $x_{n+1} = f(x_n)$. Αν η αρχική κατάσταση ή η κατάσταση εκκίνησης είναι $x_0 \in E$, η ακολουθία των καταστάσεων δίνεται από $x_0, x_1, \dots, x_n, \dots$ στο E . Ένα διακριτό σύστημα παράγει μια ροή που παριστάνεται από $\varphi_t(x)$ στο E έτσι ώστε $f(x) = \varphi_t(x), x \in E$ και t είναι ένας διακριτός χρόνος στο R . Έτσι το διακριτό δυναμικό σύστημα είναι η εξέλιξη της οικογένειας διαγραμμάτων $\{f^n\}$, $n = 0, \pm 1, \pm 2, \dots$ στο E .

2.5.2 Τροχιές

Δεδομένου ενός μονοδιάστατου χάρτη $f: R \rightarrow R$ και ένα σημείο $x_0 \in E$ η τροχιά του x_0 ορίζεται ως εξής

$$O^+(x_0) = \{f^k(x_0)\}_{k=0}^{\infty} = \{x_0, f(x_0), f^2(x_0), \dots, f^n(x_0), \dots\} \quad (1.7)$$

Ομοίως, η αντίστροφη τροχιά του x_0 ορίζεται ως εξής:

$$O^-(x_0) = \{f^{-1}(x_0), f^{-2}(x_0), \dots, f^{-n}(x_0), \dots\} \quad (1.8)$$

Η αντίστροφη τροχιά υπάρχει αν η f είναι ομοιόμορφη (συνεχής και έχει συνεχή αντίστροφο). Γενικά, η τροχιά της x_0 κάτω από έναν ομοιομορφισμό f ορίζεται ως εξής:

$$\begin{aligned} O(x_0) &= \{f^k(x_0)\}_{k=-\infty}^{\infty} \\ &= \{\dots, f^{-n}(x_0), \dots, f^{-2}(x_0), f^{-1}(x_0), x_0, f(x_0), f^2(x_0), \dots, f^n(x_0), \dots\} \end{aligned} \quad (1.9)$$

Ομοίως, μπορούν να ληφθούν τροχιές χαρτών υψηλότερης διάστασης.

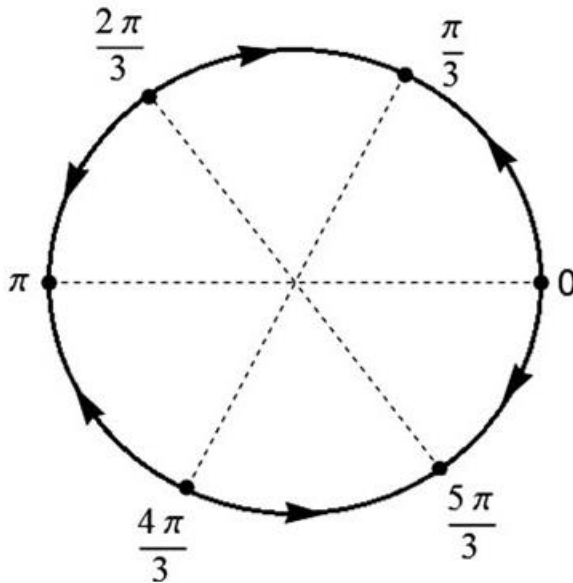
2.5.3 Διαγράμματα φάσεων

Τα διαγράμματα φάσεων χρησιμοποιούνται συχνά στα δυναμικά συστήματα για να αναπαραστήσουν γραφικά τη δυναμική ενός χάρτη (map). Το διάγραμμα φάσης είναι ένα διάγραμμα που παρουσιάζει τις πιθανές μεταβαλλόμενες θέσεις μιας συνάρτησης χάρτη και τα βέλη υποδεικνύουν την αλλαγή των θέσεων κάτω από επαναλήψεις του χάρτη. Έστω ένας απλός μονοδιάστατος χάρτης $f: [0, 2\pi] \rightarrow [0, 2\pi]$ που ορίζεται



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

από την $f(\theta) = \theta + 0.3\sin(3\theta)$. Το πορτρέτο φάσης αυτού του χάρτη απεικονίζεται στο Σχήμα 1. Το σχήμα δείχνει ότι τα έξι σημεία ικανοποιούν την σχέση $f(\theta) = \theta$. Τα βέλη δείχνουν ότι η ροή προς τα τρία σημεία $\frac{\pi}{3}, \pi, \frac{5\pi}{3}$ και η ροή απομακρύνεται από τα άλλα τρία σημεία $0, \frac{2\pi}{3}, \frac{4\pi}{3}$.



Σχήμα 1: Φασικό πορτραίτο της $f(\theta) = \theta + 0.3\sin(3\theta)$

Πηγή: (Layek, 2016)

Κεφάλαιο 3^ο: Χρονοσειρές

3.1 Ορισμός και παραδείγματα χρονοσειρών

Η ανάλυση χρονοσειρών είναι ένας ειδικός τρόπος ανάλυσης μιας ακολουθίας σημείων, δεδομένων που συλλέγονται σε ένα χρονικό διάστημα. Στην ανάλυση χρονοσειρών, οι αναλυτές καταγράφουν τα σημεία δεδομένων σε σταθερά διαστήματα κατά τη διάρκεια μιας καθορισμένης χρονικής περιόδου και δεν καταγράφουν απλώς τα σημεία δεδομένων διακεκομμένα ή τυχαία. Σύμφωνα με τον Mooris Hamburg, «μια χρονοσειρά είναι ένα σύνολο στατιστικών παρατηρήσεων τοποθετημένων σε χρονολογική σειρά», ενώ σύμφωνα με τον W.Z. Hirsch «ο κύριος στόχος της ανάλυσης χρονολογικών σειρών είναι η κατανόηση, η ερμηνεία και η αξιολόγηση της μεταβολής



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

των οικονομικών φαινομένων με την ελπίδα να προβλεφθεί ορθότερα η πορεία των μελλοντικών γεγονότων».

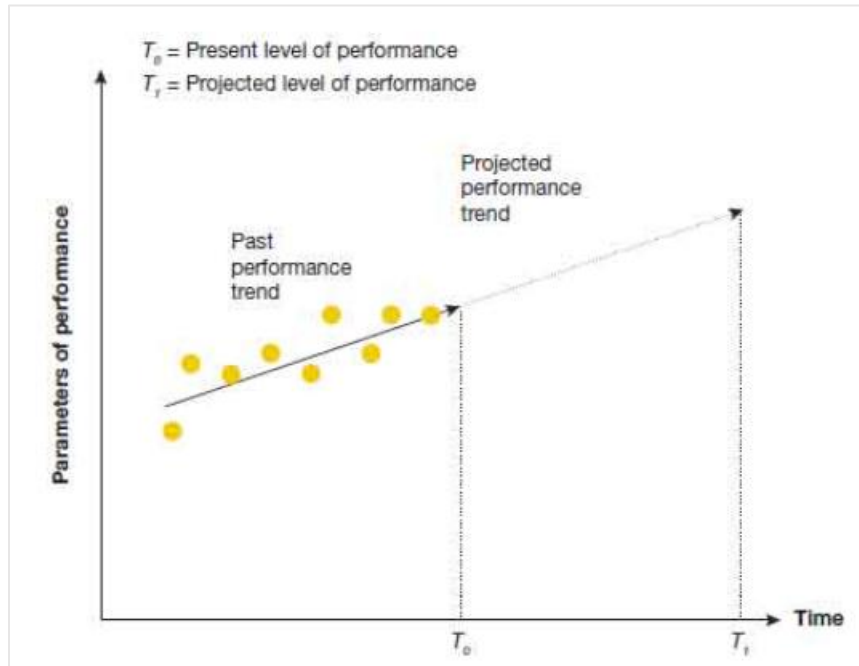
Η ανάλυση χρονοσειρών είναι μια στατιστική τεχνική που ασχολείται με την ανάλυση τάσεων και δεδομένων χρονοσειρών. Η ανάλυση χρονοσειρών εισήλθε στην ιατρική όταν εφευρέθηκαν τα πρώτα πρακτικά ηλεκτροκαρδιογραφήματα (ΗΚΓ), τα οποία μπορούν να διαγνώσουν καρδιακές παθήσεις καταγράφοντας τα ηλεκτρικά σήματα που διέρχονται από την καρδιά, το 1901. Οι πίνακες θνησιμότητας του John Graunt ήταν ένα από τα πρώτα αποτελέσματα σκέψης κοντινή στις χρονοσειρές, που εφαρμόστηκε σε ιατρικά ζητήματα. Στο βιβλίο του, ο Graunt παρουσίασε τους πρώτους πίνακες ζωής, που είναι γνωστοί ως πίνακες θνησιμότητας. Αυτοί οι πίνακες δείχνουν την πιθανότητα ένα άτομο συγκεκριμένης ηλικίας να πεθάνει πριν από τα επόμενα γενέθλιά του. Τα δεδομένα από τις χρονοσειρές είναι περιοδικές χρονικές περίοδοι που έχουν μετρηθεί σε τακτά χρονικά διαστήματα ή έχουν συγκεντρωθεί σε συγκεκριμένες χρονικές στιγμές. Επαναδιατυπώνοντας, μια χρονοσειρά είναι απλώς μια συλλογή σημείων δεδομένων τοποθετημένων χρονολογικά, και η ανάλυση χρονοσειρών είναι η πράξη της ερμηνείας αυτών των δεδομένων. Το Ακαθάριστο Εγχώριο Προϊόν (ΑΕΠ), ο Δείκτης Τιμών Καταναλωτή, ο Δείκτης SP 500 και τα ποσοστά ανεργίας είναι παραδείγματα δεδομένων χρονοσειρών στα οικονομικά. Τα δεδομένα χρονοσειρών στις κοινωνικές επιστήμες θα μπορούσαν να περιλαμβάνουν πληροφορίες σχετικά με την αύξηση του πληθυσμού, τα μεταναστευτικά πρότυπα, τα ποσοστά γεννήσεων και τις πολιτικές μεταβλητές. Στα πλαίσια της εργασίας θα παρουσιαστούν οι ορισμοί, οι παραδοχές, οι στόχοι, η εφαρμογή, τα μοντέλα χρονοσειρών και η στασιμότητα μαζί με τον τρόπο αναγνώρισής τους.

3.1.1 Παραδείγματα

Ανάλυση πωλήσεων

Η ανάλυση χρονοσειρών χρησιμοποιείται συχνά από τους λιανοπωλητές για να εξετάσουν τις μακροπρόθεσμες τάσεις των συνολικών τους πωλήσεων. Η ανάλυση χρονοσειρών είναι πολύ χρήσιμη για την εξέταση των τάσεων των πωλήσεων σε μηνιαία, εποχιακή και ετήσια βάση. Ως αποτέλεσμα, τα καταστήματα λιανικής πώλησης είναι σε θέση να προβλέπουν καλύτερα τις πωλήσεις τους για την επερχόμενη σεζόν, καθώς και τον αριθμό των εργαζομένων και των εμπορευμάτων που θα

χρειαστών σε ορισμένες χρονικές στιγμές κατά τη διάρκεια του έτους. Για παράδειγμα, το παρακάτω διάγραμμα δείχνει μια τάση των πραγματικών στοιχείων πωλήσεων για μια ορισμένη χρονική περίοδο και μια εκτιμώμενη τάση.



Διάγραμμα 2: Εξαγωγή της τρέχουσας τάσης στο μέλλον για την πρόβλεψη των πωλήσεων

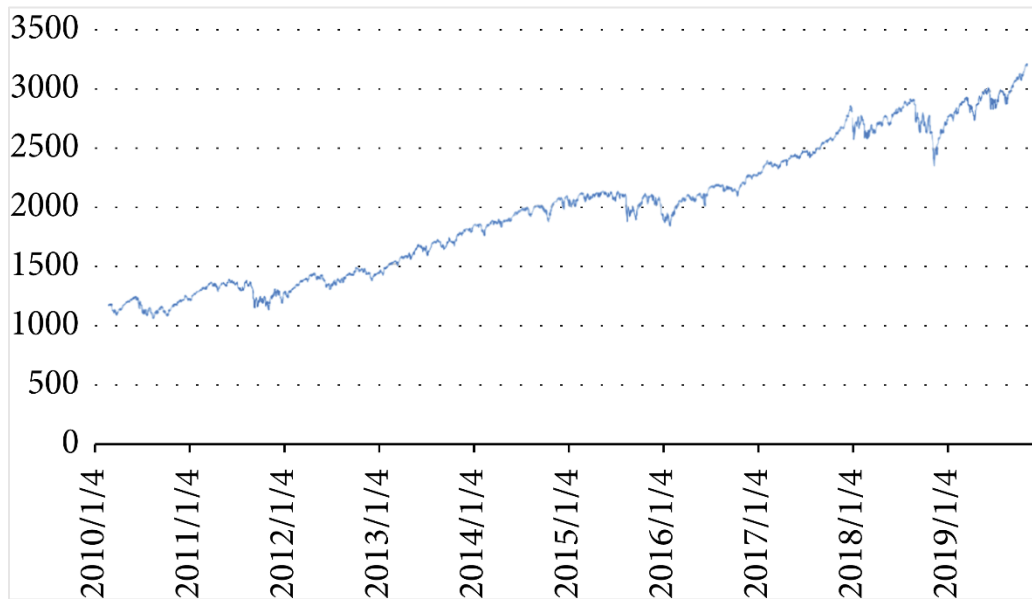
Πηγή: (Grzegorzek, 2023)

Μάρκετινγκ αποθεμάτων

Προκειμένου να κατανοήσουν καλύτερα τις τάσεις στις διάφορες τιμές των μετοχών, οι χρηματιστές χρησιμοποιούν συνήθως την ανάλυση χρονοσειρών. Ιδιαίτερα χρήσιμα είναι τα διαγράμματα χρονοσειρών, τα οποία βοηθούν τους αναλυτές και τους εμπόρους μετοχών να κατανοήσουν την τάση και την κατεύθυνση μιας συγκεκριμένης τιμής μετοχής. Για παράδειγμα, στο παρακάτω διάγραμμα, παρουσιάζεται η ακολουθία του χρηματιστηριακού δείκτη S&P 500. Από το σχήμα διαπιστώνεται ότι εντός του επιλεγμένου χρονικού διαστήματος, ο δείκτης S&P 500 παρουσιάζει γενικά μια σταθερή αυξητική τάση.



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ



Διάγραμμα 3: Ακολουθία κλεισίματος του δείκτη μετοχών S&P

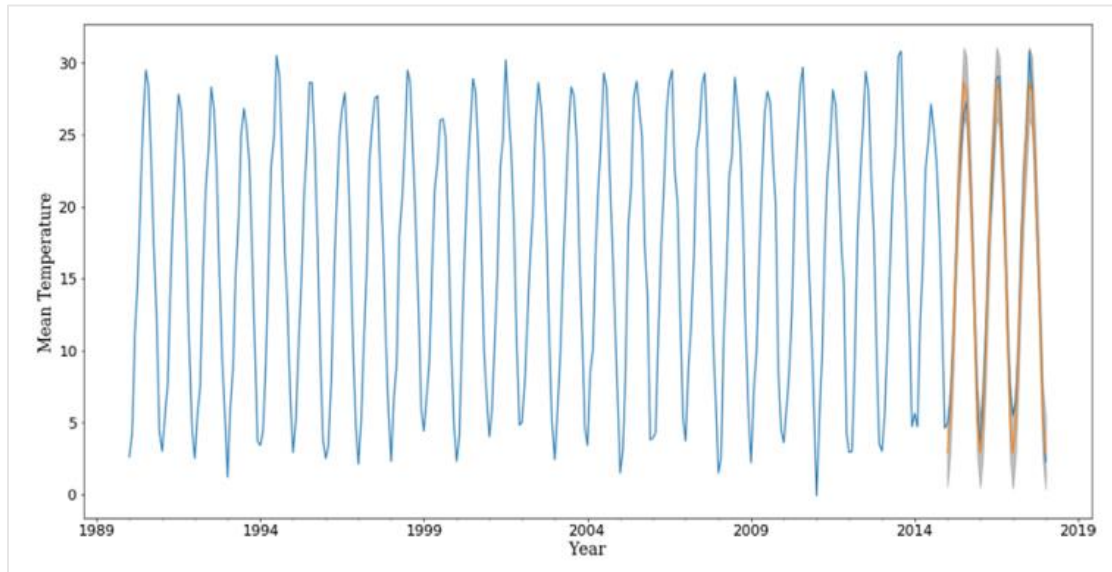
Πηγή: (Xiao & Su, 2022)

Πρόγνωση καιρού

Οι μετεωρολόγοι χρησιμοποιούν συνήθως την ανάλυση χρονοσειρών για να προβλέψουν τις θερμοκρασίες για διάφορους μήνες και εποχές του έτους. Περιλαμβάνει την εκτίμηση της θερμοκρασίας, την παρατήρηση της κλιματικής αλλαγής, τον εντοπισμό εποχιακών μεταβολών και την πρόβλεψη του καιρού. Παράδειγμα, αποτελεί το παρακάτω γράφημα, που παρουσιάζει τις προβλεπόμενες θερμοκρασίες (κόκκινη γραμμή) των επόμενων 36 μηνών στο Nanjing, με βάση τα 35 έτη δεδομένων του παρελθόντος (μπλε γραμμή).



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ



Διάγραμμα 4: Σύγκριση μεταξύ των πραγματικών τιμών και των τιμών πρόβλεψης στη θερμοκρασία στο Nanjing.

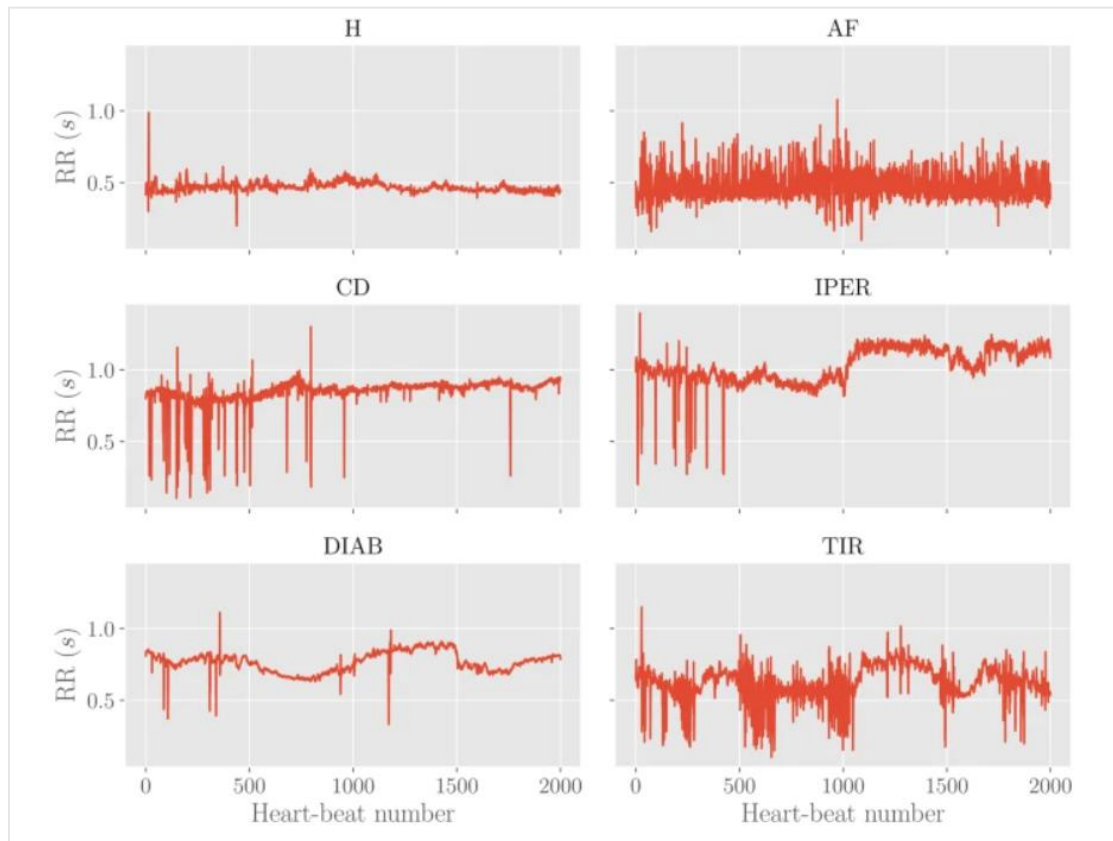
Πηγή: (Chen et al., 2018)

Καρδιακός ρυθμός

Η ανάλυση χρονοσειρών χρησιμοποιείται επίσης στον ιατρικό τομέα για την παρακολούθηση του καρδιακού ρυθμού των ασθενών που μπορεί να λαμβάνουν ορισμένα φάρμακα, ώστε να διασφαλίζεται ότι ο καρδιακός ρυθμός δεν παρουσιάζει υπερβολικές διακυμάνσεις κατά τη διάρκεια οποιασδήποτε δεδομένης ώρας της ημέρας. Για παράδειγμα στην έρευνα των Agliari et al. (2020), οι ασθενείς χωρίστηκαν σε 6 κύριες κατηγορίες: υγιείς (H), πάσχοντες από κολπική μαρμαρυγή (AF), από συμφορητική καρδιακή ανεπάρκεια (δηλ. καρδιακή αποσυμφόρηση, CD), από διαβήτη (DIAB), από υπο- ή υπερθυρεοειδισμό (TIR) ή από υπέρταση (TENS). Μετά από καταγραφές Holter, κάθε ασθενής συσχετίζεται με μια χρονοσειρά RR, δηλαδή μια σειρά χρονικών διαστημάτων μεταξύ δύο διαδοχικών καρδιακών παλμών. Όπως φαίνεται οι απεικονίσεις των διαγραμμάτων της κάθε κατηγορίας είναι αρκετά διαφορετικές μεταξύ τους.



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ



Διάγραμμα 5: Παραδείγματα χρονοσειρών RR για τις κατηγορίες ασθενών.

Κάθε διάγραμμα δείχνει τα πρώτα 2000 σημεία (δηλ. καρδιακά χτυπήματα) με τα αντίστοιχα διαστήματα RR (σε δευτερόλεπτα) για τις διάφορες εξεταζόμενες κλάσεις.

Πηγή: (Agliari et al., 2020)

3.2 Χαρακτηριστικά χρονοσειρών

Στην διαδικασία ανάλυσης χρονοσειρών εντοπίζονται τέσσερις κύριοι στόχοι, όπως η περιγραφή, η επεξήγηση, η πρόβλεψη και ο έλεγχος.

Περιγραφή

Η απεικόνιση των δεδομένων και η λήψη βασικών περιγραφικών μέτρων των θεμελιωδών χαρακτηριστικών των σειρών αποτελούν τα πρώτα βήματα της ανάλυσης. Οι μετρήσεις αυτές μπορεί να είναι τόσο απλές όσο η αναζήτηση τάσεων ή τόσο



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

σύνθετες όσο η ανάλυση εποχιακών μεταβολών. Παρόλο που η τάση αυτή των τιμών είναι ασυνεπής, το προηγούμενο σχήμα δείχνει ένα σταθερό εποχικό πρότυπο μεταβολής των τιμών. Ένα γράφημα καθιστά δυνατή την αναζήτηση "ακραίων" παρατηρήσεων, που δεν φαίνεται να είναι συνεπείς με τα υπόλοιπα δεδομένα. Η εμφάνιση σημείων καμπής όπου η ανοδική τάση μετατράπηκε απότομα σε καθοδική γίνεται δυνατή με τη γραφική παράσταση της χρονοσειράς. Αν υπάρχει σημείο καμπής, μπορεί να χρειαστεί να συνδεθούν διαφορετικά μοντέλα με τα δύο μισά της σειράς.

Επεξήγηση

Ήταν δυνατό να χρησιμοποιηθεί η διακύμανση σε μια χρονική σειρά για να εξηγηθεί η διακύμανση σε μια άλλη σειρά, δεδομένου ότι οι παρατηρήσεις έγιναν σε δύο ή περισσότερες μεταβλητές. Αυτό μπορεί να οδηγήσει σε καλύτερη κατανόηση.

Πρόβλεψη

Κάποιος μπορεί να θέλει να προβλέψει τις μελλοντικές τιμές της χρονοσειράς δεδομένης μιας παρατηρούμενης χρονοσειράς. Για παράδειγμα, η εξέταση των οικονομικών και βιομηχανικών χρονοσειρών αποτελεί καθοριστική συνθήκη για την πρόβλεψη των πωλήσεων.

Έλεγχος

Όταν δημιουργούνται χρονοσειρές για τη μέτρηση της ποιότητας μιας παραγωγικής διαδικασίας (με τη δυνατότητα ρύθμισης της διαδικασίας), υπάρχουν διάφοροι τύποι τεχνικών ελέγχου. Οι παρατηρήσεις απεικονίζονται σε ένα διάγραμμα ελέγχου κατά τη διάρκεια του ποιοτικού ελέγχου και ο ελεγκτής ανταποκρίνεται ανάλογα, εξετάζοντας τα διαγράμματα. Η σειρά προσαρμόζεται με τη χρήση ενός στοχαστικού μοντέλου. Οι αναμενόμενες μελλοντικές τιμές της σειράς χρησιμοποιούνται για την αλλαγή των μεταβλητών εισόδου της διεργασίας και τη διατήρηση της διεργασίας στην πορεία της.

3.2.1 Συνιστώσες των χρονοσειρών

Μια χρονοσειρά αποτελείται ουσιαστικά από τέσσερις συνιστώσες, την τάση, τις εποχικές διακυμάνσεις, τις κυκλικές μεταβολές και τις τυχαίες παραλλαγές.

Τάση



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

Μια τάση αποκαλύπτει ένα τακτικό μοτίβο στα δεδομένα. Κατά τη διάρκεια μιας συγκεκριμένης, εκτεταμένης χρονικής περιόδου, μπορεί να αυξηθεί και να μειωθεί. Η τάση είναι μια συνεπής, μακροπρόθεσμη γενική κατεύθυνση της κίνησης των δεδομένων. Τα δεδομένα δεν χρειάζεται να κινούνται με τον ίδιο τρόπο για να υπάρχει τάση. Κατά τη διάρκεια μιας μακράς χρονικής περιόδου, η κίνηση ή η κατεύθυνση μπορεί να αυξομειώνεται, αλλά η τάση δεν πρέπει να αλλάζει σε ένα μοτίβο. Ο αριθμός των σχολείων, η γεωργική παραγωγή, η αύξηση του πληθυσμού κ.λπ. είναι μερικές περιπτώσεις τάσεων. Αξίζει να σημειωθεί ότι η τάση μπορεί να αλλάζει κατεύθυνση κατά τη διάρκεια διαφορετικών χρονικών περιόδων, είτε προς τα πάνω, είτε προς τα κάτω, είτε σταθερά. Μια τάση μπορεί να είναι είτε γραμμική είτε μη γραμμική.

Εποχιακές διακυμάνσεις

Οι εποχιακές διακυμάνσεις είναι βραχυπρόθεσμες μεταβολές σε χρονοσειρές, συνήθως εντός χρονικού διαστήματος μικρότερου του έτους. Κατά τη διάρκεια της 12μηνιαίας διάρκειας των χρονοσειρών, εκείνες παρουσιάζουν συνήθως την ίδια τάση εξέλιξης, είτε ανοδική είτε καθοδική. Τα χρονοδιαγράμματα σε ωριαία, ημερήσια, εβδομαδιαία, τριμηνιαία και μηνιαία βάση χρησιμοποιούνται συχνά για την παρακολούθηση αυτών των διακυμάνσεων. Φυσικοί ή τεχνητοί παράγοντες ή αλλαγές μπορούν να προκαλέσουν εποχικές διαφορές. Οι εποχικές μεταβολές επηρεάζονται σε μεγάλο βαθμό από τις πολλαπλές εποχές και τις τεχνητές μεταβολές. Η χρήση ομπρελών αυξάνεται δραματικά κατά την περίοδο των βροχών, οι πωλήσεις κλιματιστικών αυξάνονται κατά τη διάρκεια του καλοκαιριού και οι συγκομιδές βασίζονται στην εποχή είναι παραδείγματα εποχικότητας. Ορισμένες νόρμες που δημιουργήθηκαν από τον άνθρωπο εμφανίζουν ξεκάθαρα εποχικές διακυμάνσεις. Οι εποχιακές διαφορές επηρεάζονται από τις διακοπές, τις παραδόσεις, τα στυλ, τις συνήθειες και τα ποικίλα γεγονότα όπως οι γάμοι. Η περίοδος εποχιακών διακυμάνσεων δεν πρέπει να χρησιμοποιείται για να κρίνει κάποιος αν μια επιχείρηση τα πάει καλύτερα ή χειρότερα.

Κυκλικές μεταβολές

Οι μεταβολές στις χρονοσειρές που εμφανίζονται για διάστημα μεγαλύτερο του ενός έτους ονομάζονται κυκλικές μεταβολές. Τέτοιες ταλαντωτικές κινήσεις του χρόνου έχουν συχνά διάρκεια μεγαλύτερη από ένα έτος. Μια πλήρης περίοδος λειτουργίας



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

ονομάζεται είτε κύκλος είτε «επιχειρηματικός κύκλος». Οι κυκλικές μεταβολές περιέχουν τέσσερις φάσεις, την ευημερία, την υποχώρηση, την ύφεση και την ανάκαμψη. Μπορεί να έχουν τακτικό ή μη περιοδικό χαρακτήρα. Συνήθως, οι κυκλικές διακυμάνσεις οφείλονται στο συνδυασμό δύο ή περισσότερων οικονομικών δυνάμεων και των αλληλεπιδράσεών τους.

Τυχαίες παραλλαγές

Υπάρχει ένα άλλο είδος κίνησης που μπορεί να παρατηρηθεί στην περίπτωση των χρονοσειρών. Πρόκειται για καθαρή ακανόνιστη και τυχαία κίνηση. Όπως υποδηλώνει το όνομα, καμία υπόθεση ή τάση δεν μπορεί να χρησιμοποιηθεί για να υποδηλώσει ακανόνιστες ή τυχαίες κινήσεις σε μια χρονοσειρά. Αυτά τα αποτελέσματα είναι απρόβλεπτα, ακανόνιστα και ανεξέλεγκτα από τη φύση τους. Οι σεισμοί, ο πόλεμος, η πείνα και οι πλημμύρες είναι μερικά παραδείγματα τυχαίων στοιχείων χρονοσειρών.

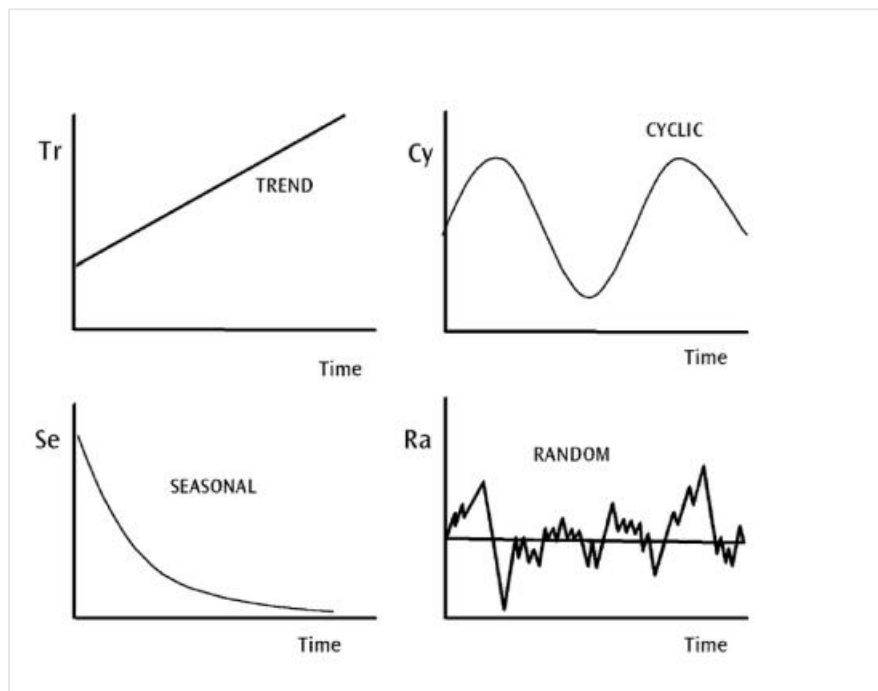
3.2.2 Διαγράμματα χρονοσειρών

Ένα διάγραμμα χρονοσειράς είναι ένα γράφημα που δείχνει πληροφορίες που συλλέγονται με την πάροδο του χρόνου από οποιαδήποτε διαδικασία. Το γράφημα μπορεί να χρησιμοποιηθεί για τον προσδιορισμό των ιστορικών τάσεων στα δεδομένα, καθώς και για το αν τα σημεία δεδομένων είναι τυχαία ή παρουσιάζουν κάποιο μοτίβο. Για παράδειγμα, για να ελέγξει κάποιος αν ο όγκος των κλήσεων που εισέρχονται σε ένα τηλεφωνικό κέντρο είναι σταθερός από μήνα σε μήνα ή αν υπάρχει κάποιο συγκεκριμένο μοτίβο στα δεδομένα, όπως μια τάση μείωσης ή μια αυξητική τάση μπορεί να χρησιμοποιήσει ένα διάγραμμα χρονοσειρών. Ένα διάγραμμα χρονοσειρών μπορεί επίσης να χρησιμοποιηθεί για τον οπτικό εντοπισμό σημείων σταθερότητας σε μια διαδικασία, προκειμένου να διαπιστωθεί η σταθερότητά της.

Τα δεδομένα σειράς είναι μια ακολουθία σημείων δεδομένων σε χρονολογική σειρά που χρησιμοποιείται από τις επιχειρήσεις για την ανάλυση δεδομένων του παρελθόντος και την πραγματοποίηση μελλοντικών προβλέψεων. Τα εν λόγω σημεία δεδομένων είναι ένα σύνολο παρατηρήσεων σε συγκεκριμένες χρονικές στιγμές και ίσα διαστήματα, συνήθως με δείκτη χρονολογίας και αντίστοιχη τιμή.

Στο παρακάτω διάγραμμα, απεικονίζονται οι τέσσερις παραλλαγές μεταβολών των δεδομένων χρονοσειράς. Η μεταβολή τάσης (trend) κινείται προς τα πάνω ή προς τα

κάτω με ένα λογικά προβλέψιμο μοτίβο για μεγάλο χρονικό διάστημα. Η μεταβολή εποχικότητας (seasonal) μπορεί να είναι τακτική και περιοδική, δηλαδή να επαναλαμβάνεται κατά τη διάρκεια μιας συγκεκριμένης περιόδου, όπως μια ημέρα, μια εβδομάδα, ένας μήνας, μια εποχή κ.ά. Η κυκλική (cyclic) μεταβολή αντιστοιχεί σε επιχειρηματικούς ή οικονομικούς κύκλους «άνθησης-κατάρρευσης» ή είναι κυκλική με κάποια άλλη μορφή. Η τυχαία (random) μεταβολή δεν εμπίπτει σε καμία από τις τρεις παραπάνω ταξινομήσεις.



Διάγραμμα 6: Διάγραμμα χρονοσειράς με διαφορετικές συνιστώσες

Πηγή: (Jose, 2022)

3.3 Μοντέλα χρονοσειρών

3.3.1 Προσθετικό μοντέλο

Στο προσθετικό μοντέλο, αναπαρίσταται μια συγκεκριμένη παρατήρηση σε μια χρονοσειρά ως το άθροισμα αυτών των τεσσάρων συνιστωσών.

$$O = M + S + C + E$$

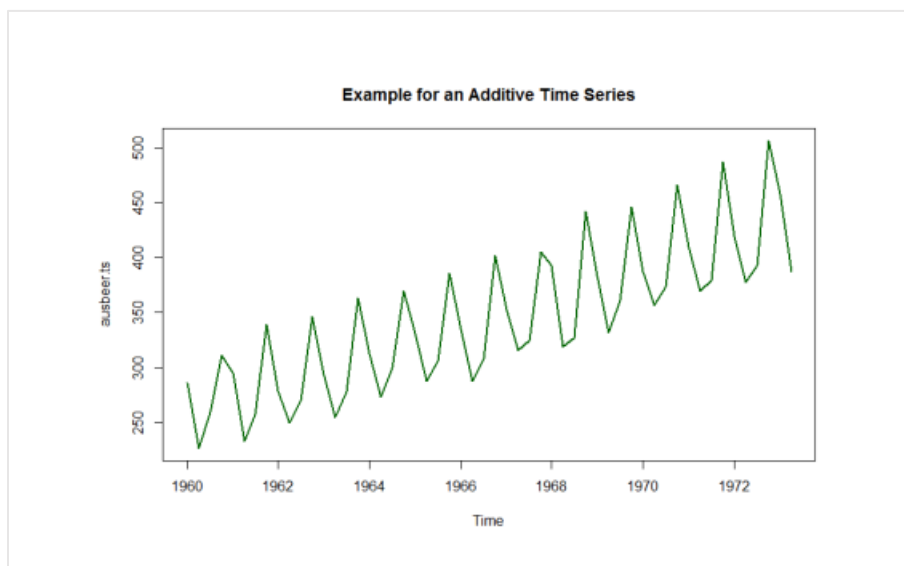


ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

όπου το O αντιπροσωπεύει τα αρχικά δεδομένα, το M αντιπροσωπεύει την τάση. Το S αντιπροσωπεύει τις εποχιακές μεταβολές, το C αντιπροσωπεύει τις κυκλικές μεταβολές και το E αντιπροσωπεύει τις ακανόνιστες μεταβολές. Μπορεί επίσης να γραφεί ότι:

$$Z(t) = M(t) + S(t) + C(t) + E(t)$$

(1.10)



Διάγραμμα 7: Προσθετικό μοντέλο με $M(t) + S(t) + E(t)$

Πηγή: (Jose, 2022)

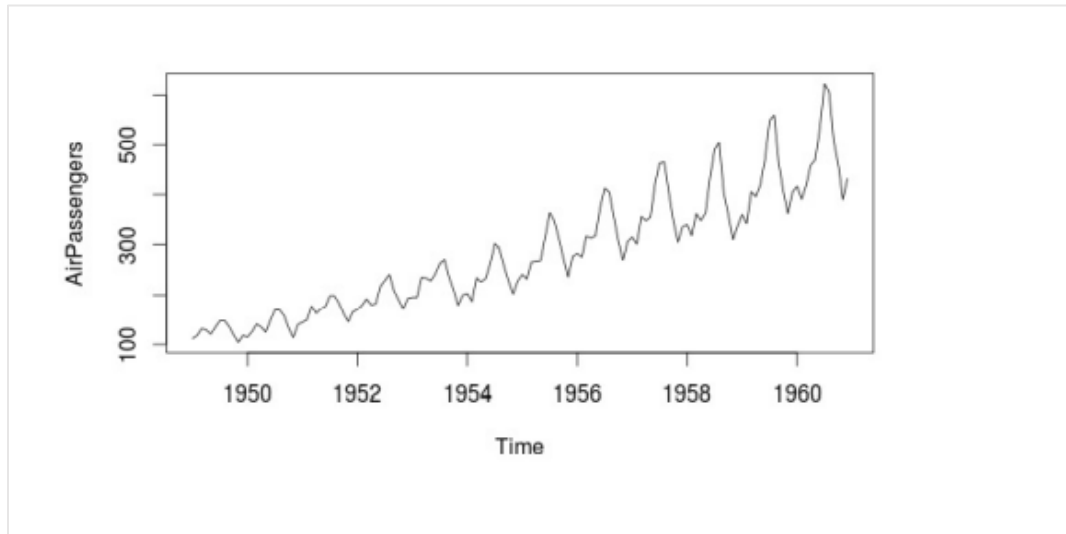
3.3.2 Πολλαπλασιαστικό μοντέλο

Σε αυτό το μοντέλο, τέσσερις συνιστώσες έχουν πολλαπλασιαστική σχέση. Έτσι, μια συγκεκριμένη παρατήρηση σε μια χρονοσειρά αναπαρίσταται ως το γινόμενο αυτών των τεσσάρων συνιστωσών:

$$O = MSCE$$



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ



Διάγραμμα 8: Παράδειγμα για το πολλαπλασιαστικό μοντέλο όπου O, M, S, C και E αντιπροσωπεύουν τους όρους όπως στο προσθετικό μοντέλο

Πηγή: (Jose, 2022)

Με άλλο τρόπο, μπορεί να γραφεί:

$$Z(t) = M(t)S(t)C(t)E(t) \quad (1.11)$$

Το μοντέλο αυτό είναι το πιο συχνά χρησιμοποιούμενο για την αποσύνθεση χρονολογικών σειρών. Για να εξαλειφθεί κάθε αμφιβολία μεταξύ των δύο μοντέλων, θα πρέπει να διευκρινιστεί ότι στο πολλαπλασιαστικό μοντέλο τα S, C και E είναι δείκτες που εκφράζονται ως δεκαδικά ποσοστά, ενώ στο προσθετικό μοντέλο τα S, C και E είναι ποσοτικές αποκλίσεις μιας τάσης που μπορεί να εκφραστεί ως εποχιακή, κυκλική και ακανόνιστη στη φύση της.

3.4 Στασιμότητα

Η στασιμότητα αναφέρεται στη χρονική διακύμανση ορισμένων ή όλων των στατιστικών στοιχείων μιας τυχαίας διαδικασίας, όπως ο μέσος όρος, η αυτοσυσχέτιση, η κατανομή n -τάξης. Μπορούν να οριστούν δύο τύποι στασιμότητας με αυστηρή έννοια σταθερή και ευρεία αίσθηση σταθερή. Μια χρονοσειρά λέγεται ότι είναι στάσιμη εάν οι παρατηρήσεις δεν εξαρτώνται από το χρόνο, δηλαδή οι στατιστικές της ιδιότητες δεν θα μεταβάλλονται με το χρόνο, άρα θα έχουν σταθερή μέση τιμή,



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

διακύμανση και συνδιακύμανση. Όταν εμφανίζεται τάση ή εποχικότητα σε μια χρονοσειρά τότε εκείνη δεν μπορεί να είναι στάσιμη. Η εμφάνιση τάσης συνεπάγεται την αλλαγή στα δεδομένα με την πάροδο του χρόνου και συνεπώς τη μεταβολή των στατιστικών χαρακτηριστικών που παράγει μια χρονοσειρά, όπως είναι ο μέσος όρος. Αντίστοιχα και στην περίπτωση της εμφάνισης περιοδικότητας μεταβάλλεται το μοτίβο των δεδομένων για τακτά χρονικά διαστήματα και κατ' επέκταση η διακύμανση. Τα δεδομένα της χρονοσειράς θα πρέπει να είναι στάσιμα προκειμένου να μπορεί κάποιος να προβλέψει. Υπάρχουν διάφοροι τρόποι για να επικυρωθεί η υπόθεση της στασιμότητας και να μετατραπούν τα δεδομένα χρονοσειρών σε στάσιμα.

3.4.1 Τρόποι ελέγχου της στασιμότητας

Έλεγχος μέσου όρου

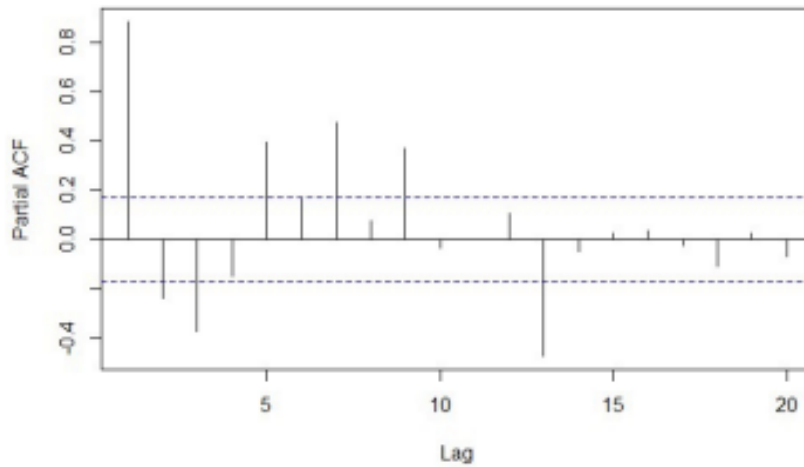
Η μέθοδος αυτή ελέγχει πρώτα το μέσο όρο της χρονοσειράς για το σύνολο των παρατηρήσεων. Στη συνέχεια παίρνει μια σειρά παρατηρήσεων από τις αρχικές παρατηρήσεις για μια συγκεκριμένη χρονική περίοδο. αν οι μέσοι όροι διαφέρουν σημαντικά μπορούμε να συμπεράνουμε ότι τα δεδομένα της χρονοσειράς δεν είναι στάσιμα. Εάν οι μέσοι όροι εξισώνονται τότε μπορεί να υπάρχει στασιμότητα στα δεδομένα της χρονοσειράς.

Διαγράμματα αυτοσυσχέτισης (Acf)

Η αυτοσυσχέτιση αντιπροσωπεύει το βαθμό ομοιότητας μεταξύ μιας δεδομένης χρονοσειράς και μιας καθυστερημένης εκδοχής της σε διαδοχικά χρονικά διαστήματα. Μετράει τη σχέση μεταξύ μιας μεταβλητής τρέχουσας αξίας και των προηγούμενων τιμών της. Μια αυτοσυσχέτιση με τιμή +1 αντιπροσωπεύει μια τέλεια θετική συσχέτιση, ενώ μια αυτοσυσχέτιση αρνητική -1 αντιπροσωπεύει μια τέλεια αρνητική συσχέτιση. Τα διαγράμματα ACF είναι γνωστά ως διαγράμματα συνάρτησης αυτοσυσχέτισης, μια απεικόνιση της σειριακής συσχέτισης σε δεδομένα που μεταβάλλονται με την πάροδο του χρόνου. Το διάγραμμα ACF μπορεί να δημιουργηθεί εύκολα με τη χρήση της συνάρτησης ACF. Από το διάγραμμα, εάν η αυτοσυσχέτιση είναι μεγαλύτερη από το όριο της, τότε τα δεδομένα της χρονοσειράς διαθέτουν στοιχεία μη στασιμότητας.



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ



Διάγραμμα 9 : Διάγραμμα αυτοσυσχέτισης (Acf)

Πηγή: (Jose, 2022)

Έλεγχος Dickey-Filler (Adf)

Ο επαυξημένος έλεγχος Dickey-Filler είναι γνωστός ως έλεγχος ADF. Η στασιμότητα της χρονοσειράς εξετάζεται με τη χρήση αυτής της τεχνικής. Έλεγχοι, όπως ο ADF, είναι απαραίτητοι δεδομένου ότι όλα τα μοντέλα πρόβλεψης απαιτούν η χρονοσειρά να είναι στάσιμη. Η απουσία σταθερής χρονοσειράς είναι η μηδενική υπόθεση στον ADF. Τα στατιστικά στοιχεία του ελέγχου θα πρέπει να έχουν τιμή κάτω από τις κρίσιμες τιμές προκειμένου να απορριφθεί η μηδενική υπόθεση (1 τοις εκατό, 5 τοις εκατό, 10 τοις εκατό). Εάν ικανοποιείται αυτό το κριτήριο, η σειρά μπορεί να θεωρηθεί στάσιμη και η μηδενική υπόθεση μπορεί να απορριφθεί. Η σημαντικότητα των στατιστικών στοιχείων ελέγχου επηρεάζει το επίπεδο εμπιστοσύνης. Για παράδειγμα, εάν η τιμή είναι μικρότερη από την κρίσιμη τιμή του 5 τοις εκατό, μπορεί να δηλωθεί με βεβαιότητα 95 τοις εκατό ότι η χρονοσειρά είναι στάσιμη.

3.4.2 Τρόποι απαλοιφής τάσης

Για να μελετηθούν δεδομένα χρονολογικών σειρών, τα δεδομένα θα πρέπει να είναι σταθερής φύσης. Από τη φύση της, κάθε χρονοσειρά δεδομένων διαθέτει διάφορες συνιστώσες όπως τάση, εποχικότητα, κυκλικές και τυχαίες διακυμάνσεις. Έτσι, για να επιτευχθεί στασιμότητα θα πρέπει να εξαλειφθεί η τάση και η εποχιακή συνιστώσα

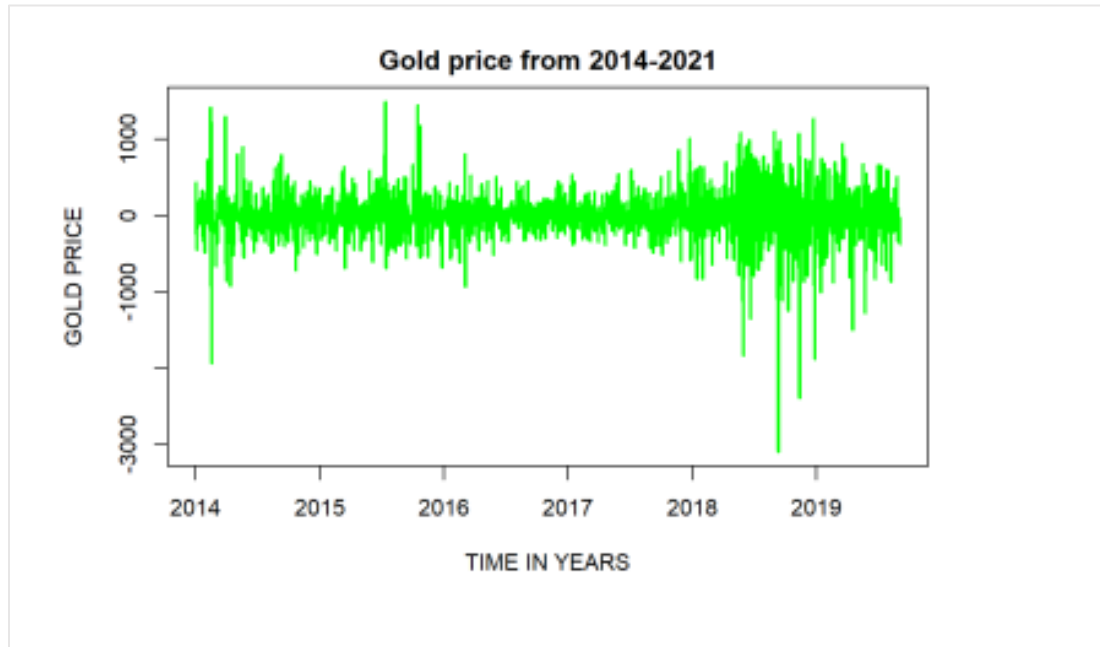


ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

μαζί με την τυχαία διακύμανση. Αν και η τυχαία διακύμανση δεν μπορεί να εξαλειφθεί εντελώς, μπορεί όμως να ελαχιστοποιηθεί. Οι μέθοδοι όπως η μέθοδος των διαφορών, η μέθοδος των ελαχίστων τετραγώνων βοηθούν στην εξάλειψη της συνιστώσας τάσης. Η μέθοδος του κινητού μέσου όρου εξαλείφει τη συνιστώσα τάσης που υπάρχει στη χρονοσειρά.

Μέθοδος της διαφοράς (Detrending)

Η επίτευξη στασιμότητας με μη στάσιμα δεδομένα γίνεται συχνότερα με αυτόν τον τρόπο. Με τη διαφοροποίηση, δημιουργείται τυπικά ένα νέο σύνολο δεδομένων που περιέχει τις μεταβολές μεταξύ των παρατηρήσεων που γίνονται σε μια δεδομένη χρονική στιγμή και σε μια προηγούμενη. Παρόλο που μερικές φορές χρειάζεται να διαφοροποιηθούν τα δεδομένα δύο φορές για να διαπιστωθεί η στασιμότητα, συνήθως αρκεί μία διαφοροποίηση. Με μεγαλύτερα σύνολα δεδομένων, η μία παρατήρηση που αφαιρεί κάθε διαφορά από το σύνολο δεδομένων καθίσταται ασήμαντη.



Διάγραμμα 10: Διάγραμμα της τιμής του χρυσού μετά την εξάλειψη της τάσης με τη μέθοδο της διαφοράς

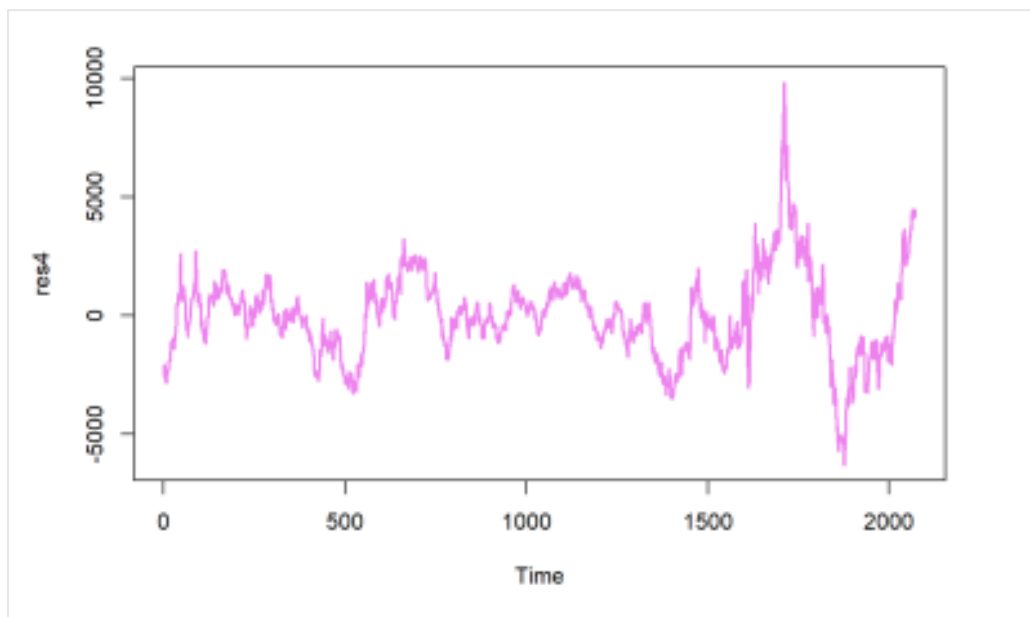
Πηγή: (Jose, 2022)



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

Μέθοδος των ελαχίστων τετραγώνων (Detrending)

Με τη μέθοδο αυτή εξαλείφεται η συνιστώσα της τάσης, καθιστώντας έτσι τα δεδομένα στάσιμα. Πολλές από τις χρονικές μεταβλητές που συναντώνται κατά την ανάλυση χρονοσειρών εξαρτώνται η μία από την άλλη, γεγονός που προκαλεί μη στασιμότητα. Ως εκ τούτου, χρησιμοποιείται η μέθοδος των ελαχίστων τετραγώνων για να εξαλειφθεί η μη στασιμότητα που προκαλείται από την τάση και άλλες συνιστώσες. Η σχέση μεταξύ μιας γνωστής ανεξάρτητης μεταβλητής και μιας άγνωστης εξαρτημένης μεταβλητής αντιπροσωπεύεται από κάθε σημείο της προσαρμοσμένης καμπύλης. Προκειμένου να προσαρμοστούν μέσω των παρεχόμενων σημείων με τη μέθοδο των ελαχίστων τετραγώνων, τα δεδομένα της χρονοσειράς συχνά αποσυντονίζονται. Εάν η καμπύλη που παράγεται από τα ελάχιστα τετράγωνα ικανοποιεί το διάγραμμα ACF για στασιμότητα, σταματάει η διαδικασία, προτείνοντας νέα μοντέλα με διάφορους βαθμούς. Σε αντίθετη περίπτωση, η διαδικασία συνεχίζει μέχρι να επιτευχθεί η παραπάνω συνθήκη.



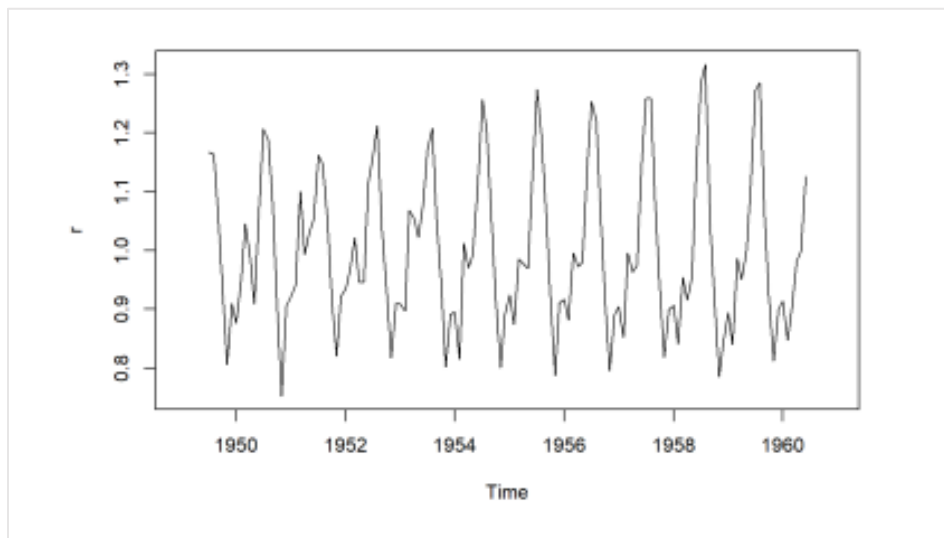
Διάγραμμα 11: Διάγραμμα της τιμής του χρυσού μετά την εξάλειψη της τάσης με τη μέθοδο των ελαχίστων τετραγώνων

Μέθοδος των κινητών μέσων όρων (Deseasonalisation)



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

Αυτή η τεχνική εξομάλυνσης μιας χρονοσειράς είναι αρκετά αποτελεσματική. Είναι ζωτικής σημασίας να εξαλειφθεί κάθε σημαντική διακύμανση των δεδομένων πριν από την πρόβλεψη, ιδίως το εποχικό φαινόμενο. Η μέθοδος των κινητών μέσων όρων εξισορροπεί τις βραχυπρόθεσμες διακυμάνσεις και αποσαφηνίζει την τάση. Είναι ζωτικής σημασίας, πριν από τη χρήση αυτής της διαδικασίας, να είναι γνωστή η χρονική διάρκεια των δεδομένων. Εάν υπάρχει χρονική διάρκεια, χρησιμοποιείται ένας κινητός μέσος όρος n περιόδων. Ο κινητός μέσος όρος υπολογίζεται αυτόματα όταν το n είναι περιττό, με κέντρο τα σημεία των δεδομένων, αλλά όταν το n είναι άρτιο, πρέπει πρώτα τα δεδομένα να κεντραριστούν.



Διάγραμμα 12: Διάγραμμα των επιβατών με τη μέθοδο των κινητών μέσων όρων

Κεφάλαιο 4^ο: Ανίχνευση ανωμαλιών

Η ανίχνευση ανωμαλιών αναφέρεται στο πρόβλημα της εύρεσης μοτίβων σε δεδομένα που δεν συμμορφώνονται με την αναμενόμενη συμπεριφορά. Αυτά τα μη συμμορφούμενα μοτίβα αναφέρονται συχνά ως ανωμαλίες, ακραίες τιμές, ασύμφωνες παρατηρήσεις, εξαιρέσεις, εκτροπές, εκπλήξεις ή ιδιαιτερότητες σε διάφορους τομείς εφαρμογών. Από αυτές, οι ανωμαλίες και οι ακραίες τιμές είναι δύο όροι που χρησιμοποιούνται συχνότερα στο πλαίσιο της ανίχνευσης ανωμαλιών. Η ανίχνευση ανωμαλιών βρίσκει εκτεταμένη χρήση σε ένα ευρύ φάσμα εφαρμογών, όπως η ανίχνευση απάτης για πιστωτικές κάρτες, ασφάλειες ή υγειονομική περίθαλψη, η ανίχνευση εισβολών για την ασφάλεια στον κυβερνοχώρο, η ανίχνευση σφαλμάτων σε



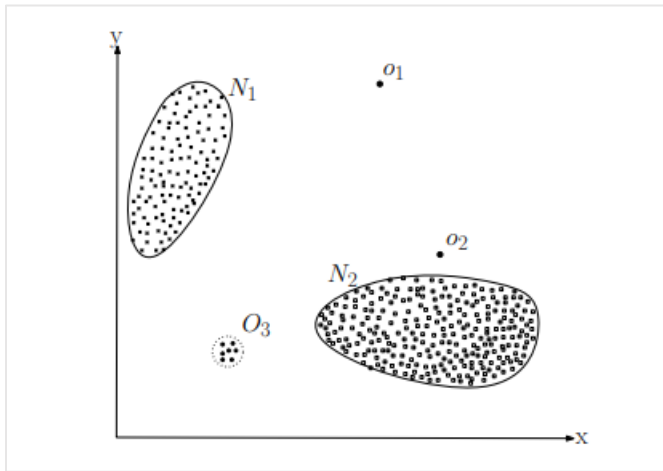
ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

κρίσιμα συστήματα ασφαλείας και η στρατιωτική επιτήρηση για εχθρικές δραστηριότητες.

Η σπουδαιότητα της ανίχνευσης ανωμαλιών οφείλεται στο γεγονός ότι οι ανωμαλίες στα δεδομένα μεταφράζονται σε σημαντικές (και συχνά κρίσιμες) πληροφορίες που μπορούν να αναληφθούν σε μια ευρεία ποικιλία τομέων εφαρμογών. Για παράδειγμα, ένα ανώμαλο μοτίβο κίνησης σε έναν υπολογιστή δικτύου θα μπορούσε να σημαίνει ότι πρόκειται για έναν παραβιασμένο υπολογιστή, που στέλνει ευαίσθητα δεδομένα σε μη εξουσιοδοτημένο προορισμό (Kumar, 2005). Μια ανώμαλη εικόνα μαγνητικής τομογραφίας μπορεί να υποδηλώνει την παρουσία κακοήθων όγκων (Spence et al., 2001). Ανωμαλίες σε δεδομένα συναλλαγών με πιστωτικές κάρτες θα μπορούσαν να υποδηλώνουν κλοπή πιστωτικής κάρτας ή ταυτότητας ή ανώμαλες ενδείξεις από αισθητήρα διαστημικού σκάφους θα μπορούσαν να υποδηλώνουν βλάβη σε κάποιο εξάρτημα του διαστημικού σκάφους.

Η ανίχνευση ακραίων τιμών ή ανωμαλιών στα δεδομένα έχει μελετηθεί στην κοινότητα της στατιστικής ήδη από τον 19ο αιώνα. Με την πάροδο του χρόνου, έχουν αναπτυχθεί διάφορες τεχνικές ανίχνευσης ανωμαλιών σε διάφορες ερευνητικές κοινότητες. Πολλές από αυτές τις τεχνικές έχουν αναπτυχθεί ειδικά για ορισμένους τομείς εφαρμογών, ενώ άλλες είναι πιο γενικές.

Οι ανωμαλίες είναι μοτίβα στα δεδομένα που δεν συμμορφώνονται με μια σαφώς καθορισμένη έννοια της κανονικής συμπεριφοράς. Το Διάγραμμα 13 απεικονίζει τις ανωμαλίες σε ένα απλό δισδιάστατο σύνολο δεδομένων. Τα δεδομένα έχουν δύο κανονικές περιοχές, N_1 και N_2 , καθώς οι περισσότερες παρατηρήσεις βρίσκονται σε αυτές τις δύο περιοχές. Σημεία που βρίσκονται αρκετά μακριά από τις περιοχές αυτές, π.χ. τα σημεία o_1 και o_2 , και τα σημεία στην περιοχή O_3 , αποτελούν ανωμαλίες.



Διάγραμμα 13: Ένα απλό παράδειγμα ανωμαλιών σε ένα δισδιάστατο σύνολο δεδομένων

Πηγή: (Chandola et a., 2009)

Οι ανωμαλίες μπορεί να προκληθούν στα δεδομένα για διάφορους λόγους, όπως κακόβουλη δραστηριότητα, π.χ. απάτη με πιστωτικές κάρτες, εισβολή στον κυβερνοχώρο, τρομοκρατική δραστηριότητα ή φθορά στις οδικές υποδομές, αλλά όλοι οι λόγοι έχουν ένα κοινό χαρακτηριστικό, ότι ενδιαφέρουν τον αναλυτή. Το «ενδιαφέρον» ή η συνάφεια των ανωμαλιών με την πραγματική ζωή είναι ένα βασικό χαρακτηριστικό της ανίχνευσης ανωμαλιών.

Η ανίχνευση ανωμαλιών σχετίζεται, αλλά διαφέρει από την αφαίρεση θορύβου (Teng et al., 1990) και την προσαρμογή στο θόρυβο (Rousseeuw & Leroy, 1987), θεωρίες οι οποίες ασχολούνται με τον ανεπιθύμητο θόρυβο στα δεδομένα. Ο θόρυβος μπορεί να οριστεί ως ένα φαινόμενο στα δεδομένα το οποίο δεν ενδιαφέρει τον αναλυτή, αλλά δρα ως εμπόδιο στην ανάλυση των δεδομένων. Η απομάκρυνση του θορύβου οδηγείται από την ανάγκη απομάκρυνσης των ανεπιθύμητων αντικειμένων πριν από την εκτέλεση οποιασδήποτε ανάλυσης δεδομένων στα δεδομένα. Η προσαρμογή στο θόρυβο αναφέρεται στην ανοσοποίηση της εκτίμησης ενός στατιστικού μοντέλου έναντι ανώμαλων παρατηρήσεων.

Ένα άλλο θέμα που σχετίζεται με την ανίχνευση ανωμαλιών είναι η ανίχνευση καινοτομίας (Markou & Singh 2003), η οποία αποσκοπεί στην ανίχνευση προηγουμένως μη παρατηρημένων (αναδυόμενων, νέων) προτύπων στα δεδομένα. Η



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

διάκριση μεταξύ των νέων προτύπων και των ανωμαλιών είναι ότι τα νέα πρότυπα συνήθως ενσωματώνονται στο κανονικό μοντέλο μετά την ανίχνευσή τους.

4.1 Προβλήματα ανίχνευσης

Σε αφηρημένο επίπεδο, μια ανωμαλία ορίζεται ως ένα πρότυπο που δεν συμμορφώνεται με την αναμενόμενη κανονική συμπεριφορά. Μια απλή προσέγγιση ανίχνευσης ανωμαλιών, επομένως, είναι ο ορισμός μιας περιοχής που αντιπροσωπεύει την κανονική συμπεριφορά και η ανακήρυξη κάθε παρατήρησης στα δεδομένα που δεν ανήκει σε αυτή την κανονική περιοχή ως ανωμαλία. Όμως, διάφοροι παράγοντες καθιστούν αυτή την φαινομενικά απλή προσέγγιση πολύ δύσκολη:

- Ο ορισμός μιας κανονικής περιοχής που να περιλαμβάνει κάθε πιθανή κανονική συμπεριφορά είναι πολύ δύσκολος. Επιπλέον, το όριο μεταξύ κανονικής και ανώμαλης συμπεριφοράς δεν είναι συχνά ακριβές. Έτσι, μια ανώμαλη παρατήρηση που βρίσκεται κοντά στο όριο μπορεί στην πραγματικότητα να είναι κανονική και το αντίστροφο.
- Όταν οι ανωμαλίες είναι αποτέλεσμα κακόβουλων ενεργειών, οι κακόβουλοι αντίπαλοι συχνά προσαρμόζονται για να κάνουν τις ανώμαλες παρατηρήσεις να φαίνονται φυσιολογικές, καθιστώντας έτσι το έργο του ορισμού της φυσιολογικής συμπεριφοράς πιο δύσκολο.
- Η ακριβής έννοια της ανωμαλίας είναι διαφορετική για διαφορετικούς τομείς εφαρμογών. Για παράδειγμα, στην μετεωρολογία μια απόκλιση από το φυσιολογικό (π.χ. διακυμάνσεις στη θερμοκρασία εκτός εποχής) μπορεί να αποτελεί ανωμαλία, ενώ παρόμοια απόκλιση στον τομέα του χρηματιστηρίου (π.χ. διακυμάνσεις στην αξία μιας μετοχής) μπορεί να θεωρείται φυσιολογική. Συνεπώς, η εφαρμογή μιας τεχνικής που αναπτύχθηκε σε έναν τομέα σε έναν άλλο δεν είναι απλή.
- Συχνά τα δεδομένα περιέχουν θόρυβο που τείνει να είναι παρόμοιος με τις πραγματικές ανωμαλίες και, ως εκ τούτου, είναι δύσκολο να διακριθεί και να αφαιρεθεί.

Λόγω των παραπάνω προκλήσεων, το πρόβλημα της ανίχνευσης ανωμαλιών, στην πιο γενική του μορφή, δεν είναι εύκολο να επιλυθεί. Στην πραγματικότητα, οι περισσότερες από τις υπάρχουσες τεχνικές ανίχνευσης ανωμαλιών επιλύουν μια



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

συγκεκριμένη διατύπωση του προβλήματος. Η διατύπωση αυτή προκαλείται από διάφορους παράγοντες, όπως η φύση των δεδομένων, η διαθεσιμότητα δεδομένων, ο τύπος των ανωμαλιών που πρέπει να ανιχνευθούν κ.λπ. Συχνά, οι παράγοντες αυτοί καθορίζονται από τον τομέα εφαρμογής στον οποίο εντοπίζονται οι ανωμαλίες. Οι ερευνητές έχουν υιοθετήσει έννοιες από διάφορους άλλους τομείς όπως η στατιστική, η μηχανική μάθηση, η εξόρυξη δεδομένων, η θεωρία πληροφοριών, η φασματική θεωρία, και τις έχουν εφαρμόσει σε συγκεκριμένες διατυπώσεις προβλημάτων.

Όπως αναφέρθηκε προηγουμένως, η συγκεκριμένη διατύπωση του προβλήματος καθορίζεται από διάφορους παράγοντες, όπως η φύση των δεδομένων εισόδου, η διαθεσιμότητα (ή μη διαθεσιμότητα) των ετικετών, καθώς και οι περιορισμοί και οι απαιτήσεις που προκαλούνται από τον τομέα εφαρμογής.

4.2 Φύση των δεδομένων εισόδου

Μια βασική πτυχή κάθε τεχνικής ανίχνευσης ανωμαλιών είναι η φύση των δεδομένων εισόδου. Η είσοδος είναι γενικά μια συλλογή από περιπτώσεις δεδομένων (που αναφέρονται επίσης ως αντικείμενο, εγγραφή, σημείο, διάνυσμα, μοτίβο, γεγονός, περίπτωση, δείγμα, παρατήρηση, οντότητα) (Tan et al., 2005). Κάθε περίπτωση δεδομένων μπορεί να περιγραφεί χρησιμοποιώντας ένα σύνολο χαρακτηριστικών (που αναφέρονται επίσης ως μεταβλητή, χαρακτηριστικό, γνώρισμα, πεδίο, διάσταση). Τα χαρακτηριστικά μπορεί να είναι διαφόρων τύπων, όπως δυαδικά, κατηγορικά ή συνεχή. Κάθε περίπτωση δεδομένων μπορεί να αποτελείται από ένα μόνο χαρακτηριστικό (μονομεταβλητό) ή πολλαπλά χαρακτηριστικά (πολυμεταβλητό). Στην περίπτωση πολυμεταβλητών περιπτώσεων δεδομένων, όλα τα χαρακτηριστικά μπορεί να είναι του ίδιου τύπου ή να αποτελούν μείγμα διαφορετικών τύπων δεδομένων.

Η φύση των χαρακτηριστικών καθορίζει τη δυνατότητα εφαρμογής των τεχνικών ανίχνευσης ανωμαλιών. Για παράδειγμα, για τις στατιστικές τεχνικές πρέπει να χρησιμοποιούνται διαφορετικά στατιστικά μοντέλα για συνεχή και κατηγορικά δεδομένα. Ομοίως, για τις τεχνικές που βασίζονται στον πλησιέστερο γείτονα, η φύση των χαρακτηριστικών θα καθορίσει το μέτρο απόστασης που θα χρησιμοποιηθεί. Συχνά, αντί των πραγματικών δεδομένων, η απόσταση ανά ζεύγος μεταξύ των περιπτώσεων μπορεί να παρέχεται με τη μορφή ενός πίνακα απόστασης (ή ομοιότητας).



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

Σε τέτοιες περιπτώσεις, οι τεχνικές που απαιτούν πρωτότυπα δεδομένα δεν είναι εφαρμόσιμες, π.χ. πολλές τεχνικές που βασίζονται στη στατιστική και την ταξινόμηση.

Τα δεδομένα εισόδου μπορούν επίσης να κατηγοριοποιηθούν με βάση τη σχέση που υπάρχει μεταξύ των περιπτώσεων δεδομένων (Tan et al., 2005). Οι περισσότερες από τις υπάρχουσες τεχνικές ανίχνευσης ανωμαλιών ασχολούνται με δεδομένα εγγραφής (ή σημειακά δεδομένα), στα οποία δεν υποτίθεται ότι υπάρχει σχέση μεταξύ των περιπτώσεων δεδομένων.

Γενικά, οι περιπτώσεις δεδομένων μπορούν να συσχετίζονται μεταξύ τους. Ορισμένα παραδείγματα είναι τα δεδομένα ακολουθίας, τα χωρικά δεδομένα και τα δεδομένα γραφημάτων. Στα δεδομένα ακολουθίας, οι περιπτώσεις δεδομένων είναι γραμμικά ταξινομημένες, π.χ. δεδομένα χρονοσειρών, ακολουθίες γονιδιώματος, ακολουθίες πρωτεϊνών. Στα χωρικά δεδομένα, κάθε περίπτωση δεδομένων σχετίζεται με τις γειτονικές της περιπτώσεις, π.χ. δεδομένα κυκλοφορίας οχημάτων, οικολογικά δεδομένα. Όταν τα χωρικά δεδομένα έχουν μια χρονική (διαδοχική) συνιστώσα, αναφέρονται ως χωροχρονικά δεδομένα, π.χ. κλιματικά δεδομένα.

4.3 Τύπος ανωμαλίας

Μια σημαντική πτυχή μιας τεχνικής ανίχνευσης ανωμαλιών είναι η φύση της επιθυμητής ανωμαλίας. Οι ανωμαλίες μπορούν να ταξινομηθούν στις ακόλουθες τρεις κατηγορίες:

4.3.1 Σημειακές ανωμαλίες.

Εάν μια μεμονωμένη περίπτωση δεδομένων μπορεί να θεωρηθεί ανώμαλη σε σχέση με τα υπόλοιπα δεδομένα, τότε η περίπτωση αυτή ονομάζεται σημειακή ανωμαλία. Αυτός είναι ο απλούστερος τύπος ανωμαλίας και αποτελεί το επίκεντρο της πλειονότητας των ερευνών για την ανίχνευση ανωμαλιών.

Για παράδειγμα, στο Σχήμα 3, τα σημεία o_1 και o_2 καθώς και τα σημεία στην περιοχή O_3 βρίσκονται εκτός των ορίων των κανονικών περιοχών και, ως εκ τούτου, αποτελούν σημειακές ανωμαλίες, δεδομένου ότι διαφέρουν από τα κανονικά σημεία δεδομένων.

Ένα άλλο παράδειγμα, αποτελεί η ανίχνευση απάτης με πιστωτικές κάρτες. Έστω ότι το σύνολο δεδομένων αντιστοιχεί στις συναλλαγές μιας πιστωτικής κάρτας ενός



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

ατόμου. Για λόγους απλότητας, θεωρείται ότι τα δεδομένα ορίζονται χρησιμοποιώντας μόνο ένα χαρακτηριστικό, το ποσό που δαπανήθηκε. Μια συναλλαγή για την οποία το ποσό που δαπανήθηκε είναι πολύ υψηλό σε σύγκριση με το κανονικό εύρος δαπανών για το συγκεκριμένο άτομο θα αποτελεί μια σημειακή ανωμαλία.

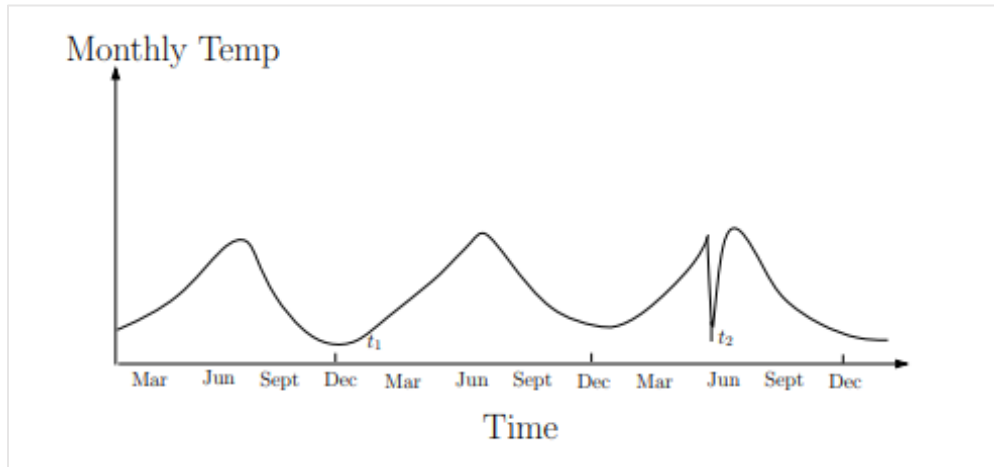
4.3.2 Ανωμαλίες πλαισίου

Εάν μια περίπτωση δεδομένων είναι ανώμαλη σε ένα συγκεκριμένο πλαίσιο (αλλά όχι αλλιώς), τότε ονομάζεται ανωμαλία πλαισίου, αναφέρεται επίσης ως ανωμαλία υπό συνθήκη (Song et al., 2007).

Η έννοια του πλαισίου προκύπτει από τη δομή του συνόλου δεδομένων και πρέπει να προσδιορίζεται ως μέρος της διατύπωσης του προβλήματος. Κάθε περίπτωση δεδομένων ορίζεται χρησιμοποιώντας τα ακόλουθα δύο σύνολα χαρακτηριστικών:

- i. Συναφή χαρακτηριστικά. Τα συναφή χαρακτηριστικά χρησιμοποιούνται για τον προσδιορισμό του πλαισίου (ή της γειτονιάς) για την εν λόγω περίπτωση. Για παράδειγμα, σε σύνολα χωρικών δεδομένων, το γεωγραφικό μήκος και πλάτος μιας τοποθεσίας είναι τα χαρακτηριστικά πλαισίου. Στα δεδομένα χρονοσειρών, ο χρόνος είναι ένα χαρακτηριστικό πλαίσιο που καθορίζει τη θέση μιας περίπτωσης σε ολόκληρη την ακολουθία.
- ii. Χαρακτηριστικά συμπεριφοράς. Τα χαρακτηριστικά συμπεριφοράς καθορίζουν τα μη πλαισιωμένα χαρακτηριστικά μιας περίπτωσης. Για παράδειγμα, σε ένα σύνολο χωρικών δεδομένων που περιγράφει τη μέση βροχόπτωση ολόκληρου του κόσμου, το ποσό της βροχόπτωσης σε κάθε τοποθεσία είναι ένα χαρακτηριστικό συμπεριφοράς.

Η ανώμαλη συμπεριφορά προσδιορίζεται χρησιμοποιώντας τις τιμές των χαρακτηριστικών συμπεριφοράς σε ένα συγκεκριμένο πλαίσιο. Μια περίπτωση δεδομένων μπορεί να είναι μια ανωμαλία σε ένα συγκεκριμένο πλαίσιο, αλλά μια πανομοιότυπη περίπτωση δεδομένων (από την άποψη των χαρακτηριστικών συμπεριφοράς) μπορεί να θεωρηθεί φυσιολογική σε ένα διαφορετικό πλαίσιο. Αυτή η ιδιότητα είναι το κλειδί για τον προσδιορισμό των χαρακτηριστικών πλαισίου και συμπεριφοράς για μια τεχνική ανίχνευσης ανωμαλιών πλαισίου.



Διάγραμμα 14. Συναφής ανωμαλία t_2 σε χρονοσειρά θερμοκρασίας. Η θερμοκρασία τη στιγμή t_1 είναι η ίδια με εκείνη τη στιγμή t_2 , αλλά εμφανίζεται σε διαφορετικό πλαίσιο και ως εκ τούτου δεν θεωρείται ως ανωμαλία

Πηγή: (Chandola et a., 2009)

Οι συγκυριακές ανωμαλίες έχουν διερευνηθεί συχνότερα σε δεδομένα χρονοσειρών και σε χωρικά δεδομένα (Kou et al., 2006). Στο Διάγραμμα 14 παρουσιάζεται ένα τέτοιο παράδειγμα για μια χρονοσειρά θερμοκρασίας που δείχνει τη μηνιαία θερμοκρασία μιας περιοχής τα τελευταία χρόνια. Μια θερμοκρασία 35°C μπορεί να είναι φυσιολογική κατά τη διάρκεια του χειμώνα (τη χρονική στιγμή t_1) σε αυτό το μέρος, αλλά η ίδια τιμή κατά τη διάρκεια του καλοκαιριού (τη χρονική στιγμή t_2) θα αποτελούσε ανωμαλία.

Ένα παρόμοιο παράδειγμα μπορεί να βρεθεί στον τομέα της ανίχνευσης απάτης με πιστωτικές κάρτες. Ένα σχετικό χαρακτηριστικό στον τομέα των πιστωτικών καρτών μπορεί να είναι ο χρόνος αγοράς. Έστω ότι ένα άτομο έχει συνήθως εβδομαδιαίο λογαριασμό αγορών ύψους 100 δολαρίων, εκτός από την εβδομάδα των Χριστουγέννων, όταν αυτός φτάνει τα 1000 δολάρια. Μια νέα αγορά 1000 δολαρίων σε μια εβδομάδα του Ιουλίου θα θεωρηθεί μια πλαισιακή ανωμαλία, καθώς δεν συνάδει με τη συνήθη συμπεριφορά του ατόμου στο πλαίσιο του χρόνου, παρόλο που το ίδιο ποσό που δαπανήθηκε κατά τη διάρκεια της εβδομάδας των Χριστουγέννων θα θεωρηθεί φυσιολογικό.

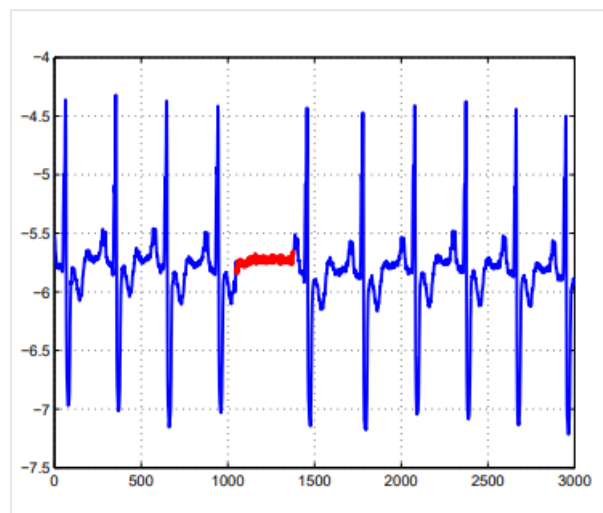


ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

Η επιλογή της εφαρμογής μιας τεχνικής ανίχνευσης ανωμαλιών πλαισίου καθορίζεται από τη σημασία που έχουν οι ανωμαλίες πλαισίου στον τομέα της εφαρμογής-στόχου. Ένας άλλος βασικός παράγοντας είναι η διαθεσιμότητα των χαρακτηριστικών του πλαισίου. Σε αρκετές περιπτώσεις ο ορισμός ενός πλαισίου είναι απλός και, ως εκ τούτου, η εφαρμογή μιας τεχνικής ανίχνευσης ανωμαλιών βάσει πλαισίου έχει νόημα. Σε άλλες περιπτώσεις, ο ορισμός ενός πλαισίου δεν είναι εύκολος, γεγονός που καθιστά δύσκολη την εφαρμογή τέτοιων τεχνικών.

4.3.3 Συλλογικές ανωμαλίες

Εάν μια συλλογή συναφών περιπτώσεων δεδομένων είναι ανώμαλη σε σχέση με το σύνολο των δεδομένων, τότε ονομάζεται συλλογική ανωμαλία. Οι μεμονωμένες περιπτώσεις δεδομένων σε μια συλλογική ανωμαλία μπορεί να μην είναι ανωμαλίες από μόνες τους, αλλά η εμφάνισή τους μαζί ως συλλογή είναι ανωμαλία. Το Διάγραμμα 9 απεικονίζει ένα παράδειγμα το οποίο δείχνει την έξοδο ενός ανθρώπινου ηλεκτροκαρδιογραφήματος (Goldberger et al., 2000). Η επισημασμένη περιοχή υποδηλώνει μια ανωμαλία επειδή η ίδια χαμηλή τιμή υπάρχει για ασυνήθιστα μεγάλο χρονικό διάστημα (που αντιστοιχεί σε μια πρόωρη κοιλιακή συστολή). Αυτή η χαμηλή τιμή από μόνη της δεν αποτελεί ανωμαλία.



Διάγραμμα 15. Συλλογική ανωμαλία που αντιστοιχεί σε μια πρόωρη κοιλιακή συστολή στην έξοδο ενός ανθρώπινου ηλεκτροκαρδιογραφήματος

Πηγή: (Chandola et al., 2009)



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

Οι συλλογικές ανωμαλίες έχουν διερευνηθεί για δεδομένα ακολουθιών, δεδομένα γραφημάτων και χωρικά δεδομένα.

Θα πρέπει να σημειωθεί ότι, ενώ οι σημειακές ανωμαλίες μπορούν να εμφανιστούν σε οποιοδήποτε σύνολο δεδομένων, οι συλλογικές ανωμαλίες μπορούν να εμφανιστούν μόνο σε σύνολα δεδομένων στα οποία οι περιπτώσεις δεδομένων σχετίζονται μεταξύ τους. Αντίθετα, η εμφάνιση ανωμαλιών πλαισίου εξαρτάται από τη διαθεσιμότητα χαρακτηριστικών πλαισίου στα δεδομένα.

Κεφάλαιο 5ο : Μέθοδοι ανίχνευσης ανωμαλιών

5.1 Τεχνικές ανίχνευσης ανωμαλιών με βάση την ταξινόμηση

Η ταξινόμηση (Tan et al. 2005) χρησιμοποιείται για την εκμάθηση ενός μοντέλου (ταξινομητή) από ένα σύνολο χαρακτηρισμένων περιπτώσεων δεδομένων. Στη συνέχεια, για την ταξινόμηση μιας περίπτωσης δοκιμής σε μία από τις κλάσεις χρησιμοποιείται το μοντέλο ταξινομητής που προέκυψε, από την προηγούμενη διαδικασία. Οι τεχνικές ανίχνευσης ανωμαλιών που βασίζονται στην ταξινόμηση λειτουργούν με παρόμοιο τρόπο δύο φάσεων. Η φάση εκπαίδευσης περιλαμβάνει τη δημιουργία ενός ταξινομητή (μοντέλου) χρησιμοποιώντας τα διαθέσιμα χαρακτηρισμένα δεδομένα εκπαίδευσης. Η φάση δοκιμής ταξινομεί μια περίπτωση δοκιμής ως κανονική ή ανώμαλη χρησιμοποιώντας τον ταξινομητή.

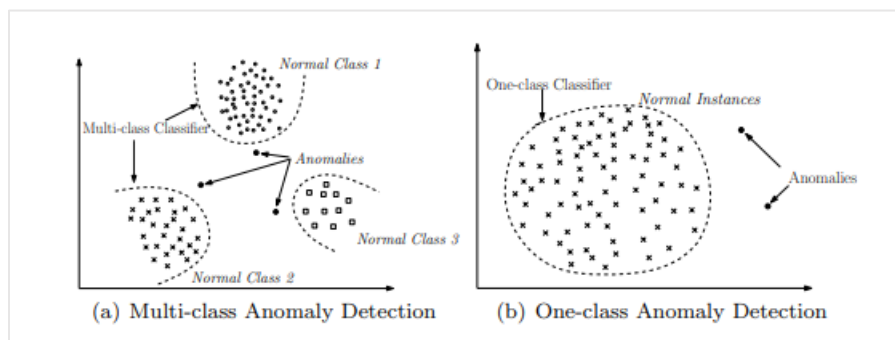
Οι τεχνικές ανίχνευσης ανωμαλιών που βασίζονται στην ταξινόμηση λειτουργούν με την ακόλουθη γενική παραδοχή:

Υπόθεση: Ένας ταξινομητής που μπορεί να διακρίνει μεταξύ φυσιολογικών και ανώμαλων κλάσεων μπορεί να προκύψει σε ένα συγκεκριμένο χώρο χαρακτηριστικών.

Σε αυτό το στάδιο, οι τεχνικές ανίχνευσης ανωμαλιών που βασίζονται στην ταξινόμηση μπορούν να ομαδοποιηθούν σε δύο ομάδες, στις τεχνικές ανίχνευσης ανωμαλιών πολλαπλών κατηγοριών και στις τεχνικές ανίχνευσης μιας κατηγορίας.

Οι τεχνικές ανίχνευσης ανωμαλιών που βασίζονται σε ταξινόμηση πολλαπλών κλάσεων υποθέτουν ότι τα δεδομένα εκπαίδευσης περιέχουν επισημασμένες

περιπτώσεις που ανήκουν σε πολλαπλές κανονικές κλάσεις (Stefano et al., 2000). Τέτοιες τεχνικές ανίχνευσης ανωμαλιών χρησιμοποιούν έναν ταξινομητή για να διακρίνουν κάθε κανονική κλάση από τις υπόλοιπες κλάσεις. Μια περίπτωση δοκιμής θεωρείται ανώμαλη εάν δεν ταξινομείται ως κανονική από κανέναν από τους ταξινομητές. Ορισμένες τεχνικές αυτής της υποκατηγορίας συσχετίζουν ένα σκορ εμπιστοσύνης με την πρόβλεψη που κάνει ο ταξινομητής. Εάν κανένας από τους ταξινομητές δεν είναι σίγουρος ότι θα ταξινομήσει την περίπτωση δοκιμής ως κανονική, η περίπτωση θεωρείται ανώμαλη.



Διάγραμμα 16: Χρήση ταξινόμησης για ανίχνευση ανωμαλιών

Πηγή: (Chandola et al., 2009)

Οι τεχνικές ανίχνευσης ανωμαλιών που βασίζονται στην ταξινόμηση μίας κατηγορίας υποθέτουν ότι όλα τα παραδείγματα εκπαίδευσης έχουν μόνο μία ετικέτα κλάσης (ομαλά ή ανώμαλα). Τέτοιες τεχνικές μαθαίνουν ένα διακριτικό όριο γύρω από τις κανονικές περιπτώσεις χρησιμοποιώντας έναν αλγόριθμο ταξινόμησης μιας κατηγορίας (Scholkopf et al., 2001). Κάθε περίπτωση δοκιμής που δεν εμπίπτει στο όριο που μαθαίνεται δηλώνεται ως ανώμαλη.

Στις υποενότητες που ακολουθούν, παρουσιάζονται μια ποικιλία τεχνικών ανίχνευσης ανωμαλιών που χρησιμοποιούν διαφορετικούς αλγορίθμους ταξινόμησης για τη δημιουργία ταξινομητών.

5.1.1 Νευρωνικά δίκτυα

Τα νευρωνικά δίκτυα, επίσης γνωστά ως τεχνητά νευρωνικά δίκτυα (ANN) ή προσομοιωμένα νευρωνικά δίκτυα (SNN), αποτελούν ένα υποσύνολο της μηχανικής μάθησης και βρίσκονται στο επίκεντρο των αλγορίθμων βαθιάς μάθησης. Το όνομα

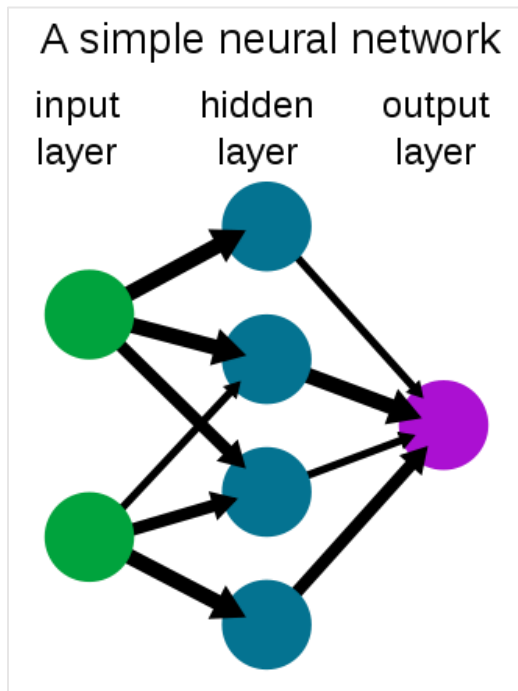


ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

και η δομή τους είναι εμπνευσμένα από τον ανθρώπινο εγκέφαλο, μιμούμενα τον τρόπο με τον οποίο οι βιολογικοί νευρώνες στέλνουν σήματα ο ένας στον άλλο.

Τα τεχνητά νευρωνικά δίκτυα (ΤΝΔ) αποτελούνται από στρώματα κόμβων, που περιέχουν ένα στρώμα εισόδου, ένα ή περισσότερα κρυφά στρώματα και ένα στρώμα εξόδου. Κάθε κόμβος, ή τεχνητός νευρώνας, συνδέεται με έναν άλλο και έχει ένα σχετικό βάρος και κατώφλι. Εάν η έξοδος οποιουδήποτε μεμονωμένου κόμβου είναι πάνω από την καθορισμένη τιμή κατωφλίου, ο εν λόγω κόμβος ενεργοποιείται, στέλνοντας δεδομένα στο επόμενο στρώμα του δικτύου. Διαφορετικά, δεν διαβιβάζονται δεδομένα στο επόμενο επίπεδο του δικτύου.

Τα νευρωνικά δίκτυα βασίζονται σε δεδομένα εκπαίδευσης για να μαθαίνουν και να βελτιώνουν την ακρίβειά τους με την πάροδο του χρόνου. Ωστόσο, μόλις αυτοί οι αλγόριθμοι μάθησης ρυθμιστούν με ακρίβεια, αποτελούν ισχυρά εργαλεία στην επιστήμη των υπολογιστών και την τεχνητή νοημοσύνη, επιτρέποντάς μας να ταξινομήσουμε και να ομαδοποιήσουμε δεδομένα με μεγάλη ταχύτητα. Οι εργασίες στην αναγνώριση ομιλίας ή στην αναγνώριση εικόνας μπορούν να διαρκέσουν λίγα λεπτά έναντι ωρών σε σύγκριση με τη χειροκίνητη αναγνώριση από ανθρώπινους εμπειρογνώμονες. Ένα από τα πιο γνωστά νευρωνικά δίκτυα είναι ο αλγόριθμος αναζήτησης της Google.



Σχήμα 2: Απλουστευμένη άποψη ενός τεχνητού νευρωνικού δικτύου τροφοδότησης

Τα νευρωνικά δίκτυα έχουν εφαρμοστεί στην ανίχνευση ανωμαλιών σε περιβάλλον πολλαπλών τάξεων καθώς και σε περιβάλλον μίας τάξης.

Πίνακας 1. Ορισμένα παραδείγματα τεχνικών ανίχνευσης ανωμαλιών με βάση την ταξινόμηση με χρήση νευρωνικών δικτύων

Χρήση Νευρωνικών Δικτύων	Αναφορές
Πολυεπίπεδα αντιληπτικά συστήματα	(Augusteijn & Folkert 2002; Cun et al., 1990)
Νευρωνικά δέντρα	(Martinez, 1998)
Αυτο-συνδεδετικά δίκτυα	(Byungho & Sungzoon, 1999)
Προσαρμοστικός συντονισμός με βάση τη θεωρία	(Moya et al., 1993)
Ακτινική βάση, βάσει συνάρτησης	(Albrecht et al., 2000; Bishop, 1994)
Δίκτυα Hopfield	(Crook & Hayes, 2001; Crook et al., 2002)
Δίκτυα ταλάντωσης	(Ho & Rouat, 1998; Kojima & Ito, 1999)

Μια βασική τεχνική ανίχνευσης ανωμαλιών πολλαπλών κατηγοριών που χρησιμοποιεί νευρωνικά δίκτυα λειτουργεί σε δύο βήματα. Πρώτον, ένα νευρωνικό δίκτυο



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

εκπαιδεύεται στα κανονικά δεδομένα εκπαίδευσης για να μάθει τις διάφορες κανονικές κλάσεις. Δεύτερον, κάθε περίπτωση δοκιμής παρέχεται ως είσοδος στο νευρωνικό δίκτυο. Εάν το δίκτυο αποδεχθεί την είσοδο δοκιμής, πρόκειται για κανονική και εάν το δίκτυο απορρίψει μια είσοδο δοκιμής, πρόκειται για ανωμαλία (Stefano et al., 2000). Έχουν προταθεί διάφορες παραλλαγές της βασικής τεχνικής των νευρωνικών δικτύων που χρησιμοποιούν διαφορετικούς τύπους νευρωνικών δικτύων, όπως συνοψίζονται στον Πίνακα 1.

Τα νευρωνικά δίκτυα αντιγραφής έχουν χρησιμοποιηθεί για την ανίχνευση ανωμαλιών μίας κατηγορίας (Hawkins et al., 2002). Κατασκευάζεται ένα νευρωνικό δίκτυο πολλαπλών επιπέδων που έχει τον ίδιο αριθμό νευρώνων εισόδου και εξόδου (που αντιστοιχούν στα χαρακτηριστικά των δεδομένων). Η εκπαίδευση περιλαμβάνει τη συμπίεση των δεδομένων σε τρία κρυφά στρώματα. Η φάση δοκιμής περιλαμβάνει την ανακατασκευή κάθε περίπτωσης δεδομένων x_i χρησιμοποιώντας το δίκτυο που μαθαίνεται για να ληφθεί η ανακατασκευασμένη έξοδος o_i . Το σφάλμα ανακατασκευής δ_i για την περίπτωση δοκιμής x_i υπολογίζεται στη συνέχεια ως εξής:

$$\delta_i = \frac{1}{n} \sum_{j=1}^n (x_{ij} - o_{ij})^2 \quad (1.12)$$

όπου n είναι ο αριθμός των χαρακτηριστικών στα οποία ορίζονται τα δεδομένα. Το σφάλμα ανακατασκευής δ_i χρησιμοποιείται άμεσα ως βαθμολογία ανωμαλίας για την περίπτωση δοκιμής.

5.1.2 Δίκτυα Bayesian

Ένα δίκτυο Bayes είναι ένα πιθανοτικό γραφικό μοντέλο (ένα είδος στατιστικού μοντέλου), το οποίο αντιπροσωπεύει ένα σύνολο τυχαίων μεταβλητών, όπως και τις εξαρτήσεις τους, μέσω ενός άκυκλου κατευθυνόμενου γράφου. Σε αυτούς τους άκυκλους κατευθυνόμενους γράφους, οι κόμβοι αποτελούν τυχαίες μεταβλητές στην Bayesian λογική: μπορεί να είναι παρατηρήσιμες ποσότητες, λανθάνουσες μεταβλητές, άγνωστες παράμετροι ή υποθέσεις. Οι ακμές αντιπροσωπεύουν τις εξαρτήσεις, ενώ οι κόμβοι οι οποίοι δεν είναι συνδεδεμένοι αντιπροσωπεύουν



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

μεταβλητές οι οποίες είναι δυναμικά ανεξάρτητες μεταξύ τους. Κάθε κόμβος συνδέεται με μια συνάρτηση πιθανότητας η οποία λαμβάνει ως είσοδο ένα συγκεκριμένο σύνολο τιμών για τις μεταβλητές του γονέα του κόμβου και δίνει την πιθανότητα της μεταβλητής η οποία αντιπροσωπεύεται από τον κόμβο.

Μια βασική τεχνική που χρησιμοποιεί ένα δίκτυο Bayes για ένα μονοπαραγοντικό σετ κατηγορικών δεδομένων, υπολογίζει την μεταγενέστερη πιθανότητα παρατήρησης μιας κλάσης (από ένα σύνολο κανονικών κλάσεων και την ανώμαλη κλάση), εφόσον έχει πάρει προηγουμένως ένα στιγμιότυπο από τα δεδομένα δοκιμής. Η κλάση με τη μεγαλύτερη μεταγενέστερη πιθανότητα επιλέγεται ως η προβλεπόμενη για το δεδομένο στιγμιότυπο δοκιμής. Η πιθανότητα παρατήρησης του στιγμιότυπου, δεδομένης μιας κλάσης και οι προηγούμενες πιθανότητες της κλάσης, υπολογίζονται από το σύνολο δεδομένων εκπαίδευσης.

Η βασική τεχνική μπορεί να γενικευτεί σε πολυμεταβλητό σύνολο κατηγορικών δεδομένων με τη συγκέντρωση των μεταγενέστερων πιθανοτήτων ανά χαρακτηριστικό για κάθε περίπτωση δοκιμής και τη χρήση της συγκεντρωτικής τιμής για την ανάθεση μιας ετικέτας κλάσης στην περίπτωση δοκιμής. Διάφορες παραλλαγές της βασικής τεχνικής έχουν προταθεί για την ανίχνευση εισβολών σε δίκτυα (Sebyala et al., 2002), για την ανίχνευση καινοτομιών στην παρακολούθηση βίντεο (Diehl & Hampshire, 2002), για την ανίχνευση ανωμαλιών σε δεδομένα κειμένου (Baker et al., 1999) και για την ανίχνευση επιδημιών (Wong et al., 2003).

Η βασική τεχνική που περιγράφεται παραπάνω προϋποθέτει ανεξαρτησία μεταξύ των διαφόρων χαρακτηριστικών. Έχουν προταθεί διάφορες παραλλαγές της βασικής τεχνικής που αποτυπώνουν τις εξαρτήσεις υπό όρους μεταξύ των διαφόρων χαρακτηριστικών χρησιμοποιώντας πιο σύνθετα δίκτυα Bayes (Das & Schneider, 2007).

5.1.3 Μηχανές διανυσμάτων υποστήριξης

Ο βασικός SVM, αφού πάρει ένα σύνολο δεδομένων εισόδου, προβλέπει, για κάθε είσοδο, ποιά από τις δύο δυνατές τάξεις αποτελεί την έξοδο, καθιστώντας τον ως ένα γραμμικό δυαδικό ταξινομητή. Λαμβάνοντας υπόψη ένα σύνολο στιγμιότυπων εκπαίδευσης, κάθε ένα από τα οποία ανήκει σε μία από τις δύο κατηγορίες, ένα SVM



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

χτίζει ένα μοντέλο το οποίο εκχωρεί νέα παραδείγματα στη μία κατηγορία ή την άλλη. Ένα μοντέλο SVM αποτελεί μια αναπαράσταση των παραδειγμάτων ως σημεία στο χώρο, τα οποία χαρτογραφούνται, έτσι ώστε τα παραδείγματα των ξεχωριστών κατηγοριών να χωρίζονται από ένα σαφές κενό το οποίο είναι όσο το δυνατόν ευρύτερο. Τα νέα παραδείγματα στη συνέχεια χαρτογραφούνται στον ίδιο χώρο. Επιπλέον, προβλέπεται αν ανήκουν σε μια κατηγορία με βάση την πλευρά του διακένου στην οποία «πέφτουν».

Παραλλαγές της βασικής τεχνικής έχουν προταθεί για την ανίχνευση ανωμαλιών σε δεδομένα ηχητικών σημάτων (Davy and Godsill, 2002), την ανίχνευση καινοτομιών σε σταθμούς παραγωγής ενέργειας (King et al., 2002) και την ανίχνευση εισβολών σε κλήσεις συστήματος (Heller et al., 2003). Η βασική τεχνική έχει επίσης επεκταθεί για την ανίχνευση ανωμαλιών σε χρονικές ακολουθίες (Ma & Perkins, 2003).

5.1.4 Βασισμένες σε κανόνες

Οι βασισμένες σε κανόνες τεχνικές ανίχνευσης ανωμαλιών, μαθαίνουν κανόνες που απεικονίζουν την κανονική συμπεριφορά ενός συστήματος. Κάποιο στιγμιότυπο το οποίο δεν καλύπτεται από τους εν λόγω κανόνες θεωρείται ανωμαλία. Οι τεχνικές αυτές έχουν εφαρμοστεί σε multi-class καθώς και σε one-class περιπτώσεις. Στην περίπτωση multi-class η τεχνική αποτελείται από δύο βήματα. Στο πρώτο βήμα συντελείται η μάθηση των κανόνων από τα δεδομένα εκπαίδευσης χρησιμοποιώντας έναν αλγόριθμο μάθησης όπως ο RIPPER, τα δέντρα αποφάσεων, κλπ. Κάθε κανόνας έχει μια σχετική τιμή εμπιστοσύνης η οποία είναι ανάλογη με το λόγο του αριθμού των περιπτώσεων εκπαίδευσης οι οποίες ταξινομούνται ορθώς από τον κανόνα, και του συνολικού αριθμού περιπτώσεων εκπαίδευσης. Στο δεύτερο βήμα επιχειρείται η εύρεση, για κάθε παράδειγμα δοκιμής, του κανόνα ο οποίος συλλαμβάνει καλύτερα το παράδειγμα δοκιμής, έτσι ώστε το αντίστροφο της εμπιστοσύνης η οποία συνδέεται με τον καλύτερο κανόνα να είναι το σκορ ανωμαλίας του παραδείγματος δοκιμής

Τεχνικές βασισμένες στην εξόρυξη κανόνων συσχέτισης έχουν χρησιμοποιηθεί για την ανίχνευση εισβολών σε δίκτυα (Mahoney & Chan, 2002), την ανίχνευση εισβολών σε κλήσεις συστήματος (Lee et al., 2000), την ανίχνευση απάτης με πιστωτικές κάρτες



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

(Brause et al., 1999) και την ανίχνευση απάτης σε δεδομένα διαστημικών σκαφών (Yairi et al., 2001).

5.1.5 Υπολογιστική πολυπλοκότητα

Η υπολογιστική πολυπλοκότητα των τεχνικών που βασίζονται στην ταξινόμηση εξαρτάται από τον αλγόριθμο ταξινόμησης που χρησιμοποιείται (Kearns, 1990). Γενικά, η εκπαίδευση των δέντρων απόφασης τείνει να είναι ταχύτερη, ενώ οι τεχνικές που περιλαμβάνουν τετραγωνική βελτιστοποίηση, όπως οι SVM, είναι πιο κοστοβόρες, αν και έχουν προταθεί SVM γραμμικού χρόνου που έχουν γραμμικό χρόνο εκπαίδευσης. Η φάση δοκιμής των τεχνικών ταξινόμησης είναι συνήθως πολύ γρήγορη, δεδομένου ότι η φάση δοκιμής χρησιμοποιεί ένα γνωστό μοντέλο για την ταξινόμηση. Τα πλεονεκτήματα των τεχνικών που βασίζονται στην ταξινόμηση είναι τα εξής:

- Οι τεχνικές που βασίζονται στην ταξινόμηση, ιδίως οι τεχνικές πολλαπλών κατηγοριών, μπορούν να κάνουν χρήση ισχυρών αλγορίθμων για να διακρίνουν ανωμαλίες μεταξύ περιπτώσεων που ανήκουν σε διαφορετικές κλάσεις.
- Η φάση δοκιμής των τεχνικών που βασίζονται στην ταξινόμηση είναι γρήγορη, δεδομένου ότι κάθε περίπτωση δοκιμής πρέπει να συγκρίνεται με το προ-υπολογισμένο μοντέλο.

Τα μειονεκτήματα των τεχνικών που βασίζονται στην ταξινόμηση είναι τα εξής:

- Οι τεχνικές βασισμένες σε ταξινόμηση πολλαπλών κλάσεων βασίζονται στη διαθεσιμότητα ακριβών ετικετών για διάφορες κανονικές κλάσεις, πράγμα που συχνά δεν είναι δυνατό.
- Οι τεχνικές που βασίζονται στην ταξινόμηση αποδίδουν μια ετικέτα σε κάθε περίπτωση δοκιμής, γεγονός που μπορεί επίσης να αποτελέσει μειονέκτημα, όταν είναι επιθυμητό ένα ουσιαστικό αποτέλεσμα ανωμαλίας για τις περιπτώσεις δοκιμής. Ορισμένες τεχνικές ταξινόμησης που λαμβάνουν μια πιθανολογική πρόβλεψη βαθμολόγησης από την έξοδο ενός ταξινομητή, μπορούν να χρησιμοποιηθούν για την αντιμετώπιση αυτού του προβλήματος (Platt, 2000).



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

5.2 Τεχνικές ανίχνευσης ανωμαλιών με βάση τον πλησιέστερο γείτονα

Η έννοια της ανάλυσης του πλησιέστερου γείτονα έχει χρησιμοποιηθεί σε διάφορες τεχνικές ανίχνευσης ανωμαλιών. Οι τεχνικές αυτές βασίζονται στην ακόλουθη βασική υπόθεση:

Υπόθεση: οι κανονικές περιπτώσεις δεδομένων εμφανίζονται σε πυκνές γειτονιές, ενώ οι ανωμαλίες εμφανίζονται μακριά από τους πλησιέστερους γείτονές τους.

Οι τεχνικές ανίχνευσης ανωμαλιών που βασίζονται στον πλησιέστερο γείτονα απαιτούν ένα μέτρο απόστασης ή ομοιότητας που ορίζεται μεταξύ δύο περιπτώσεων δεδομένων. Η απόσταση (ή η ομοιότητα) μεταξύ δύο περιπτώσεων δεδομένων μπορεί να υπολογιστεί με διάφορους τρόπους. Για συνεχή χαρακτηριστικά, η ευκλείδεια απόσταση είναι μια δημοφιλής επιλογή, αλλά μπορούν να χρησιμοποιηθούν και άλλα μέτρα (Tan et al., 2005). Για κατηγορικά χαρακτηριστικά, χρησιμοποιείται συχνά ένας απλός συντελεστής αντιστοιχίας, αλλά μπορούν να χρησιμοποιηθούν και πιο σύνθετα μέτρα απόστασης (Chandola et al., 2008). Για πολυμεταβλητές περιπτώσεις δεδομένων, η απόσταση ή η ομοιότητα υπολογίζεται συνήθως για κάθε χαρακτηριστικό και στη συνέχεια συνδυάζεται.

Οι περισσότερες από τις τεχνικές που θα συζητηθούν στην παρούσα ενότητα, δεν απαιτούν το μέτρο απόστασης να είναι αυστηρά μετρικό. Τα μέτρα απαιτείται συνήθως να είναι θετικά ορισμένα και συμμετρικά, αλλά δεν απαιτείται να ικανοποιούν την τριγωνική ανισότητα. Οι τεχνικές ανίχνευσης ανωμαλιών που βασίζονται στον πλησιέστερο γείτονα μπορούν να ομαδοποιηθούν σε γενικές γραμμές σε δύο κατηγορίες:

1. Τεχνικές που χρησιμοποιούν την απόσταση μιας περίπτωσης δεδομένων από τον k^{th} πλησιέστερο γείτονά της ως βαθμίδα ανωμαλίας.
2. Τεχνικές που υπολογίζουν τη σχετική πυκνότητα κάθε περίπτωσης δεδομένων για τον υπολογισμό της βαθμολογίας ανωμαλίας.

Επιπλέον, υπάρχουν ορισμένες τεχνικές που χρησιμοποιούν την απόσταση μεταξύ των περιπτώσεων δεδομένων με διαφορετικό τρόπο για τον εντοπισμό ανωμαλιών και θα συζητηθούν εν συντομία στη συνέχεια.



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

5.2.1 Χρήση της απόστασης από τον k^{th} πλησιέστερο γείτονα

Μια βασική τεχνική ανίχνευσης ανωμαλιών από τον πλησιέστερο γείτονα βασίζεται στον ακόλουθο ορισμό. Ο βαθμός ανωμαλίας μιας περίπτωσης δεδομένων ορίζεται ως η απόστασή της από τον k^{th} πλησιέστερο γείτονά της σε ένα δεδομένο σύνολο δεδομένων. Αυτή η βασική τεχνική έχει εφαρμοστεί για την ανίχνευση ναρκών από δορυφορικές εικόνες εδάφους (Byers & Raftery, 1998) και για την ανίχνευση βραχυκυκλωμένων στροφών (ανωμαλιών) στο πεδίο συνεχούς ρεύματος μεγάλων σύγχρονων στροβιλογεννητριών (Guttormsson et al., 1999). Συνήθως, ένα κατώφλι εφαρμόζεται στη συνέχεια στο σκορ ανωμαλίας για να καθοριστεί αν μια περίπτωση δοκιμής είναι ανώμαλη ή όχι. Οι Ramaswamy et al. (2000), από την άλλη πλευρά, επιλέγουν n περιπτώσεις με τα μεγαλύτερα σκορ ανωμαλίας ως ανωμαλίες.

Η βασική τεχνική έχει επεκταθεί από τους ερευνητές με τρεις διαφορετικούς τρόπους. Το πρώτο σύνολο παραλλαγών τροποποιεί τον παραπάνω ορισμό για να ληφθεί η βαθμολογία ανωμαλίας μιας περίπτωσης δεδομένων. Το δεύτερο σύνολο παραλλαγών χρησιμοποιεί διαφορετικά μέτρα απόστασης/ομοιότητας για να χειριστεί διαφορετικούς τύπους δεδομένων. Το τρίτο σύνολο παραλλαγών επικεντρώνεται στη βελτίωση της αποτελεσματικότητας της βασικής τεχνικής με διαφορετικούς τρόπους.

Οι Eskin et al. (2002) υπολογίζουν το σκορ ανωμαλίας μιας περίπτωσης δεδομένων ως το άθροισμα των αποστάσεών της από τους k πλησιέστερους γείτονές της. Μια παρόμοια τεχνική έχει εφαρμοστεί για την ανίχνευση απάτης με πιστωτικές κάρτες από τους Bolton & Hand (1999) με την ονομασία Peer Group Analysis.

Ένας διαφορετικός τρόπος για τον υπολογισμό της βαθμολογίας ανωμαλίας μιας περίπτωσης δεδομένων είναι η καταμέτρηση του αριθμού των πλησιέστερων γειτόνων (n) που δεν απέχουν περισσότερο από d απόσταση από τη συγκεκριμένη περίπτωση δεδομένων (Knoig et al., 2000). Αυτή η μέθοδος μπορεί επίσης να θεωρηθεί ως εκτίμηση της ολικής πυκνότητας πιθανότητας για κάθε περίπτωση δεδομένων, δεδομένου ότι περιλαμβάνει την καταμέτρηση του αριθμού των γειτόνων σε μια υπερσφαίρα ακτίνας d . Για παράδειγμα, σε ένα σύνολο δεδομένων 2-D, η πυκνότητα μιας περίπτωσης δεδομένων $= \frac{n}{\pi d^2}$. Το αντίστροφο της πυκνότητας είναι η βαθμολογία ανωμαλίας για την περίπτωση δεδομένων. Αντί του υπολογισμού της πραγματικής



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

πυκνότητας, αρκετές τεχνικές καθορίζουν την ακτίνα d και χρησιμοποιούν το d ως βαθμολογία ανωμαλίας, ενώ αρκετές τεχνικές καθορίζουν το n και χρησιμοποιούν το 1 ως βαθμολογία ανωμαλίας.

Έχουν προταθεί διάφορες παραλλαγές της βασικής τεχνικής για τη βελτίωση της αποτελεσματικότητας. Ορισμένες τεχνικές συρρικνώνουν το χώρο αναζήτησης είτε αγνοώντας περιπτώσεις που δεν μπορούν να είναι ανώμαλες είτε εστιάζοντας σε περιπτώσεις που είναι πιο πιθανό να είναι ανώμαλες. Μετά τον υπολογισμό των πλησιέστερων γειτόνων για μια περίπτωση δεδομένων, ο αλγόριθμος θέτει το κατώφλι ανωμαλίας για κάθε περίπτωση δεδομένων στη βαθμολογία της πιο αδύναμης ανωμαλίας που έχει βρεθεί μέχρι στιγμής. Χρησιμοποιώντας αυτή τη διαδικασία κλαδέματος, η τεχνική απορρίπτει περιπτώσεις που είναι κοντά και, ως εκ τούτου, δεν παρουσιάζουν ενδιαφέρον.

5.2.2 Χρήση της σχετικής πυκνότητας

Οι τεχνικές ανίχνευσης ανωμαλιών με βάση την πυκνότητα εκτιμούν την πυκνότητα της γειτονιάς κάθε περίπτωσης δεδομένων. Μια περίπτωση που βρίσκεται σε μια γειτονιά με χαμηλή πυκνότητα δηλώνεται ως ανώμαλη, ενώ μια περίπτωση που βρίσκεται σε μια πυκνή γειτονιά δηλώνεται ως κανονική.

Για μια δεδομένη περίπτωση δεδομένων, η απόσταση από τον k^{th} πλησιέστερο γείτονα της ισοδυναμεί με την ακτίνα μιας υπερσφαίρας, με κέντρο τη δεδομένη περίπτωση δεδομένων, η οποία περιέχει k άλλες περιπτώσεις. Έτσι, η απόσταση από τον k^{th} πλησιέστερο γείτονα για μια δεδομένη περίπτωση δεδομένων μπορεί να θεωρηθεί ως μια εκτίμηση του αντιστρόφου της πυκνότητας της περίπτωσης στο σύνολο δεδομένων και η βασική τεχνική με βάση τον πλησιέστερο γείτονα που περιγράφηκε στο προηγούμενο υποκεφάλαιο μπορεί να θεωρηθεί ως τεχνική ανίχνευσης ανωμαλιών με βάση την πυκνότητα.

Οι τεχνικές που βασίζονται στην πυκνότητα έχουν κακή απόδοση εάν τα δεδομένα έχουν περιοχές με διαφορετικές πυκνότητες. Για παράδειγμα, έστω ένα σύνολο δεδομένων 2 διαστάσεων που παρουσιάζεται στο Σχήμα 3. Λόγω της χαμηλής πυκνότητας της συστάδας C_1 είναι προφανές ότι για κάθε περίπτωση q εντός της συστάδας C_1 , η απόσταση μεταξύ της περίπτωσης q και του πλησιέστερου γείτονα της

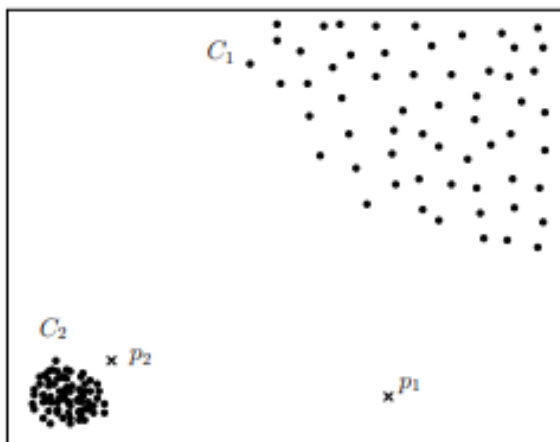


ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

είναι μεγαλύτερη από την απόσταση μεταξύ της περίπτωσης p_2 και του πλησιέστερου γείτονα από τη συστάδα C_2 , και η περίπτωση p_2 δεν θα θεωρηθεί ως ανωμαλία. Ως εκ τούτου, η βασική τεχνική θα αποτύχει να διακρίνει μεταξύ του p_2 και των περιπτώσεων στο C_1 . Ωστόσο, η περίπτωση p_1 μπορεί να εντοπιστεί.

Για να αντιμετωπιστεί το ζήτημα των διαφορετικών πυκνοτήτων στο σύνολο δεδομένων, έχει προταθεί ένα σύνολο τεχνικών για τον υπολογισμό της πυκνότητας των περιπτώσεων σε σχέση με την πυκνότητα των γειτόνων τους.

Οι Breunig et al. (2000) αποδίδουν μια βαθμολογία ανωμαλίας σε μια δεδομένη περίπτωση δεδομένων, γνωστή ως Local Outlier Factor (LOF). Για κάθε δεδομένη περίπτωση δεδομένων, η βαθμολογία LOF ισούται με τον λόγο της μέσης τοπικής πυκνότητας των k πλησιέστερων γειτόνων της περίπτωσης και της τοπικής πυκνότητας της ίδιας της περίπτωσης δεδομένων. Για να βρουν την τοπική πυκνότητα για μια περίπτωση δεδομένων, οι συγγραφείς βρίσκουν πρώτα την ακτίνα της μικρότερης υπερσφαιράς με κέντρο την περίπτωση δεδομένων, η οποία περιέχει τους k πλησιέστερους γείτονές της. Στη συνέχεια, η τοπική πυκνότητα υπολογίζεται διαιρώντας το k με τον όγκο αυτής της υπερσφαιράς. Για μια κανονική περίπτωση που βρίσκεται σε μια πυκνή περιοχή, η τοπική πυκνότητά της θα είναι παρόμοια με εκείνη των γειτόνων της, ενώ για μια ανώμαλη περίπτωση, η τοπική πυκνότητά της θα είναι χαμηλότερη από εκείνη των πλησιέστερων γειτόνων. Ως εκ τούτου, η ανώμαλη περίπτωση θα λάβει υψηλότερη βαθμολογία LOF.



Σχήμα 3: Πλεονέκτημα της τεχνολογίας με βάση την τοπική πυκνότητα



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

Πηγή: (Chandola et a., 2009)

5.2.3 Υπολογιστική πολυπλοκότητα

Ένα μειονέκτημα της βασικής τεχνικής που βασίζεται στον πλησιέστερο γείτονα και της τεχνικής LOF είναι η απαιτούμενη πολυπλοκότητα $O(N^2)$. Δεδομένου ότι αυτές οι τεχνικές περιλαμβάνουν την εύρεση των πλησιέστερων γειτόνων για κάθε περίπτωση, μπορούν να χρησιμοποιηθούν αποδοτικές δομές δεδομένων, όπως τα δέντρα k-d (Bentley, 1975) και τα δέντρα R (Roussopoulos et al., 1995). Αλλά αυτές οι τεχνικές δεν κλιμακώνονται καλά καθώς αυξάνεται ο αριθμός των χαρακτηριστικών. Αρκετές τεχνικές έχουν βελτιστοποιήσει άμεσα την τεχνική ανίχνευσης ανωμαλιών με την υπόθεση ότι μόνο οι λίγες κορυφαίες ανωμαλίες είναι ενδιαφέρουσες. Εάν απαιτείται ένα σκορ ανωμαλίας για κάθε περίπτωση δοκιμής, αυτές οι τεχνικές δεν είναι εφαρμόσιμες. Οι τεχνικές που χωρίζουν το χώρο των χαρακτηριστικών σε ένα υπερ-πλέγμα, είναι γραμμικές ως προς το μέγεθος των δεδομένων, αλλά εκθετικές ως προς τον αριθμό των χαρακτηριστικών και, ως εκ τούτου, δεν είναι κατάλληλες για μεγάλο αριθμό χαρακτηριστικών. Οι τεχνικές δειγματοληψίας προσπαθούν να αντιμετωπίσουν το πρόβλημα πολυπλοκότητας $O(N^2)$ προσδιορίζοντας τους πλησιέστερους γείτονες σε ένα μικρό δείγμα του συνόλου δεδομένων. Όμως, η δειγματοληψία μπορεί να οδηγήσει σε λανθασμένες βαθμολογίες ανωμαλίας εάν το μέγεθος του δείγματος είναι πολύ μικρό.

Τα πλεονεκτήματα των τεχνικών που βασίζονται στον πλησιέστερο γείτονα είναι τα εξής:

- Ένα βασικό πλεονέκτημα των τεχνικών που βασίζονται στον πλησιέστερο γείτονα είναι ότι είναι μη εποπτευόμενες από τη φύση τους και δεν κάνουν καμία υπόθεση σχετικά με τη γενεσιουργό κατανομή των δεδομένων. Αντίθετα, είναι καθαρά καθοδηγούμενες από τα δεδομένα.
- Η προσαρμογή των τεχνικών που βασίζονται στον πλησιέστερο γείτονα σε διαφορετικό τύπο δεδομένων είναι απλή και απαιτεί κυρίως τον ορισμό ενός κατάλληλου μέτρου απόστασης για τα συγκεκριμένα δεδομένα.

Τα μειονεκτήματα των τεχνικών που βασίζονται στον πλησιέστερο γείτονα είναι τα εξής:



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

- Για τις τεχνικές χωρίς επίβλεψη, εάν τα δεδομένα έχουν κανονικές περιπτώσεις που δεν έχουν αρκετούς στενούς γείτονες ή αν τα δεδομένα έχουν ανωμαλίες που έχουν αρκετούς στενούς γείτονες, η τεχνική αποτυγχάνει να τις επισημάνει σωστά, με αποτέλεσμα να μην εντοπίζονται ανωμαλίες.
- Για τις τεχνικές με ημι-επίβλεψη, εάν οι κανονικές περιπτώσεις στα δεδομένα δοκιμής δεν έχουν αρκετές παρόμοιες κανονικές περιπτώσεις στα δεδομένα εκπαίδευσης, το ποσοστό ψευδώς θετικών αποτελεσμάτων για τέτοιες τεχνικές είναι υψηλό.
- Η υπολογιστική πολυπλοκότητα της φάσης δοκιμής αποτελεί επίσης σημαντική πρόκληση, καθώς περιλαμβάνει τον υπολογισμό της απόστασης κάθε περίπτωσης δοκιμής με όλες τις περιπτώσεις που ανήκουν είτε στα ίδια τα δεδομένα δοκιμής είτε στα δεδομένα εκπαίδευσης, ώστε να υπολογιστούν οι πλησιέστεροι γείτονες.
- Η απόδοση μιας τεχνικής που βασίζεται στον πλησιέστερο γείτονα εξαρτάται σε μεγάλο βαθμό από ένα μέτρο απόστασης, που ορίζεται μεταξύ ενός ζεύγους περιπτώσεων δεδομένων, το οποίο μπορεί να διακρίνει αποτελεσματικά μεταξύ κανονικών και ανώμαλων περιπτώσεων. Ο ορισμός μέτρων απόστασης μεταξύ περιπτώσεων μπορεί να είναι δύσκολος όταν τα δεδομένα είναι πολύπλοκα, π.χ. γράφοι, ακολουθίες κ.λπ.

5.3 Τεχνικές ανίχνευσης ανωμαλιών με βάση την ομαδοποίηση

Η ομαδοποίηση (Tan et al., 2005) χρησιμοποιείται για την ομαδοποίηση παρόμοιων περιπτώσεων δεδομένων σε ομάδες. Η συσταδοποίηση είναι κυρίως μια μη επιβλεπόμενη τεχνική, αν και έχει επίσης διερευνηθεί η ημι-επιβλεπόμενη συσταδοποίηση (Basu et al., 2004). Παρόλο που η συσταδοποίηση και η ανίχνευση ανωμαλιών φαίνεται να είναι θεμελιωδώς διαφορετικές μεταξύ τους, έχουν αναπτυχθεί αρκετές τεχνικές ανίχνευσης ανωμαλιών που βασίζονται στην συσταδοποίηση. Οι τεχνικές ανίχνευσης ανωμαλιών με βάση τη συσταδοποίηση μπορούν να ομαδοποιηθούν σε τρεις κατηγορίες.

Η πρώτη κατηγορία τεχνικών που βασίζονται στην ομαδοποίηση βασίζεται στην ακόλουθη υπόθεση:



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

Υπόθεση: Οι κανονικές περιπτώσεις δεδομένων ανήκουν σε μια συστάδα στα δεδομένα, ενώ οι ανωμαλίες δεν ανήκουν σε καμία.

Οι τεχνικές που βασίζονται στην παραπάνω υπόθεση εφαρμόζουν έναν γνωστό αλγόριθμο βασισμένο στην ομαδοποίηση στο σύνολο δεδομένων και δηλώνουν κάθε περίπτωση δεδομένων που δεν ανήκει σε καμία ομάδα ως ανώμαλη. Μπορούν να χρησιμοποιηθούν διάφοροι αλγόριθμοι συσταδοποίησης που δεν υποχρεώνουν κάθε περίπτωση δεδομένων να ανήκει σε μια συστάδα, όπως οι DBSCAN, ROCK και SNN (Ertoz et al., 2003). Ο αλγόριθμος FindOut είναι μια επέκταση του αλγορίθμου WaveCluster (Sheikholeslami et al., 1998), στον οποίο οι ανιχνευμένες συστάδες αφαιρούνται από τα δεδομένα και οι εναπομείνουσες περιπτώσεις δηλώνονται ως ανωμαλίες.

Ένα μειονέκτημα αυτών των τεχνικών είναι ότι δεν είναι βελτιστοποιημένες για την εύρεση ανωμαλιών, δεδομένου ότι ο κύριος στόχος του υποκείμενου αλγορίθμου ομαδοποίησης είναι η εύρεση συστάδων.

Η δεύτερη κατηγορία τεχνικών που βασίζονται στην ομαδοποίηση βασίζεται στην ακόλουθη υπόθεση:

Υπόθεση: Κανονικές περιπτώσεις δεδομένων βρίσκονται κοντά στο πλησιέστερο κεντροειδές της συστάδας τους, ενώ οι ανωμαλίες απέχουν πολύ από το πλησιέστερο κεντροειδές της συστάδας τους.

Οι τεχνικές που βασίζονται στην παραπάνω υπόθεση αποτελούνται από δύο βήματα. Στο πρώτο βήμα, τα δεδομένα ομαδοποιούνται χρησιμοποιώντας έναν αλγόριθμο ομαδοποίησης. Στο δεύτερο βήμα, για κάθε περίπτωση δεδομένων, η απόστασή της από το πλησιέστερο κεντροειδές της συστάδας υπολογίζεται ως η βαθμολογία ανωμαλίας.

Οι τεχνικές που βασίζονται στη δεύτερη παραδοχή μπορούν επίσης να λειτουργήσουν με ημι-επιβλεπόμενη λειτουργία, κατά την οποία τα δεδομένα εκπαίδευσης ομαδοποιούνται σε ομάδες και οι περιπτώσεις που ανήκουν στα δεδομένα δοκιμής συγκρίνονται με τις ομάδες για να ληφθεί μια βαθμολογία ανωμαλίας για την περίπτωση των δεδομένων δοκιμής (Vinueza & Grudic, 2004). Εάν τα δεδομένα



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

εκπαίδευσης έχουν περιπτώσεις που ανήκουν σε πολλαπλές κλάσεις, μπορεί να εφαρμοστεί ημι-επιβλεπόμενη συσταδοποίηση για τη βελτίωση των συστάδων. Οι He et al. (2002) ενσωματώνουν τη γνώση των ετικετών για να βελτιώσουν την τεχνική ανίχνευσης ανωμαλιών που βασίζεται στη μη επιβλεπόμενη συσταδοποίηση, υπολογίζοντας ένα μέτρο που ονομάζεται παράγοντας σημασιολογικής ανωμαλίας (semantic anomaly factor), ο οποίος είναι υψηλός εάν η ετικέτα κλάσης ενός αντικειμένου σε μια συστάδα διαφέρει από την πλειοψηφία των ετικετών κλάσης στην εν λόγω συστάδα.

Αν οι ανωμαλίες στα δεδομένα σχηματίζουν από μόνες τους συστάδες, οι προαναφερθείσες τεχνικές δεν θα είναι σε θέση να ανιχνεύσουν τέτοιες ανωμαλίες. Για την αντιμετώπιση αυτού του ζητήματος έχει προταθεί μια τρίτη κατηγορία τεχνικών που βασίζονται στην ομαδοποίηση και βασίζονται στην ακόλουθη υπόθεση:

Υπόθεση: Οι κανονικές περιπτώσεις δεδομένων ανήκουν σε μεγάλες και πυκνές συστάδες, ενώ οι ανωμαλίες είτε ανήκουν σε μικρές είτε σε αραιές συστάδες.

Οι τεχνικές που βασίζονται στην παραπάνω υπόθεση δηλώνουν ως ανώμαλες τις περιπτώσεις που ανήκουν σε συστάδες των οποίων το μέγεθος ή/και η πυκνότητα είναι κάτω από ένα κατώφλι.

Έχουν προταθεί διάφορες παραλλαγές της τρίτης κατηγορίας τεχνικών. Η τεχνική που προτείνεται από τους He et al. (2002), η οποία ονομάζεται FindCBLOF, αποδίδει μια βαθμολογία ανωμαλίας γνωστή ως Cluster-Based Local Outlier Factor (CBLOF) για κάθε περίπτωση δεδομένων. Η βαθμολογία CBLOF καταγράφει το μέγεθος της συστάδας στην οποία ανήκει η περίπτωση δεδομένων, καθώς και την απόσταση της περίπτωσης δεδομένων από το κεντροειδές της συστάδας της.

5.3.1 Υπολογιστική πολυπλοκότητα

Η υπολογιστική πολυπλοκότητα της εκπαίδευσης μιας τεχνικής ανίχνευσης ανωμαλιών που βασίζεται στην ομαδοποίηση εξαρτάται από τον αλγόριθμο ομαδοποίησης που χρησιμοποιείται για τη δημιουργία ομάδων από τα δεδομένα. Έτσι, τέτοιες τεχνικές μπορεί να έχουν τετραγωνική πολυπλοκότητα εάν η τεχνική ομαδοποίησης απαιτεί τον υπολογισμό αποστάσεων ανά ζεύγη για όλες τις περιπτώσεις δεδομένων, ή γραμμική όταν χρησιμοποιούνται τεχνικές που βασίζονται σε ευρετικές μεθόδους ή



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

προσεγγιστικές τεχνικές ομαδοποίησης (Eskin et al., 2002). Η φάση δοκιμής των τεχνικών που βασίζονται στην ομαδοποίηση είναι γρήγορη, δεδομένου ότι περιλαμβάνει τη σύγκριση μιας δοκιμαστικής περίπτωσης με έναν μικρό αριθμό συστάδων.

Τα πλεονεκτήματα των τεχνικών που βασίζονται στην ομαδοποίηση είναι τα εξής:

- Οι τεχνικές που βασίζονται στην ομαδοποίηση μπορούν να λειτουργούν χωρίς επίβλεψη.
- Τέτοιες τεχνικές μπορούν συχνά να προσαρμοστούν σε άλλους πολύπλοκους τύπους δεδομένων, απλά προσθέτοντας έναν αλγόριθμο ομαδοποίησης που μπορεί να χειριστεί τον συγκεκριμένο τύπο δεδομένων.
- Η φάση της δοκιμής για τις τεχνικές που βασίζονται στην ομαδοποίηση είναι γρήγορη, δεδομένου ότι ο αριθμός των ομάδων με τις οποίες πρέπει να συγκριθεί κάθε περίπτωση δοκιμής είναι μια μικρή σταθερά.

Τα μειονεκτήματα των τεχνικών που βασίζονται στην ομαδοποίηση είναι τα εξής:

- Η απόδοση των τεχνικών που βασίζονται στην ομαδοποίηση εξαρτάται σε μεγάλο βαθμό από την αποτελεσματικότητα του αλγορίθμου ομαδοποίησης όσον αφορά την καταγραφή της δομής των ομάδων των κανονικών περιπτώσεων.
- Πολλές τεχνικές ανιχνεύουν τις ανωμαλίες ως υποπροϊόν της ομαδοποίησης και, ως εκ τούτου, δεν είναι βελτιστοποιημένες για την ανίχνευση ανωμαλιών.
- Αρκετοί αλγόριθμοι ομαδοποίησης αναγκάζουν κάθε περίπτωση να ανατεθεί σε κάποια ομάδα. Αυτό μπορεί να έχει ως αποτέλεσμα οι ανωμαλίες να αντιστοιχίζονται σε μια μεγάλη συστάδα και να θεωρούνται ως κανονικές περιπτώσεις από τεχνικές που λειτουργούν με την υπόθεση ότι οι ανωμαλίες δεν ανήκουν σε καμία συστάδα.
- Διάφορες τεχνικές που βασίζονται στην ομαδοποίηση είναι αποτελεσματικές μόνο όταν οι ανωμαλίες δεν σχηματίζουν σημαντικές ομάδες μεταξύ τους.

5.4 Στατιστικές τεχνικές ανίχνευσης ανωμαλιών

Η βασική αρχή κάθε στατιστικής τεχνικής ανίχνευσης ανωμαλιών είναι:



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

«Μια ανωμαλία είναι μια παρατήρηση για την οποία υπάρχει υποψία ότι είναι εν μέρει ή εξ ολοκλήρου άσχετη, επειδή δεν παράγεται από το στοχαστικό μοντέλο που υποτίθεται» (Anscombe & Guttman, 1960). Οι στατιστικές τεχνικές ανίχνευσης ανωμαλιών βασίζονται στην ακόλουθη βασική υπόθεση:

Υπόθεση: Οι κανονικές περιπτώσεις δεδομένων εμφανίζονται σε περιοχές υψηλής πιθανότητας ενός στοχαστικού μοντέλου, ενώ οι ανωμαλίες εμφανίζονται στις περιοχές χαμηλής πιθανότητας του στοχαστικού μοντέλου.

Οι στατιστικές τεχνικές προσαρμόζουν ένα στατιστικό μοντέλο (συνήθως για κανονική συμπεριφορά) στα δεδομένα και στη συνέχεια εφαρμόζουν μια δοκιμή στατιστικής συμπερασματολογίας για να καθορίσουν αν μια αθέατη περίπτωση ανήκει σε αυτό το μοντέλο ή όχι. Οι περιπτώσεις που έχουν χαμηλή πιθανότητα να προέρχονται από το μοντέλο που μαθαίνεται, με βάση το στατιστικό τεστ που εφαρμόζεται, δηλώνονται ως ανωμαλίες. Για την προσαρμογή ενός στατιστικού μοντέλου έχουν εφαρμοστεί τόσο παραμετρικές όσο και μη παραμετρικές τεχνικές. Ενώ οι παραμετρικές τεχνικές προϋποθέτουν τη γνώση της υποκείμενης κατανομής και εκτιμούν τις παραμέτρους από τα δεδομένα που δίνονται, οι μη παραμετρικές τεχνικές δεν προϋποθέτουν γενικά τη γνώση της υποκείμενης κατανομής.

5.4.1 Παραμετρικές τεχνικές

Όπως αναφέρθηκε προηγουμένως, οι παραμετρικές τεχνικές υποθέτουν ότι τα κανονικά δεδομένα παράγονται από μια παραμετρική κατανομή με παραμέτρους θ και συνάρτηση πυκνότητας πιθανότητας $f(x, \theta)$, όπου x είναι μια παρατήρηση. Η βαθμολογία ανωμαλίας μιας περίπτωσης δοκιμής (ή παρατήρησης) x είναι το αντίστροφο της συνάρτησης πυκνότητας πιθανότητας, $f(x, \theta)$. Οι παράμετροι θ εκτιμώνται από τα δεδομένα που δίνονται.

Εναλλακτικά, μπορεί να χρησιμοποιηθεί ένας στατιστικός έλεγχος υποθέσεων (που αναφέρεται επίσης ως έλεγχος ασυμφωνίας στη βιβλιογραφία για την ανίχνευση στατιστικών ακραίων τιμών. Η μηδενική υπόθεση (H_0) για τέτοιους ελέγχους είναι ότι η περίπτωση δεδομένων x έχει παραχθεί χρησιμοποιώντας την εκτιμώμενη κατανομή (με παραμέτρους θ). Εάν ο στατιστικός έλεγχος απορρίψει την H_0 , το x δηλώνεται ως ανωμαλία. Ένας στατιστικός έλεγχος υποθέσεων σχετίζεται με μια στατιστική δοκιμής,



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

η οποία μπορεί να χρησιμοποιηθεί για να ληφθεί μια πιθανολογική βαθμολογία ανωμαλίας για την περίπτωση δεδομένων x .

Με βάση τον τύπο της υποτιθέμενης κατανομής, οι παραμετρικές τεχνικές μπορούν να κατηγοριοποιηθούν περαιτέρω ως εξής:

5.4.1.1 Με βάση το μοντέλο Gauss.

Οι τεχνικές αυτές υποθέτουν ότι τα δεδομένα παράγονται από μια κατανομή Gauss. Οι παράμετροι εκτιμώνται χρησιμοποιώντας εκτιμήσεις μέγιστης πιθανοφάνειας (MLE). Η απόσταση μιας περίπτωσης δεδομένων από την εκτιμώμενη μέση τιμή είναι η βαθμολογία ανωμαλίας για την εν λόγω περίπτωση. Ένα κατώφλι εφαρμόζεται στις βαθμολογίες ανωμαλίας για τον προσδιορισμό των ανωμαλιών. Οι διάφορες τεχνικές αυτής της κατηγορίας υπολογίζουν την απόσταση από τον μέσο όρο και το κατώφλι με διαφορετικούς τρόπους.

5.4.1.2 Μοντέλο παλινδρόμησης.

Η ανίχνευση ανωμαλιών με χρήση παλινδρόμησης έχει διερευνηθεί εκτενώς για δεδομένα χρονοσειρών (Abraham & Chuang, 1989). Η βασική τεχνική ανίχνευσης ανωμαλιών με βάση το μοντέλο παλινδρόμησης αποτελείται από δύο βήματα. Στο πρώτο βήμα, ένα μοντέλο παλινδρόμησης προσαρμόζεται στα δεδομένα. Στο δεύτερο βήμα, για κάθε περίπτωση δοκιμής, το υπόλειμμα για την περίπτωση δοκιμής χρησιμοποιείται για τον προσδιορισμό του σκορ ανωμαλίας. Το υπόλειμμα είναι το μέρος της περίπτωσης που δεν εξηγείται από το μοντέλο παλινδρόμησης. Το μέγεθος του υπολοίπου μπορεί να χρησιμοποιηθεί ως βαθμολογία ανωμαλίας για την περίπτωση δοκιμής, αν και έχουν προταθεί στατιστικές δοκιμές για τον προσδιορισμό ανωμαλιών με ορισμένη εμπιστοσύνη (Torr & Murray, 1993).

Έχουν προταθεί παραλλαγές της βασικής τεχνικής που βασίζεται σε μοντέλα παλινδρόμησης για τη διαχείριση πολυμεταβλητών δεδομένων χρονοσειρών. Οι Tsay et al. (2000) συζητούν την πρόσθετη πολυπλοκότητα των πολυμεταβλητών χρονοσειρών σε σχέση με τις μονομεταβλητές χρονοσειρές και καταλήγουν σε στατιστικά στοιχεία που μπορούν να εφαρμοστούν για την ανίχνευση ανωμαλιών σε πολυμεταβλητά μοντέλα ARIMA.



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

Μια άλλη παραλλαγή που ανιχνεύει ανωμαλίες σε πολυμεταβλητές χρονοσειρές δεδομένων που παράγονται από ένα μοντέλο αυτοπαλινδρομικού κινητού μέσου (ARMA), προτάθηκε από τους Galeano et al. (2004). Σε αυτή την τεχνική οι συγγραφείς μετατρέπουν τις πολυμεταβλητές χρονοσειρές σε μονομεταβλητές χρονοσειρές συνδυάζοντας γραμμικά τις συνιστώσες των πολυμεταβλητών χρονοσειρών.

5.4.2 Μη παραμετρικές τεχνικές

Οι τεχνικές ανίχνευσης ανωμαλιών αυτής της κατηγορίας χρησιμοποιούν μη παραμετρικά στατιστικά μοντέλα, έτσι ώστε η δομή του μοντέλου να μην ορίζεται εκ των προτέρων, αλλά να καθορίζεται από τα δεδομένα. Τέτοιες τεχνικές κάνουν συνήθως λιγότερες υποθέσεις σχετικά με τα δεδομένα, όπως η ομαλότητα της πυκνότητας, σε σύγκριση με τις παραμετρικές τεχνικές.

5.4.2.1 Με βάση το ιστόγραμμα.

Η απλούστερη μη παραμετρική στατιστική τεχνική είναι η χρήση ιστογραμμάτων για τη διατήρηση ενός προφίλ των κανονικών δεδομένων. Τέτοιες τεχνικές αναφέρονται επίσης ως βασισμένες στη συχνότητα ή στη μέτρηση. Οι τεχνικές που βασίζονται στα ιστογράμματα είναι ιδιαίτερα δημοφιλείς στην κοινότητα της ανίχνευσης εισβολών (Eskin, 2000) και της ανίχνευσης απάτης (Fawcett & Provost, 1999), δεδομένου ότι η συμπεριφορά των δεδομένων διέπεται από ορισμένα προφίλ (χρήστη ή λογισμικού ή συστήματος) που μπορούν να καταγραφούν αποτελεσματικά με τη χρήση του μοντέλου ιστογράμματος.

Μια βασική τεχνική ανίχνευσης ανωμαλιών με βάση το ιστόγραμμα για μονομεταβλητά δεδομένα αποτελείται από δύο βήματα. Το πρώτο βήμα περιλαμβάνει τη δημιουργία ενός ιστογράμματος με βάση τις διαφορετικές τιμές που λαμβάνει το συγκεκριμένο χαρακτηριστικό στα δεδομένα εκπαίδευσης. Στο δεύτερο βήμα, λέγεται αν μια δοκιμαστική περίπτωση εμπίπτει σε οποιοδήποτε από τα δοχεία του ιστογράμματος. Εάν ναι, η περίπτωση δοκιμής είναι φυσιολογική, διαφορετικά είναι ανώμαλη.

Για πολυμεταβλητά δεδομένα, μια βασική τεχνική είναι η κατασκευή ιστογραμμάτων ανά χαρακτηριστικό. Κατά τη διάρκεια της δοκιμής, για κάθε περίπτωση δοκιμής, η



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

βαθμολογία ανωμαλίας για κάθε τιμή χαρακτηριστικού της περίπτωσης δοκιμής υπολογίζεται ως το ύψος του δοχείου που περιέχει την τιμή του χαρακτηριστικού. Οι βαθμολογίες ανωμαλίας ανά χαρακτηριστικό αθροίζονται για να ληφθεί μια συνολική βαθμολογία ανωμαλίας για την περίπτωση δοκιμής.

5.4.2.2 Με βάση τη συνάρτηση Kernel.

Μια μη παραμετρική τεχνική για την εκτίμηση πυκνότητας πιθανότητας είναι η εκτίμηση παραθύρων Parzen (Parzen, 1962). Αυτή περιλαμβάνει τη χρήση συναρτήσεων πυρήνα για την προσέγγιση της πραγματικής πυκνότητας. Οι τεχνικές ανίχνευσης ανωμαλιών που βασίζονται σε συναρτήσεις πυρήνα είναι παρόμοιες με τις παραμετρικές μεθόδους που περιγράφηκαν προηγουμένως. Η μόνη διαφορά είναι η τεχνική εκτίμησης πυκνότητας που χρησιμοποιείται. Οι Desforges (1998) πρότειναν μια στατιστική τεχνική με ημι-επίβλεψη για την ανίχνευση ανωμαλιών, η οποία χρησιμοποιεί συναρτήσεις πυρήνα για την εκτίμηση της συνάρτησης κατανομής πιθανότητας για τις κανονικές περιπτώσεις. Μια νέα περίπτωση που βρίσκεται στην περιοχή χαμηλής πιθανότητας αυτής της συνάρτησης κατανομής πιθανότητας δηλώνεται ως ανώμαλη. Παρόμοια εφαρμογή των παραθύρων parzen προτείνεται για την ανίχνευση εισβολών σε δίκτυα (Chow & Yeung, 2002), για την ανίχνευση καινοτομιών σε δεδομένα ροής πετρελαίου (Bishop, 1994) και για την ανάλυση μαστογραφικών εικόνων (Tarassenko, 1995).

5.4.3 Υπολογιστική πολυπλοκότητα

Η υπολογιστική πολυπλοκότητα των στατιστικών τεχνικών ανίχνευσης ανωμαλιών εξαρτάται από τη φύση του στατιστικού μοντέλου που απαιτείται να προσαρμοστεί στα δεδομένα. Τα πλεονεκτήματα των στατιστικών τεχνικών είναι:

- Εάν οι υποθέσεις σχετικά με την υποκείμενη κατανομή των δεδομένων ισχύουν, οι στατιστικές τεχνικές παρέχουν μια στατιστικά δικαιολογημένη λύση για την ανίχνευση ανωμαλιών.
- Η βαθμολογία ανωμαλίας που παρέχεται από μια στατιστική τεχνική συνδέεται με ένα διάστημα εμπιστοσύνης, το οποίο μπορεί να χρησιμοποιηθεί ως πρόσθετη πληροφορία κατά τη λήψη απόφασης σχετικά με κάθε περίπτωση δοκιμής.



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

- Εάν το βήμα εκτίμησης της κατανομής είναι ανθεκτικό στις ανωμαλίες των δεδομένων, οι στατιστικές τεχνικές μπορούν να λειτουργήσουν σε ένα περιβάλλον χωρίς επίβλεψη, χωρίς να χρειάζονται δεδομένα εκπαίδευσης με ετικέτες.

Τα μειονεκτήματα των στατιστικών τεχνικών είναι:

- Το βασικό μειονέκτημα των στατιστικών τεχνικών είναι ότι βασίζονται στην υπόθεση ότι τα δεδομένα προέρχονται από μια συγκεκριμένη κατανομή. Η υπόθεση αυτή συχνά δεν ισχύει, ιδίως για σύνολα πραγματικών δεδομένων υψηλής διάστασης.
- Ακόμα και όταν η στατιστική υπόθεση μπορεί να δικαιολογηθεί λογικά, υπάρχουν πολλά στατιστικά στοιχεία ελέγχου υποθέσεων που μπορούν να εφαρμοστούν για τον εντοπισμό ανωμαλιών, η επιλογή του καλύτερου στατιστικού στοιχείου συχνά δεν είναι απλή υπόθεση. Ειδικότερα, η κατασκευή ελέγχων στατιστικών υποθέσεων για πολύπλοκες κατανομές πολλών διαστάσεων είναι προβληματική.

5.5 Χειρισμός ανωμαλιών του πλαισίου

Οι ανωμαλίες του πλαισίου απαιτούν τα δεδομένα να διαθέτουν ένα σύνολο από χαρακτηριστικά του πλαισίου (για τον ορισμό ενός πλαισίου) και ένα σύνολο από χαρακτηριστικά συμπεριφοράς (για την ανίχνευση ανωμαλιών σε ένα πλαίσιο). Οι Song et al. (2007) χρησιμοποιούν τους όρους περιβαλλοντικά χαρακτηριστικά και χαρακτηριστικά δείκτη. Ορισμένοι από τους τρόπους με τους οποίους μπορούν να οριστούν τα χαρακτηριστικά περιβάλλοντος είναι οι εξής:

(1) Χωροταξικά: Τα δεδομένα έχουν χωρικά χαρακτηριστικά, τα οποία καθορίζουν τη θέση μιας περίπτωσης δεδομένων και, ως εκ τούτου, μια χωρική γειτονιά. Για δεδομένα με χωρικά δεδομένα έχουν προταθεί διάφορες τεχνικές ανίχνευσης ανωμαλιών με βάση το πλαίσιο (Lu et al., 2003).

(2) Γραφήματα: Οι ακμές που συνδέουν κόμβους (περιπτώσεις δεδομένων) ορίζουν τη γειτονιά για κάθε κόμβο. Οι τεχνικές ανίχνευσης ανωμαλιών με βάση το πλαίσιο έχουν εφαρμοστεί σε δεδομένα που βασίζονται σε γραφήματα από τους Sun et al. (2005).



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

(3) Διαδοχική: Δηλαδή, τα συμφραζόμενα χαρακτηριστικά μιας περίπτωσης δεδομένων είναι η θέση της στην ακολουθία. Τα δεδομένα χρονοσειρών έχουν διερευνηθεί εκτενώς στην κατηγορία της ανίχνευσης ανωμαλιών με βάση το πλαίσιο (Salvador & Chan, 2003). Μια άλλη μορφή διαδοχικών δεδομένων για την οποία έχουν αναπτυχθεί τεχνικές ανίχνευσης ανωμαλιών είναι τα δεδομένα συμβάντων, στα οποία κάθε συμβάν έχει μια χρονοσφραγίδα (όπως τα δεδομένα κλήσεων του λειτουργικού συστήματος ή τα δεδομένα ιστού). Η διαφορά μεταξύ των δεδομένων χρονοσειρών και των ακολουθιών συμβάντων είναι ότι για τα τελευταία, ο χρόνος μεταξύ των διαδοχικών συμβάντων είναι άνισος.

(4) Προφίλ: Συχνά τα δεδομένα μπορεί να μην έχουν σαφή χωρική ή διαδοχική δομή, αλλά μπορούν να τμηματοποιηθούν ή να ομαδοποιηθούν σε συνιστώσες χρησιμοποιώντας ένα σύνολο χαρακτηριστικών πλαισίου. Αυτά τα χαρακτηριστικά χρησιμοποιούνται συνήθως για τη δημιουργία προφίλ και την ομαδοποίηση των χρηστών σε συστήματα παρακολούθησης δραστηριοτήτων, όπως η ανίχνευση απάτης σε κινητά τηλέφωνα (Fawcett & Provost, 1999), οι βάσεις δεδομένων CRM (He et al., 2004) και η ανίχνευση απάτης με πιστωτικές κάρτες (Bolton & Hand, 1999).

Σε σύγκριση με την πλούσια βιβλιογραφία σχετικά με τις τεχνικές ανίχνευσης σημειακών ανωμαλιών, η έρευνα σχετικά με την ανίχνευση ανωμαλιών βάσει πλαισίου είναι περιορισμένη. Σε γενικές γραμμές, οι τεχνικές αυτές μπορούν να ταξινομηθούν σε δύο κατηγορίες. Η πρώτη κατηγορία τεχνικών ανάγει ένα πρόβλημα ανίχνευσης ανωμαλιών με βάση το πλαίσιο σε πρόβλημα ανίχνευσης σημειακών ανωμαλιών, ενώ η δεύτερη κατηγορία τεχνικών μοντελοποιεί τη δομή των δεδομένων και χρησιμοποιεί το μοντέλο για την ανίχνευση ανωμαλιών.

5.5.1 Αξιοποίηση της δομής των δεδομένων

Σε διάφορα σενάρια, ο διαχωρισμός των δεδομένων σε πλαίσια δεν είναι απλός. Αυτό ισχύει συνήθως για δεδομένα χρονοσειρών και δεδομένα ακολουθίας γεγονότων. Σε τέτοιες περιπτώσεις, οι τεχνικές μοντελοποίησης χρονοσειρών και μοντελοποίησης ακολουθιών επεκτείνονται για την ανίχνευση ανωμαλιών πλαισίου στα δεδομένα.

Μια γενική τεχνική αυτής της κατηγορίας μπορεί να περιγραφεί ως εξής. Από τα δεδομένα εκπαίδευσης μαθαίνεται ένα μοντέλο το οποίο μπορεί να προβλέψει την



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

αναμενόμενη συμπεριφορά σε σχέση με ένα δεδομένο πλαίσιο. Εάν η αναμενόμενη συμπεριφορά διαφέρει σημαντικά από την παρατηρούμενη συμπεριφορά, δηλώνεται μια ανωμαλία. Ένα απλό παράδειγμα αυτής της γενικής τεχνικής είναι η παλινδρόμηση στην οποία τα χαρακτηριστικά του πλαισίου μπορούν να χρησιμοποιηθούν για την πρόβλεψη του χαρακτηριστικού συμπεριφοράς με την προσαρμογή μιας γραμμής παλινδρόμησης στα δεδομένα.

Ένα μοντέλο για την ανίχνευση εισβολών σε τηλεφωνικά δίκτυα προτάθηκε από τον Scott (2001) και για τη μοντελοποίηση δεδομένων διαδικτυακών κλικ από τους Ihler et al. (2006). Και οι δύο εργασίες ακολουθούν μια τεχνική κατά την οποία υποθέτουν ότι η κανονική συμπεριφορά σε μια χρονοσειρά παράγεται από μια μη σταθερή διαδικασία Poisson, ενώ οι ανωμαλίες παράγονται από μια ομοιογενή διαδικασία Poisson. Η μετάβαση μεταξύ κανονικής και ανώμαλης συμπεριφοράς μοντελοποιείται με τη χρήση μιας διαδικασίας Markov. Οι προτεινόμενες τεχνικές σε κάθε μία από αυτές τις εργασίες χρησιμοποιούν την τεχνική εκτίμησης Markov Chain Monte Carlo (MCMC) για την εκτίμηση των παραμέτρων αυτών των διαδικασιών. Για τον έλεγχο, μια χρονοσειρά μοντελοποιείται με τη χρήση αυτής της διαδικασίας και οι χρονικές στιγμές για τις οποίες η ανώμαλη συμπεριφορά ήταν ενεργή θεωρούνται ως ανωμαλίες.

5.5.2 Υπολογιστική πολυπλοκότητα

Η υπολογιστική πολυπλοκότητα της φάσης εκπαίδευσης στις τεχνικές ανίχνευσης ανωμαλιών με βάση το πλαίσιο που χρησιμοποιούν τη δομή των δεδομένων για τη δημιουργία μοντέλων, είναι συνήθως υψηλότερη από εκείνη των τεχνικών που περιορίζουν το πρόβλημα στην ανίχνευση σημειακών ανωμαλιών. Ένα πλεονέκτημα για τέτοιες τεχνικές είναι ότι η φάση δοκιμής είναι σχετικά γρήγορη, δεδομένου ότι κάθε περίπτωση απλώς συγκρίνεται με το ενιαίο μοντέλο και της αποδίδεται μια βαθμολογία ανωμαλίας ή μια ετικέτα ανωμαλίας. Το βασικό πλεονέκτημα των τεχνικών ανίχνευσης ανωμαλιών με βάση το πλαίσιο είναι ότι επιτρέπουν έναν φυσικό ορισμό μιας ανωμαλίας σε πολλές εφαρμογές της πραγματικής ζωής, όπου οι περιπτώσεις δεδομένων τείνουν να είναι παρόμοιες μέσα σε ένα πλαίσιο. Τέτοιες τεχνικές είναι σε θέση να ανιχνεύουν ανωμαλίες που μπορεί να μην ανιχνεύονται από τεχνικές ανίχνευσης σημειακών ανωμαλιών που λαμβάνουν μια σφαιρική άποψη των



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

δεδομένων. Το μειονέκτημα των τεχνικών ανίχνευσης ανωμαλιών με βάση το πλαίσιο είναι ότι εφαρμόζονται μόνο όταν μπορεί να οριστεί ένα πλαίσιο.

Κεφάλαιο 6ο : Εφαρμογές ανίχνευσης ανωμαλιών

Στο παρόν κεφάλαιο θα μελετηθούν δύο εργασίες στις οποίες εφαρμόζονται σε πραγματικό χρόνο μέθοδοι ανίχνευσης ανωμαλιών, ώστε να καταστεί πιο σαφής η λειτουργία των εφαρμογών ανίχνευσης ανωμαλιών και το πως συνδράμουν στην επίλυση κρίσιμων ζητημάτων σε διαφορετικούς τομείς. Ειδικότερα, θα μελετηθεί η περίπτωση των (Méndez et al., 2022), όπου εφαρμόζονται δύο μέθοδοι βαθιάς μάθησης, η CNN και η BiLSTM για την ανίχνευση ανωμαλιών σε δεδομένα κυκλοφοριακής ροής στην πόλη της Μαδρίτης. Η άλλη εργασία που θα μελετηθεί είναι των (Zhang & Lei, 2021) και αφορά αμιγώς ένα πρόβλημα που εμπίπτει στον κλάδο των πολιτικών μηχανικών, τα δομικά δεδομένα. Και σε αυτή την εργασία εφαρμόζεται η μέθοδος μονοδιάστατης CNN, σε στατικά δεδομένα επιτάχυνσης για μία γέφυρα στην Κίνα. Παρά την ίδια μέθοδο βαθιάς μάθησης που χρησιμοποιείται, στην εργασία των ...2 ακολουθείται διαφορετική επεξεργασία στα δεδομένα, πριν χρησιμοποιηθούν στο νευρωνικό δίκτυο.

6.1 Ανίχνευση ανωμαλιών στην κυκλοφοριακή ροή

Ένα σημαντικό πρόβλημα στη διαχείριση της κυκλοφορίας στους αυτοκινητοδρόμους είναι η ανίχνευση συμβάντων/ανωμαλιών, δηλαδή η παρακολούθηση των ροών δεδομένων από ενσωματωμένους αισθητήρες που καταγράφουν την κυκλοφορία σε πραγματικό χρόνο με στόχο να ειδοποιηθούν οι αρμόδιοι φορείς για πιθανά τα προβλήματα σε πραγματικό χρόνο. Ο τομέας αυτός παρουσιάζει ένα κύριο πρόβλημα που δεν εμφανίζεται σε άλλες περιπτώσεις δεδομένων χρονοσειρών, την έλλειψη γνώσης σχετικά με το τι συνιστά, ανωμαλία. Ως εκ τούτου, η τεχνική αυτή απαιτεί την εφαρμογή μη επιβλεπόμενων μοντέλων¹. Μια πρώτη και απλοϊκή προσέγγιση, είναι να σκεφτεί κανείς ότι οι ανωμαλίες ταυτίζονται με τις ακραίες τιμές. Παρόλα αυτά,

¹ Η μη επιβλεπόμενη μάθηση, γνωστή και ως μη επιβλεπόμενη μηχανική μάθηση, χρησιμοποιεί αλγόριθμους μηχανικής μάθησης για την ανάλυση και ομαδοποίηση μη επισημασμένων συνόλων δεδομένων. Αυτοί οι αλγόριθμοι ανακαλύπτουν κρυμμένα μοτίβα ή ομαδοποιήσεις δεδομένων χωρίς την ανάγκη ανθρώπινης παρέμβασης.



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

ιδανικά, ένα μοντέλο με καλή απόδοση θα πρέπει να είναι σε θέση να ανιχνεύει όχι μόνο τις ακραίες τιμές αλλά και τις τιμές που υποδηλώνουν, για παράδειγμα, την αποτυχία ενός συγκεκριμένου προτύπου ή μια ξαφνική αύξηση ή μείωση σε σχέση με μια προηγούμενη τιμή.

Οι Mendez et al., (2022) παρουσιάζουν δύο μοντέλα για την ανίχνευση ανωμαλιών σε δεδομένα χρονοσειρών χρησιμοποιώντας προηγούμενα δεδομένα ως δεδομένα εισόδου. Προκειμένου να αξιολογήσουν τη χρησιμότητα των μοντέλων τους, τα εφαρμόζουν σε δεδομένα κυκλοφοριακής ροής που λαμβάνονται στην πόλη της Μαδρίτης. Αμφότερα τα μοντέλα τα οποία χρησιμοποιήθηκαν στην εργασία τους είναι μεταβλητού αυτόματου κωδικοποιητής (VAE)², ο οποίος είναι ένας αλγόριθμος που αποτελείται από έναν κωδικοποιητή και έναν αποκωδικοποιητή. Πρώτα, αναπτύχθηκε ένας αυτόματος κωδικοποιητής συνελκτικού νευρωνικού δικτύου (CNN)³. Εκτός από την εξαρτημένη μεταβλητή, το μοντέλο αυτό χρησιμοποίησε άλλες οκτώ ανεξάρτητες μεταβλητές που σχετίζονται με τη ροή της κυκλοφορίας, δηλαδή λειτουργούσε στο πλαίσιο μιας πολυμεταβλητής χρονοσειράς. Το δεύτερο μοντέλο ήταν ένα νέο μοντέλο αυτόματου κωδικοποιητή Αμφίδρομης Μακροχρόνιας Βραχυπρόθεσμης μνήμης ή Bidirectional (LSTM)⁴. Αυτό το μοντέλο χρησιμοποιούσε ως είσοδο μόνο τα δεδομένα της κυκλοφοριακής ροής (εξαρτημένη μεταβλητή), δηλαδή λειτουργούσε στο πλαίσιο μιας μονομεταβλητής χρονοσειράς.

6.1.1 Δεδομένα και μεθοδολογία

6.1.1.1 Δεδομένα

Στην εργασία των Mendez et al. (2022) χρησιμοποιήθηκαν ωριαία δεδομένα ροής κυκλοφορίας από τον Ιανουάριο του 2018 έως τον Ιούλιο του 2020, ενός σταθμού που βρίσκεται στη Μαδρίτη. Χρησιμοποιήθηκαν επιπρόσθετα οκτώ μεταβλητές που σχετίζονται με τα δεδομένα κυκλοφοριακής ροής για να επιτρέψουν στο μοντέλο να κατανοήσει καλύτερα το πλαίσιο. Αυτές οι οκτώ μεταβλητές είναι οι εξής:

- (1, 2) Ωριαία ροή κυκλοφορίας σε δύο κοντινούς πρόσθετους σταθμούς.
- (3, 4, 5) Μέση ημερήσια, μέγιστη και ελάχιστη θερμοκρασία.

² Variational auto-encoders

³ Convolutional Neural Network

⁴ Long short-term memory



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

- (6) Ημερήσια βροχόπτωση.
- (7) Είδος ημέρας (εργάσιμη ημέρα, Σαββατοκύριακο ή αργία).
- (8) Χρονοσφραγίδα των δεδομένων.

Οι τιμές των καιρικών μεταβλητών (3-6) μετρήθηκαν στον πλησιέστερο μετεωρολογικό σταθμό, ενώ η ροή της κυκλοφορίας που αφορά πρόσθετους σταθμούς (1-2) μετρήθηκε στους δύο πλησιέστερους σταθμούς κυκλοφορίας προς το σταθμό κυκλοφορίας-στόχο.

Σενάρια

Στην εργασία τους ορίστηκαν τρεις διαφορετικοί τύποι για τις μη φυσιολογικές τιμές των χρονοσειρών (Rajeswari et al., 2018):

- Μια τιμή θεωρείται παγκόσμια ανωμαλία (universal anomaly) εάν βρίσκεται πολύ έξω από το εύρος του συνόλου δεδομένων.
- Μια τιμή θεωρείται ανωμαλία πλαισίου (contextual anomaly) εάν αποκλίνει σημαντικά από άλλα δεδομένα σε παρόμοιο πλαίσιο.
- Ένα υποσύνολο του συνόλου δεδομένων θεωρείται συλλογική ανωμαλία (collective anomaly) εάν οι τιμές του υποσυνόλου, αποκλίνουν σημαντικά από το σύνολο δεδομένων, αλλά μεμονωμένα οι τιμές δεν θεωρούνται από μόνες τους ανώμαλες.

Για την εξυπηρέτηση του ερευνητικού στόχου, της αυτόματης ανίχνευσης μη φυσιολογικής συμπεριφοράς δεδομένων ροής κυκλοφορίας στη Μαδρίτη, χρησιμοποιήθηκαν δύο διαφορετικά σενάρια, ένα βασικό και ένα καθοδηγούμενο. Στο βασικό σενάριο έγινε χρήση όλων των δεδομένων για εκπαίδευση και στόχος των μοντέλων (CNN και Bivariate LSTM) είναι η ανάλυση και ο εντοπισμός τιμών που αντιστοιχούν σε μη φυσιολογικά δεδομένα. Η αποτελεσματικότητα των μοντέλων αξιολογήθηκε αναλύοντας γραφικά αν οι τιμές που θεωρούνται ανωμαλίες αντιστοιχούν σε κάθε είδους ανώμαλη τιμή στη χρονοσειρά και ελέγχοντας αν συνέβη κάποιο αξιοσημείωτο γεγονός στη Μαδρίτη την ημερομηνία κατά την οποία ανιχνεύονται τα ανώμαλα δεδομένα. Στο καθοδηγούμενο σενάριο ο στόχος ήταν ο ίδιος, αλλά τα δεδομένα χωρίζονται σε δύο σύνολα. Το πρώτο σύνολο περιλαμβάνει τις τιμές που είναι μεγαλύτερες από ένα υψηλό εκατοστημόριο του συνόλου



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

δεδομένων, για τις οποίες ήταν γνωστό με βεβαιότητα ότι είναι ανώμαλες τιμές. Το δεύτερο σύνολο περιλάμβανε τις υπόλοιπες τιμές των δεδομένων.

6.1.1.2 Μέθοδος VAE

Ο μεταβλητός αυτόματος κωδικοποιητής (VAE) είναι ένας αλγόριθμος που αποτελείται από έναν κωδικοποιητή και έναν αποκωδικοποιητή. Η ιδέα είναι ότι ο κωδικοποιητής μεταφράζει την είσοδο, X , από τον χώρο υψηλής διάστασης σε έναν χώρο χαμηλότερης διάστασης. Στη συνέχεια, ο αποκωδικοποιητής θα λάβει αυτή τη χαμηλότερης διάστασης αναπαράσταση της εισόδου, αυτό που συνήθως αποκαλούμε λανθάνον διάνυσμα, και θα προσπαθήσει να ανακατασκευάσει την αρχική είσοδο στον χώρο υψηλής διάστασης, παράγοντας το \tilde{X} . Λόγω της φύσης του, ο αποκωδικοποιητής έχει συνήθως μια δομή που μιμείται αυτή του κωδικοποιητή, καθώς στόχος του είναι να λάβει μια τιμή \tilde{X} όσο το δυνατόν πιο όμοια με την αρχική είσοδο, X . Η ανομοιότητα μεταξύ των αρχικών και των ανακατασκευασμένων δεδομένων ορίζεται από μια συνάρτηση απωλειών. Επομένως, στόχος του μοντέλου είναι να μειώσει, όσο το δυνατόν περισσότερο, την τιμή αυτή κατά τη διαδικασία εκπαίδευσης.

Ο βαθμός της ανομοιότητας υπολογίζεται συγκρίνοντας τα δεδομένα εισόδου με τα ανακατασκευασμένα δεδομένα εξόδου χρησιμοποιώντας ένα μέτρο ανομοιότητας. Στην συνέχεια οι βαθμολογίες ανομοιότητας χρησιμοποιούνται για τον υπολογισμό ενός κατωφλίου δ . Αυτό το κατώφλι ορίζεται συνήθως έτσι ώστε να έχει το $90 - 99^\circ$ εκατοστημόριο των βαθμολογιών ανομοιότητας ως μη φυσιολογικά δεδομένα. Στη συνέχεια, για νέα δεδομένα, εάν το μοντέλο μπορεί να ανακατασκευάσει δεδομένα εισόδου με τιμή ανομοιότητας μικρότερη από δ , τότε θα θεωρηθούν κανονικά δεδομένα- διαφορετικά, θα ταξινομηθούν ως μη φυσιολογικά δεδομένα.

Όσον αφορά τα δύο διαφορετικά σενάρια που έλεγξαν οι Mendez et al. (2022), η διαφορά τους έγκειται στην εκπαίδευση. Στο βασικό σενάριο, το μοντέλο «εκπαιδεύτηκε» με όλα τα δεδομένα (κανονικά και μη κανονικά), ενώ στο καθοδηγούμενο σενάριο το μοντέλο «εκπαιδεύτηκε» μόνο με τα δεδομένα που θεωρούνται κανονικά. Η ιδέα είναι ότι, σε αυτή τη δεύτερη περίπτωση, ο αλγόριθμος θα μάθει να αναπαριστά καλύτερα τα κανονικά δεδομένα, καθώς τα μη κανονικά δεδομένα έχουν αφαιρεθεί από την εκπαίδευση.



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

6.1.1.3 Μοντέλο CNN αυτόματου αποκωδικοποιητή

Έστω y η ροή κυκλοφορίας στο σταθμό-στόχο, x^1, \dots, x^8 οι πρόσθετες μεταβλητές του μοντέλου CNN και t ένα χρονικό βήμα. Για να αποφασίσουμε αν τα δεδομένα της ροής κυκλοφορίας στο χρονικό βήμα t , δηλαδή τα y_t είναι ανώμαλα δεδομένα, χρησιμοποιήθηκαν οι μεταβλητές y, x^1, \dots, x^8 από το χρονικό βήμα $t - 23$ έως t . Οι τιμές αυτές αναδιαμορφώθηκαν υπό μορφή πίνακα 24×9 , όπου 24 είναι ο αριθμός των τιμών που χρησιμοποιήθηκαν για να την ανίχνευση των ανωμαλιών:

$$X = \begin{bmatrix} y_{t-23} x_{t-23}^0 \dots x_{t-23}^8 \\ y_{t-22} x_{t-22}^0 \dots x_{t-22}^8 \\ \dots \\ y_t x_t^0 \dots x_t^8 \end{bmatrix} \quad (1.13)$$

Επιπλέον, χρησιμοποιήθηκε ως συνάρτηση απώλειας ανομοιότητας το μέσο απόλυτο σφάλμα:

$$Loss(X, \tilde{X}) = MAE = \frac{1}{216} \sum_{i=1}^{24} \sum_{j=1}^9 |X_{i,j} - \tilde{X}_{i,j}| \quad (1.14)$$

όπου i και j αντιπροσωπεύουν, αντίστοιχα, τις γραμμές και τις στήλες του πίνακα.

6.1.1.4 Μοντέλο BiLSTM αυτόματου αποκωδικοποιητή

Αυτό το μοντέλο αποτελείται επίσης από έναν κωδικοποιητή και έναν αποκωδικοποιητή. Στην περίπτωση αυτή, χρησιμοποιούνται στρώματα BiLSTM επειδή αξιοποιούν τα πλεονεκτήματα των αρχιτεκτονικών BiLSTM και κωδικοποιητή-αποκωδικοποιητή, αντίστοιχα. Η γενική συμπεριφορά αυτού του μοντέλου είναι παρόμοια με το μοντέλο CNN που περιγράφηκε προηγουμένως.

Προκειμένου να προσδιοριστεί αν τα δεδομένα της ροής κυκλοφορίας στο χρονικό βήμα t , δηλαδή η τιμή y_t , είναι ανώμαλα, χρησιμοποιείται σε αυτή την περίπτωση η μονομεταβλητή χρονοσειρά της ροής κυκλοφορίας, από το χρονικό βήμα $t - 23$ έως t . Χρησιμοποιούμε ως δεδομένα εισόδου ένα διάνυσμα μεγέθους 24, όπου 24 είναι



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

ο αριθμός των προηγούμενων τιμών που χρησιμοποιήθηκαν στην CNN για την ανίχνευση των ανωμαλιών:

$$X = [y_{t-23}, y_{t-22}, \dots, y_t] \quad (1.15)$$

Στην περίπτωση αυτή, η ανομοιότητα μετράται επίσης μεταξύ της ακολουθίας εισόδου και της εξόδου, δηλαδή της ανακατασκευασμένης ακολουθίας. Χρησιμοποιείται και πάλι το μέσο απόλυτο σφάλμα ως συνάρτηση απώλειας ανομοιότητας:

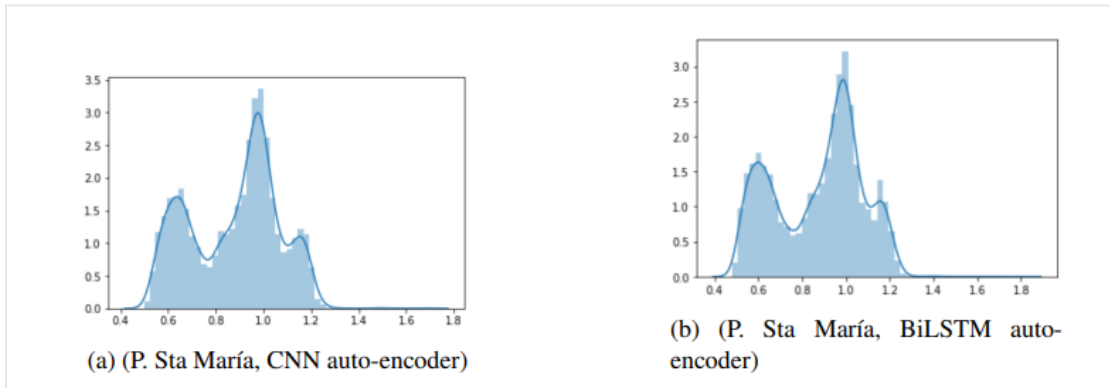
$$Loss(X, \tilde{X}) = MAE = \frac{1}{24} \sum_{t=23}^t |y_i - \tilde{y}_i| \quad (1.16)$$

6.1.2 Αποτελέσματα μοντέλων

6.1.2.1 Βασικό σενάριο

Η εφαρμογή των μοντέλων αυτόματου κωδικοποιητή CNN και BiLSTM στο βασικό σενάριο αφορούσε τον σταθμό κυκλοφορίας Paseo de Santa Maria de la Cabeza (N40.4065° E3.6946°). Σε αμφότερες τις περιπτώσεις χρησιμοποιήθηκαν ως σύνολο εκπαίδευσης δεδομένα 16.391 παρατηρήσεων και ως σύνολο δοκιμής δεδομένα 5.448 παρατηρήσεων. Εκπαιδύουμε όλα τα μοντέλα χρησιμοποιώντας 15 εποχές και 16 ως μέγεθος δέσμης.

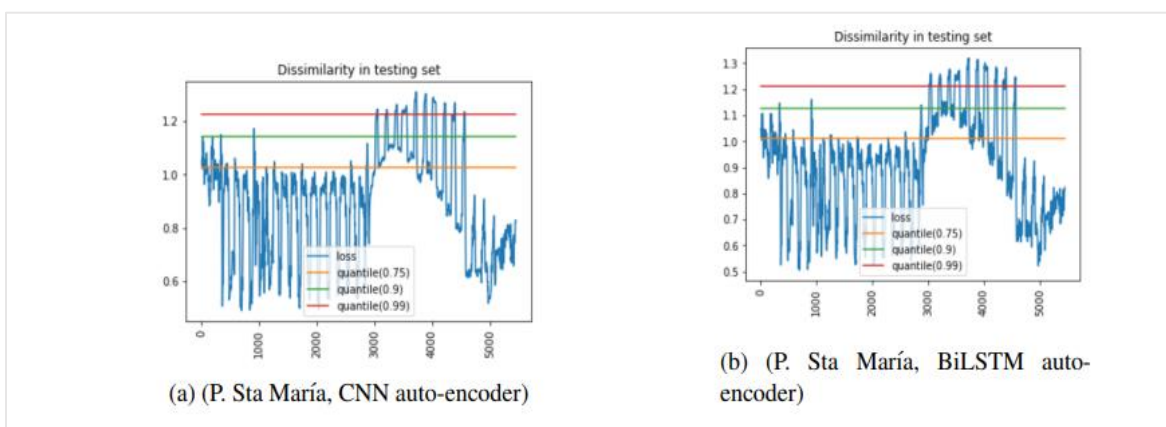
Και για τα δύο μοντέλα, οι Mendez et al. (2021) εξήγαγαν τρία διαφορετικά γραφήματα. Τα πρώτα γραφήματα (βλ. Διάγραμμα 17) είναι ιστογράμματα που αναπαριστούν την κατανομή της απώλειας ανομοιότητας στα δεδομένα εκπαίδευσης. Μια υψηλότερη συγκέντρωση της απώλειας ανομοιότητας με χαμηλή αναπαράσταση στο δεξί άκρο των γραφημάτων δίνει τη δυνατότητα επιλογής υψηλότερων κατωφλίων (thresholds) που ουσιαστικά δεν θα επιστρέφουν ψευδώς θετικά αποτελέσματα για τον εντοπισμό ανωμαλιών.



Διάγραμμα 17: Ιστόγραμμα της κατανομής της απώλειας ανομοιότητας στα δεδομένα εκπαίδευσης

Πηγή: (Mendez et al., 2021)

Τα δεύτερα γραφήματα (βλ. Διάγραμμα 18) αναπαριστούν τη βαθμολογία ανομοιότητας για τις παρατηρήσεις του συνόλου δοκιμών και ποιες από αυτές θεωρούνται ανωμαλίες ανάλογα με το επιλεγμένο κατώφλι. Ο άξονας y αναπαριστά την απώλεια ανομοιότητας και ο άξονας x αναπαριστά κάθε παρατήρηση του συνόλου δοκιμών. Τα επιλεγμένα κατώτατα όρια (αναπαριστώνται με οριζόντιες γραμμές) είναι το 75^ο, 90^ο και 99^ο τεταρτημόριο ντιλο της απώλειας ανομοιότητας στα δεδομένα εκπαίδευσης. Οι τιμές πάνω από κάθε οριζόντια γραμμή είναι οι ανωμαλίες που θεωρούνται σύμφωνα με το αντίστοιχο κατώφλι.



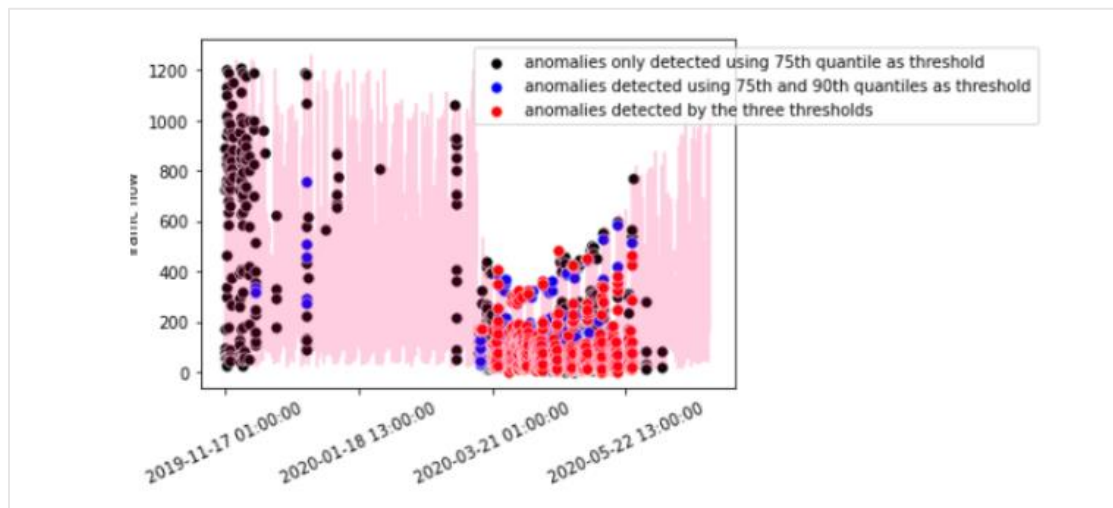
Διάγραμμα 18: Συγκεκριμένα δεδομένα που ανιχνεύονται ως ανώμαλα με βάση το κατώτατο όριο (κατώφλι)

Πηγή: (Mendez et al., 2021)



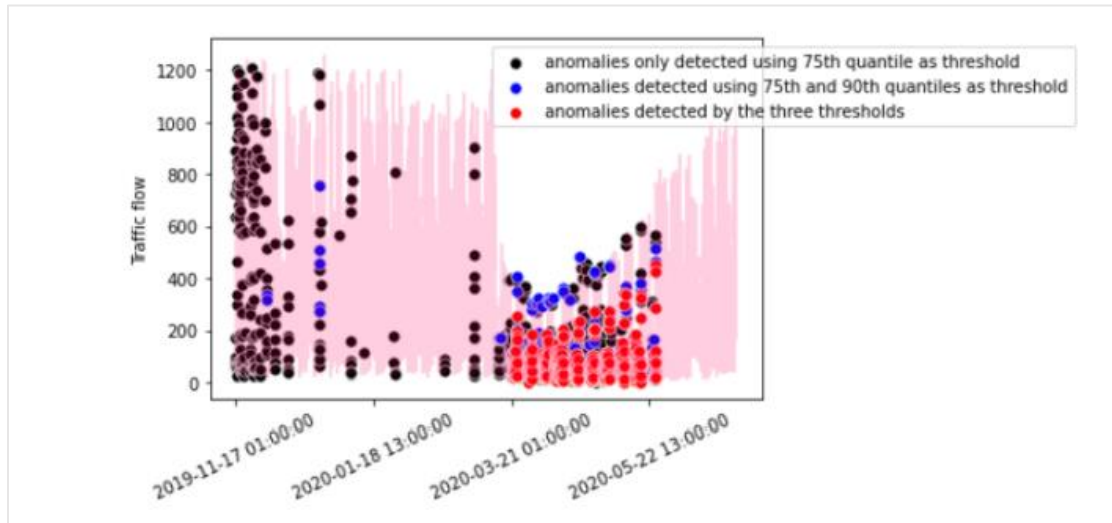
ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

Το τρίτο ζεύγος γραφημάτων (βλ. Διάγραμμα 19 και 20) παρουσιάζει τις ανωμαλίες που ανιχνεύθηκαν (απεικονίζονται στη χρονοσειρά). Ο άξονας x παριστάνει την ημερομηνία, ενώ ο άξονας y παριστάνει την ωριαία ροή κυκλοφορίας. Τα διαφορετικά χρωματισμένα σημεία δείχνουν τις τιμές που θεωρούνται ανώμαλα δεδομένα ανάλογα με το επιλεγμένο κατώφλι. Τα δεδομένα που θεωρούνται ανώμαλα με τη χρήση ενός συγκεκριμένου τεταρτημορίου ως κατώφλι θεωρούνται επίσης ανώμαλα και σε υψηλότερα κβαντίλια. Σε αυτά τα γραφήματα, αναλύονται γραφικά και ανά ημερομηνία, οι λόγοι για τους οποίους διαφορετικές παρατηρήσεις θεωρούνται ως ανωμαλίες.



Διάγραμμα 19 : Συγκεκριμένα δεδομένα που ανιχνεύονται ως ανώμαλα βάσει κατωφλίου με την μέθοδο CNN

Πηγή: (Mendez et al., 2021)



Διάγραμμα 20 : Συγκεκριμένα δεδομένα που ανιχνεύονται ως ανώμαλα βάσει κατώφλιου με την μέθοδο BiLSTM

Πηγή: (Mendez et al., 2021)

6.1.2.2 Σύγκριση

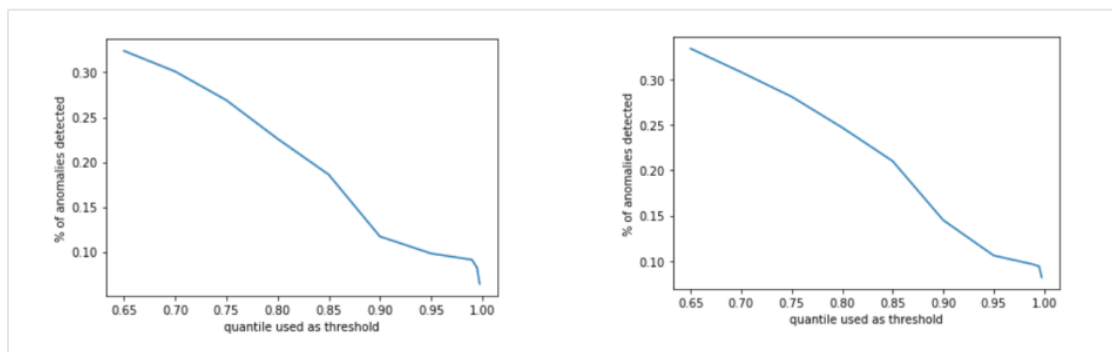
Τα αποτελέσματα του βασικού σεναρίου έδειξαν μεγάλη ομοιότητα μεταξύ των αποτελεσμάτων των δύο μοντέλων, ωστόσο, υπάρχουν ορισμένες αξιοσημείωτες διαφορές. Κατά μέσο όρο, το μέτρο ανομοιότητας είναι μικρότερο στο μοντέλο αυτόματου κωδικοποιητή BiLSTM και τα κατώφλια είναι υψηλότερα στο μοντέλο αυτόματου κωδικοποιητή CNN. Κατά συνέπεια, ο αυτόματος κωδικοποιητής BiLSTM θα ανιχνεύσει περισσότερες πιθανές ακραίες τιμές από τον αυτόματο κωδικοποιητή CNN. Με τη χρήση του 99ου τεταρτημορίου ως κατώφλι (Διαγράμματα 19 και 20), οι ανωμαλίες ανιχνεύονται σε ένα ίδιο εύρος ημερομηνιών που, όλως περιέργως, αντιστοιχούσε στην αρχή των περιορισμών που προκάλεσε η πανδημία COVID-19. Με τη χρήση του 90ου τεταρτημορίου ως κατώφλι, και τα δύο μοντέλα μπορούν επίσης να ανιχνεύσουν ως ανωμαλίες δεδομένα που αντιστοιχούν στις εορταστικές εκδηλώσεις των Χριστουγέννων και στο τελευταίο Σαββατοκύριακο του Νοεμβρίου. Ωστόσο, ο αυτόματος κωδικοποιητής BiLSTM ανιχνεύει περισσότερες ακραίες τιμές στο εύρος των ημερομηνιών που αντιστοιχούν στους περιορισμούς (με το 90^ο τεταρτημόριο ντιλο ως κατώφλι, ο αυτόματος κωδικοποιητής BiLSTM λαμβάνει συνολικά 788 ανωμαλίες έναντι των 638 που λαμβάνει ο αυτόματος κωδικοποιητής CNN). Με τη χρήση του 75ου τεταρτημορίου ως κατώφλι, υπάρχουν ορισμένες



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

ανωμαλίες που ανιχνεύονται και δεν αντιστοιχούν σε καμία σημαντική ημερομηνία ή σε οποιαδήποτε γραφική απόκλιση. Αυτό συνεπάγεται ότι δεν πρόκειται για πραγματικές ανωμαλίες.

Αυτό που φαίνεται από τα αποτελέσματα είναι ότι με τα μικρότερα κβαντίλια ως κατώφλι, ο αριθμός των ανιχνευόμενων ανωμαλιών αυξάνεται. Το διάγραμμα 21 αυτό το αποτέλεσμα και στα δύο μοντέλα. Ενδιαφέρον είναι ότι το ποσοστό των ανωμαλιών παραμένει πρακτικά σταθερό στο μοντέλο αυτόματου κωδικοποιητή CNN από το 90° έως το 99° τεταρτημόριο ντιλο, ενώ στο μοντέλο αυτόματου κωδικοποιητή BiLSTM, αυτή η σταθερότητα συμβαίνει μόνο από το 95° έως το 99° τεταρτημόριο ντιλο. Η παρουσία αυτής της σταθερότητας σε μεγαλύτερο εύρος συνεπάγεται μεγαλύτερη βεβαιότητα για τις ανωμαλίες που ανιχνεύονται με τη χρήση ως κατωφλίου του κατώτερου άκρου του εξεταζόμενου εύρους.



Διάγραμμα 21 : % της ανίχνευσης ανωμαλιών ανά τεταρτημόριο ντιλο που χρησιμοποιείται ως κατώφλι με την μέθοδο CNN (αριστερά) και % ανίχνευσης ανωμαλιών ανά τεταρτημόριο ντιλο που χρησιμοποιείται ως κατώφλι με την μέθοδο BiLSTM (δεξιά)

Πηγή: (Mendez et al., 2021)

Συμπερασματικά, και τα δύο μοντέλα λαμβάνουν παρόμοια αποτελέσματα στο βασικό σενάριο. Ωστόσο, το μοντέλο αυτόματου κωδικοποιητή BiLSTM παρουσιάζει καλύτερη ανακατασκευή της εισόδου και, ως εκ τούτου, μπορεί να ανιχνεύσει περισσότερες ανωμαλίες με το ίδιο τεταρτημόριο ντιλο ως κατώφλι. Επιπλέον, θα πρέπει να σημειωθεί ότι ο αυτόματος κωδικοποιητής CNN χρειάζεται οκτώ



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

ανεξάρτητες μεταβλητές, ενώ ο αυτόματος κωδικοποιητής BiLSTM χρειάζεται μόνο τη χρονοσειρά της μεταβλητής-στόχου.

6.1.2.3 Σενάριο «καθοδηγούμενο»

Στο σενάριο καθοδήγησης τα αποτελέσματα αναλύονται χρησιμοποιώντας τα μοντέλα αυτόματου κωδικοποιητή CNN και BiLSTM για τον ίδιο σταθμό που χρησιμοποιήθηκε και στο βασικό σενάριο. Ως ανώμαλο σύνολο θεωρήθηκε κάθε τιμή υψηλότερη από το 90% και το 99% του εύρους των δεδομένων της χρονοσειράς ροής κυκλοφορίας. Η ιδέα είναι ήταν να δειχθεί πώς η τροποποίηση του αριθμού των ανώμαλων δεδομένων που χρησιμοποιούνται για την εκπαίδευση επηρεάζει την απόδοση του αλγορίθμου. Χρησιμοποιήθηκαν ξανά το 75°, 90° και 99° τεταρτημόριο ντιλο για το κατώφλι απώλειας ανομοιογένειας.

Σε αντίθεση με το βασικό σενάριο, το σύνολο εκπαίδευσης είναι μόνο τα δεδομένα που θεωρούνται «κανονικά», δηλαδή τα σημεία δεδομένων που δεν ανήκουν στο μη κανονικό σύνολο. Έτσι, ιδανικά, τα μοντέλα θα «εκπαιδευτούν» στην αναπαράσταση των κανονικών δεδομένων και θα έχουν μικρή απώλεια ανομοιότητας. Επομένως, αναμένονταν ότι στην εξέταση του συνόλου των δοκιμών (που περιλαμβάνει, ιδίως, μη φυσιολογικά δεδομένα), η απώλεια ανομοιότητας θα είναι μεγαλύτερη.

6.2 Ανίχνευση ανωμαλιών σε κατασκευές γεφυρών με χρήση συνελκτικού νευρωνικού δικτύου (CNN)

Η δομική παρακολούθηση παρέχει πολύτιμες πληροφορίες σχετικά με την κατάσταση της δομικής υγείας (SHM)⁵, οι οποίες είναι χρήσιμες για την ανίχνευση δομικών βλαβών και την αξιολόγηση της δομικής κατάστασης. Ωστόσο, όταν οι αισθητήρες ενός έργου εκτίθενται σε σκληρές περιβαλλοντικές συνθήκες, διάφορες ανωμαλίες που προκαλούνται από αστοχία ή βλάβη του αισθητήρα οδηγούν σε ανωμαλίες των δεδομένων παρακολούθησης. Η χειροκίνητη εξάλειψη των μη φυσιολογικών δεδομένων είναι αναποτελεσματική. Είναι αναποτελεσματικό να αφαιρούνται τα μη φυσιολογικά δεδομένα με χειροκίνητη εξάλειψη λόγω του τεράστιου όγκου τους, που λαμβάνονται από τα συστήματα παρακολούθησης.

⁵ Structural Health Monitoring



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

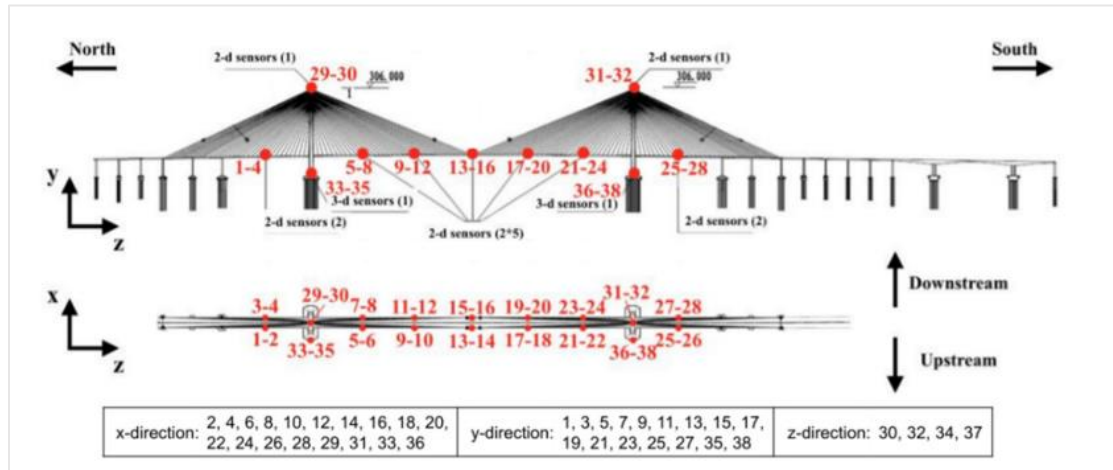
Στην εργασία των Zhang & Lei (2021) προτείνεται μια μέθοδος ανίχνευσης ανωμαλιών δεδομένων που βασίζεται σε σήματα δομικών δονήσεων και ένα συνελκτικό νευρωνικό δίκτυο (CNN), το οποίο μπορεί να εντοπίσει και να εξαλείψει αυτόματα τα μη φυσιολογικά δεδομένα. Πρώτον, το πρόβλημα ανίχνευσης ανωμαλιών μοντελοποιείται ως πρόβλημα ταξινόμησης χρονοσειρών. Για την επεξεργασία της αρχικής χρονοσειράς χρησιμοποιούνται προεπεξεργασία δεδομένων και επαύξηση δεδομένων, συμπεριλαμβανομένης της επέκτασης δεδομένων και της μείωσης της δειγματοληψίας για την κατασκευή νέων δειγμάτων. Για ένα μικρό αριθμό δειγμάτων στο σύνολο δεδομένων, χρησιμοποιούνται μέθοδοι τυχαίας αύξησης των ακραίων τιμών, συμμετρικής αναστροφής και προσθήκης θορύβου για την επέκταση δεδομένων και προστίθενται δείγματα με την ίδια ετικέτα χωρίς αύξηση των αρχικών δειγμάτων. Η μέθοδος μείωσης της δειγματοληψίας με συμμετρική εξαγωγή της μέγιστης τιμής και της ελάχιστης τιμής ταυτόχρονα μπορεί να μειώσει αποτελεσματικά τη διαστατικότητα του δείγματος εισόδου, διατηρώντας παράλληλα τα χαρακτηριστικά των δεδομένων στο μεγαλύτερο βαθμό. Χρησιμοποιώντας ρύθμιση των υπερπαραμέτρων των βαρών ταξινόμησης, το CNN είναι πιο αποτελεσματικό στην αντιμετώπιση μη ισορροπημένων συνόλων εκπαίδευσης.

Η αποτελεσματικότητα της προτεινόμενης μεθόδου αποδεικνύεται με την ανίχνευση ανωμαλιών σε δεδομένα επιτάχυνσης σε μια γέφυρα μεγάλου ανοίγματος που βρίσκεται στην Κίνα. Για το πρόβλημα ανίχνευσης ανωμαλιών που μοντελοποιείται ως πρόβλημα ταξινόμησης χρονοσειρών, η προτεινόμενη μέθοδος μπορεί να εντοπίσει αποτελεσματικά διάφορα ανώμαλα πρότυπα.

6.2.1 Επεξεργασία των δεδομένων και μεθοδολογία

Η υπό συζήτηση έρευνα των Zhang & Lei (2021) χρησιμοποιεί το σύνολο δεδομένων παρακολούθησης της υγείας μιας καλωδιωτής γέφυρας μεγάλου ανοίγματος στην Κίνα. Το κύριο άνοιγμα της γέφυρας έχει μήκος 1088 m και τα δύο πλευρικά ανοίγματα έχουν μήκος 300 m το καθένα, συμπεριλαμβανομένων δύο πύργων ύψους 306 m. Το σύστημα παρακολούθησης της δομικής υγείας της γέφυρας αποτελείται από 38 αισθητήρες. Η θέση τους στη γέφυρα φαίνεται στο παρακάτω σχήμα. Οι αισθητήρες περιλαμβάνουν επιταχυνσιόμετρα, ανεμόμετρα, μετρητές τάσης, παγκόσμια συστήματα εντοπισμού θέσης (GPS) και θερμομέτρα. Στην έρευνα τους οι Zhang & Lei (2021),

χρησιμοποίησαν δεδομένα επιτάχυνσης ενός μήνα (1 Ιανουαρίου-31 Ιανουαρίου 2012) και από τους 38 αισθητήρες του συστήματος SHM για την ανίχνευση ανωμαλιών στα δεδομένα. Η συχνότητα δειγματοληψίας του επιταχυνσιόμετρου είναι 20 Hz.



Σχήμα 4 : Η παρακολουθούμενη γέφυρα και η θέση του επιταχυνσιόμετρου στο σώμα της γέφυρας και στον πύργο

Πηγή: (Zhang & Lei, 2021)

Τα αρχικά συνεχή δεδομένα μέτρησης χωρίζονται σε χρονικές περιόδους μιας ώρας και σε μια χρονική περίοδο ενός μήνα, μέσω της μεθόδου των μη επικαλυπτόμενων παραθύρων. Λαμβάνονται 744 δεδομένα μέτρησης χρονοσειράς από κάθε αισθητήρα, ώστε να προκύψουν συνολικά 28.272 (744×38) δεδομένα. Οι διαστάσεις ενός μεμονωμένου σημείου δεδομένων είναι 1×72.000 . Ο πίνακας 2 περιγράφει την ποσότητα και τα χαρακτηριστικά των κανονικών δεδομένων και των έξι τύπων μη κανονικών δεδομένων. Κάθε σημείο δεδομένων έχει μια πραγματική ετικέτα κατηγορίας. Τα κανονικά δεδομένα μέτρησης χρονοσειρών σημειώνονται ως 1 στον πίνακα και τα άλλα έξι ανώμαλα πρότυπα δεδομένων σημειώνονται ως 2-7. Από τον πίνακα φαίνεται ότι σχεδόν το 52% των δεδομένων είναι ανώμαλα. Η «τάση» είναι το κύριο ανώμαλο μοτίβο που αποτελεί το 20% του συνόλου των δεδομένων, ακολουθούμενο από τα «ελλείποντα» και «τετράγωνα», καθένα από τα οποία αντιπροσωπεύει περίπου 10%. Από την άλλη πλευρά, το «ακραία» αντιπροσωπεύει μόνο το 1,9% του συνόλου δεδομένων, ακολουθούμενο από το «παρέκκλιση», το οποίο αντιπροσωπεύει το 2,4% των δεδομένων.



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

Πίνακας 2: Περιγραφή του κάθε τύπου του μοτίβου των δεδομένων

Πηγή: (Zhang & Lei, 2021)

No.	Μοτίβα ανωμαλίας δεδομένων	Περιγραφή	Ποσότητα
1	Κανονικά	Η χρονική απόκριση είναι κανονική καμπύλη ταλάντωσης- η απόκριση συχνότητας μοιάζει με κορυφή (μπορεί να διαφέρει μεταξύ των γεφυρών)	135715 (48%)
2	Ελλείποντα	Το μεγαλύτερο μέρος/όλη η χρονική απόκριση λείπει, γεγονός που καθιστά την απόκριση χρόνου και συχνότητας μηδενική.	2942 (10.4%>)
3	Δευτερεύοντα	Σε σχέση με τα κανονικά δεδομένα του αισθητήρα, το πλάτος είναι πολύ μικρό στο πεδίο του χρόνου των ελαστικών	1775 (6.3%)
4	Δεδομένα ακραίων τιμών	Μία ή περισσότερες ακραίες τιμές εμφανίζονται στη χρονική απόκριση	527 (1.9%)
5	Τετράγωνο	Η χρονική απόκριση είναι σαν τετραγωνικό κύμα	2996 (10.6%ο)
6	Τάση	Τα δεδομένα έχουν μια προφανή τάση στο πεδίο του χρόνου και μια προφανή τιμή αιχμής στο πεδίο της ταχύτητας.	5778 (20.4%.)
7	Δεδομένα παρέκκλισης	Η απόκριση των δονήσεων είναι μη σταθερή, με τυχαία ολίσθηση	679 (2.4%)

6.2.1.1 Συμπλήρωση-επαύξηση δεδομένων

Οι μέθοδοι επαύξησης πρέπει πάντα να επιλέγονται κατάλληλα ανά την περίπτωση που εξετάζεται. Ως εκ τούτου, η παρούσα έρευνα ασχολείται με κάθε ώρα, δηλαδή με όλο το μήκος του δείγματος. Η βελτίωση των δεδομένων περιλαμβάνει δύο βήματα:

1. Διεύρυνση των δεδομένων ενός μικρού αριθμού δειγμάτων δεδομένων
2. Μείωση της δειγματοληψίας όλων των δειγμάτων

Η επέκταση δεδομένων εφαρμόζεται σε μικρό αριθμό δειγμάτων, δηλαδή στα ακραία και στα δεδομένα παρέκκλισης, στην αριθμητική προσομοίωση. Δεν χρειάζεται να επεκταθούν όλα τα ανώμαλα δείγματα. Τα ακραία δεδομένα μπορούν να οριστούν ως



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

μεμονωμένα σημεία των κανονικών δεδομένων των οποίων το πλάτος υπερβαίνει κατά πολύ το κανονικό εύρος. Επομένως, για τα ακραία δείγματα χρησιμοποιείται μια μέθοδος επέκτασης δεδομένων που μεγεθύνει τα μεμονωμένα σημεία. Έστω x είναι ένα κανονικό δείγμα $\{x_1, x_2, \dots, x_n\}$, και η προτεινόμενη μέθοδος παρουσιάζεται στην εξίσωση:

$$x(p) = mean + b \cdot range \tag{1.17}$$

Όπου p είναι ένας τυχαίος αριθμός μεταξύ 10 και 60, $mean$ είναι η μέση τιμή του x , β είναι ένας τυχαίος αριθμός μεταξύ -2 και 2, και $range$ είναι η διαφορά μεταξύ της μέγιστης και της ελάχιστης τιμής στο x .

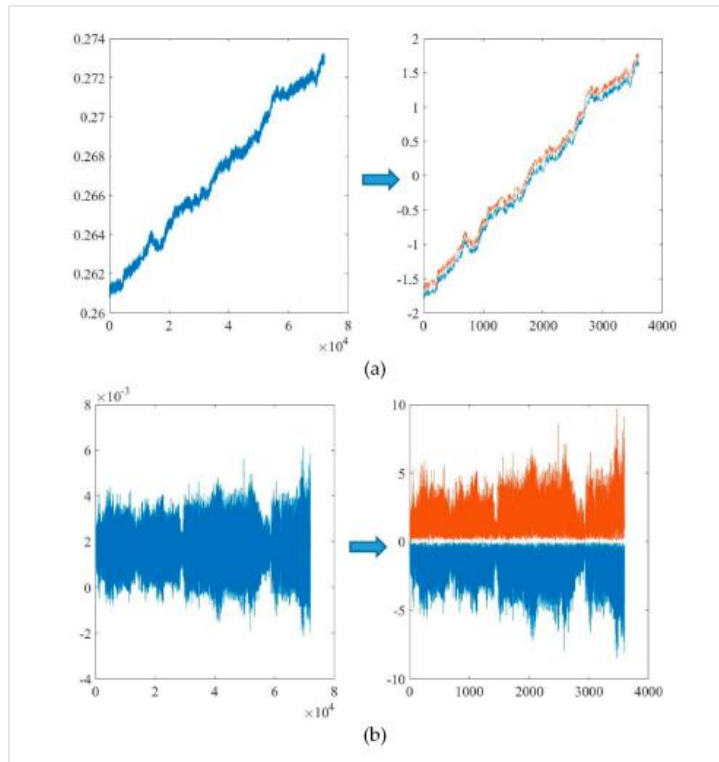
Η μέθοδος της συμμετρικής αναστροφής και της προσθήκης θορύβου χρησιμοποιείται για την επέκταση των δεδομένων παρέκκλισης. Τα δεδομένα παρέκκλισης έχουν τυχαία ολίσθηση προς τα πάνω ή προς τα κάτω. Επομένως, η μέθοδος της συμμετρικής αναστροφής προς τα πάνω και προς τα κάτω μπορεί να κατασκευάσει ένα αποτελεσματικό δείγμα. Για τη χρονοσειρά $\{x_1, x_2, \dots, x_n\}$, η συμμετρική αναστροφή μπορεί να δημιουργήσει μια νέα χρονοσειρά $\{\acute{x}_1, \acute{x}_2, \dots, \acute{x}_n\}$ με τις ίδιες ετικέτες ανωμαλίας, όπου $\acute{x}_i = -x_i$. Διαφορετικοί βαθμοί λευκού θορύβου Gauss προστίθενται στην αρχική ακολουθία για τη δημιουργία περισσότερων δειγμάτων με τις ίδιες ετικέτες ανωμαλίας.

Η διάσταση του δείγματος μιας ώρας είναι 1×72.000 , η οποία είναι σχετικά μεγάλη ως είσοδος του νευρωνικού δικτύου. Ως εκ τούτου, έγινε μείωση της δειγματοληψίας για να μειωθεί η διάσταση του δείγματος, διατηρώντας παράλληλα τα χαρακτηριστικά του δείγματος όσο το δυνατόν περισσότερο, ώστε να αυξηθεί η απόδοση του νευρωνικού δικτύου. Το άνω και το κάτω όριο ενός δείγματος είναι και τα δύο χρήσιμα χαρακτηριστικά. Ως εκ τούτου, χρησιμοποιείται μια μέθοδος μείωσης της δειγματοληψίας που χρησιμοποιεί ένα ολισθαίνον παράθυρο για τη συμμετρική εξαγωγή της μέγιστης και της ελάχιστης τιμής. Όλα τα δείγματα 1×72.000 λαμβάνονται σε δειγματοληψία προς τα κάτω σε όλο το μήκος του δείγματος. Επιλέγεται ένα μέγεθος βήματος, το οποίο είναι 20 και οι μέγιστες και ελάχιστες τιμές στην ακολουθία αφαιρούνται για κάθε σημείο δειγματοληψίας του μεγέθους βήματος.



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

Επομένως, μετά την επεξεργασία καθενός από τα 72.000 δείγματα, θα προκύψει ένα μέγεθος δείγματος 2×3600 . Το διάγραμμα σύγκρισης ορισμένων παραδειγμάτων πριν και μετά τη μείωση της δειγματοληψίας παρουσιάζεται στο Διάγραμμα 22α,β. Ο οριζόντιος άξονας αντιπροσωπεύει τον αριθμό των σημείων δειγματοληψίας και ο κατακόρυφος άξονας αντιπροσωπεύει το πλάτος επιτάχυνσης σε m/s^2 .



Διάγραμμα 22: Συγκριτικό διάγραμμα πριν και μετά την επέκταση των δεδομένων. (α) Επέκταση δεδομένων για ακραία δεδομένα (β) Επέκταση δεδομένων για δεδομένα παρέκκλισης

Πηγή: (Zhang & Lei, 2021)

6.2.1.2 Δεδομένα παρακολούθησης της γέφυρας

Σύμφωνα την προτεινόμενη μέθοδο των Zhang & Lei (2021) η ανίχνευση ανωμαλιών πραγματοποιείται στο σύνολο δεδομένων παρακολούθησης της γέφυρας. Πρώτον, πραγματοποιείται η προ-επεξεργασία των δεδομένων σε όλα τα αρχικά δείγματα, διαγράφονται οι ελλείπουσες τιμές και τα δείγματα τυποποιούνται. Προκειμένου να ελεγχθεί η ικανότητα γενίκευσης του μοντέλου, το σύνολο δεδομένων χωρίζεται σε σύνολα εκπαίδευσης και δοκιμής και το 80% των δειγμάτων επιλέγεται τυχαία ως



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

σύνολο εκπαίδευσης. Το μέγεθος του συνόλου εκπαίδευσης είναι 22.616. Το 20% των δειγμάτων επιλέγεται τυχαία ως σύνολο δοκιμής και το μέγεθος του συνόλου δοκιμής είναι 5656.

Στην διαδικασία η κατασκευή ενός ισορροπημένου συνόλου εκπαίδευσης διαφόρων κατηγοριών είναι επωφελής κυρίως για την εκπαίδευση του νευρωνικού δικτύου. Η επέκταση των δεδομένων πραγματοποιείται για τον μικρό αριθμό ανωμαλιών στο σύνολο προσδιορισμού, δηλαδή τις ακραίες τιμές και τα δεδομένα παρέκκλισης. Τα κανονικά δείγματα σε όλα τα σύνολα εκπαίδευσης επεκτείνονται σε δείγματα ακραίων τιμών με μεγέθυνση μεμονωμένων σημείων. Ο γκαουσιανός κατανεμημένος θόρυβος με 2%, 3%, 4%, 5%, 6%, 7% και 8% τυπική απόκλιση από τα σήματα προστίθεται σε κάθε δείγμα ολίσθησης μία φορά και αναστρέφεται συμμετρικά μία φορά για να ληφθεί 8 φορές ο αριθμός των δειγμάτων παρέκκλισης. Ως εκ τούτου, λαμβάνονται επιπλέον 10,86 ($13,575 \cdot 80\%$) δείγματα ακραίων τιμών και 4345 ($679 \cdot 80\% \cdot 8$) δείγματα ολίσθησης. Μετά την προσθήκη στο σύνολο εκπαίδευσης, το νέο μέγεθος του συνόλου εκπαίδευσης είναι 37.821 ($22.616 + 10.860 + 4345$). Η μείωση της δειγματοληψίας εφαρμόζεται στο σύνολο των δεδομένων και στα νέα δείγματα εκπαίδευσης και η διαστατικότητα των δειγμάτων μειώνεται από $1 \cdot 72000$ σε $2 \cdot 3600$ διατηρώντας ωστόσο τα περισσότερα από τα χαρακτηριστικά τους.

6.2.1.3 1D-CNN

Ένα συνελκτικό νευρωνικό δίκτυο (CNN) αποτελείται συνήθως από ένα στρώμα εισόδου, ένα συνελκτικό στρώμα (Conv), ένα στρώμα συγκέντρωσης (Pooling), ένα πυκνό στρώμα (Dense) και ένα στρώμα εξόδου. Στην αρχιτεκτονική του CNN, τα πρώτα στρώματα συνήθως εναλλάσσονται μεταξύ των στρωμάτων συνελκτικής και των στρωμάτων συγκέντρωσης και τα τελευταία στρώματα κοντά στο στρώμα εξόδου αποτελούνται από πυκνά στρώματα. Το CNN είναι ένα μοντέλο μεθόδου μάθησης από άκρο σε άκρο, το οποίο μπορεί να χρησιμοποιήσει τον υπάρχοντα αλγόριθμο κλίσης με επίβλεψη για την εκπαίδευση του μοντέλου. Για προβλήματα επεξεργασίας χρονοσειρών, το αποτέλεσμα ενός μονοδιάστατου νευρωνικού δικτύου συνελίξεων (1D-CNN) μπορεί να είναι συγκρίσιμο με ένα επαναλαμβανόμενο νευρωνικό δίκτυο (RNN)⁶ και το υπολογιστικό κόστος είναι πολύ μικρότερο. Για απλές εργασίες, όπως η

⁶ Recurrent neural network

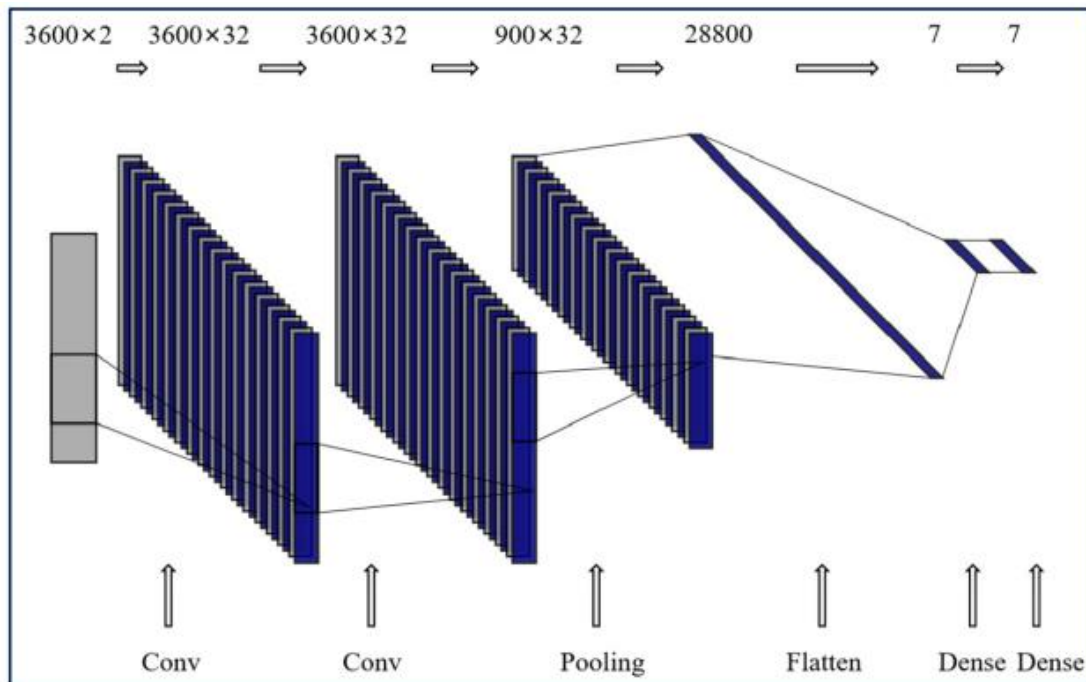


ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

ταξινόμηση χρονοσειρών, ένα μικρό μονοδιάστατο συνελκτικό δίκτυο μπορεί να αντικαταστήσει πλήρως το RNN και εκτελείται ταχύτερα (Chollet, 2018).

Ανεξάρτητα από το αν χρησιμοποιείται μονοδιάστατη ή δισδιάστατη συνέλιξη, τα νευρωνικά δίκτυα συνέλιξης έχουν παρόμοια δομή. Η δομή ξεκινά με μια στοίβα από στρώματα συνέλιξεων και συγκέντρωσης, και στη συνέχεια συνδέεται με ένα επίπεδο στρώμα για τη μετατροπή δισδιάστατων χαρακτηριστικών σε μονοδιάστατη έξοδο, και στη συνέχεια μπορούν να προστεθούν πολλαπλά πυκνά στρώματα για ταξινόμηση ή παλινδρόμηση. Ωστόσο, υπάρχει μια μικρή διαφορά μεταξύ τους: τα μονοδιάστατα νευρωνικά δίκτυα συνέλιξης μπορούν να χρησιμοποιούν μεγαλύτερους πυρήνες συνέλιξης (Chollet, 2018). Για παράδειγμα, για ένα δισδιάστατο στρώμα συνέλιξης, ένας πυρήνας συνέλιξης 3×3 περιέχει $3 \times 3 = 9$ διανύσματα συνέλιξης- ωστόσο, για ένα μονοδιάστατο στρώμα συνέλιξης, ένας πυρήνας συνέλιξης μεγέθους 3 περιέχει μόνο 3 διανύσματα συνέλιξης. Επομένως, ένας μονοδιάστατος πυρήνας συνέλιξης μεγαλύτερος ή ίσος με 9 μπορεί εύκολα να χρησιμοποιηθεί

Για να δημιουργηθεί η αρχιτεκτονική 1D-CNN, δύο μονοδιάστατα στρώματα συνελκτικής ανάλυσης στοιβάζονται για να αποκτήσουν τα βαθιά χαρακτηριστικά του δείγματος πιο αποτελεσματικά, και ένα επίπεδο και δύο πυκνά στρώματα συνδέονται για να μετατρέψουν τα δισδιάστατα χαρακτηριστικά σε μονοδιάστατη έξοδο. Το τελευταίο στρώμα του δικτύου χρησιμοποιεί τον πολυταξινομητή softmax. Εν συντομία, softmax είναι η τιμή που αντιστοιχίζει την έξοδο του προηγούμενου στρώματος στο (0,1) μέσω της συνάρτησης softmax. Το άθροισμα αυτών των τιμών είναι 1, το οποίο μπορεί να εκληφθεί ως πιθανότητα. Ο κόμβος -με τη μεγαλύτερη πιθανότητα- επιλέγεται ως ο αποκλειόμενος ανώμαλος τύπος δεδομένων. Η δομή του δικτύου παρουσιάζεται στο Σχήμα 5.



Σχήμα 5 : Σχηματική αναπαράσταση της προ προτεινόμενης αρχιτεκτονικής CNN

Πηγή: (Zhang & Lei, 2021)

Το μέσο τετραγωνικό σφάλμα (MSE)⁷ ως συνάρτηση απώλειας για την εκπαίδευση και την επικύρωση μπορεί να είναι εκφραστεί ως εξής:

$$MSE = \frac{1}{N} \sum_{i=1}^N (Y_i - Y_{0,i})^2$$

όπου το Y αντιπροσωπεύει την προβλεπόμενη από το μοντέλο τιμή και το Y_0 δείχνει την πραγματική τιμή. Το N είναι το σύνολο όλων των δειγμάτων.

6.2.2 Αποτελέσματα του μοντέλου CNN

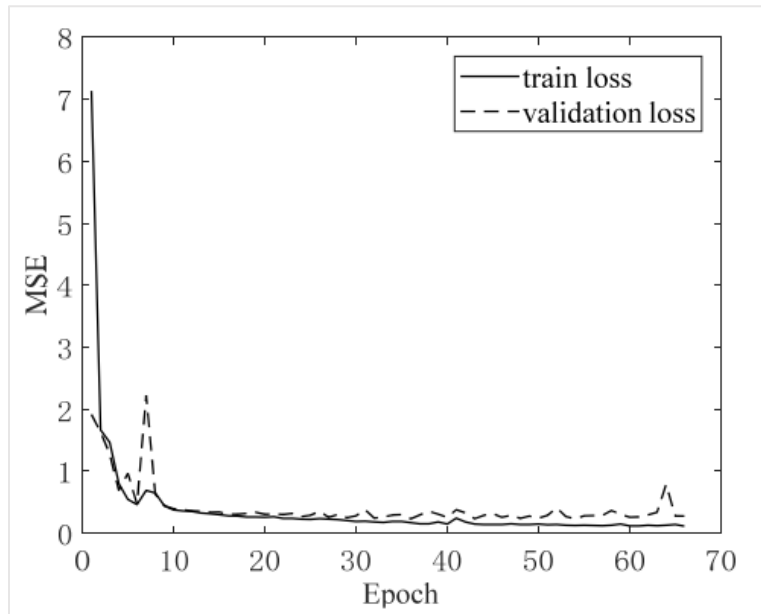
Κατά τη διαδικασία εκπαίδευσης, το σύνολο εκπαίδευσης χωρίζεται σε 12,5% ως σύνολο επαλήθευσης. Κατά τη διάρκεια της διαδικασίας εκπαίδευσης, παρακολουθείται η απώλεια εκπαίδευσης και η απώλεια επικύρωσης (MSE), καθώς και η ακρίβεια εκπαίδευσης και η ακρίβεια επαλήθευσης (Accuracy). Η μεταβολή της

⁷ Mean Squared Error



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

συνάρτησης απώλειας και η μεταβολή της ακρίβειας παρουσιάζονται στα Διαγράμματα 23 και 24.



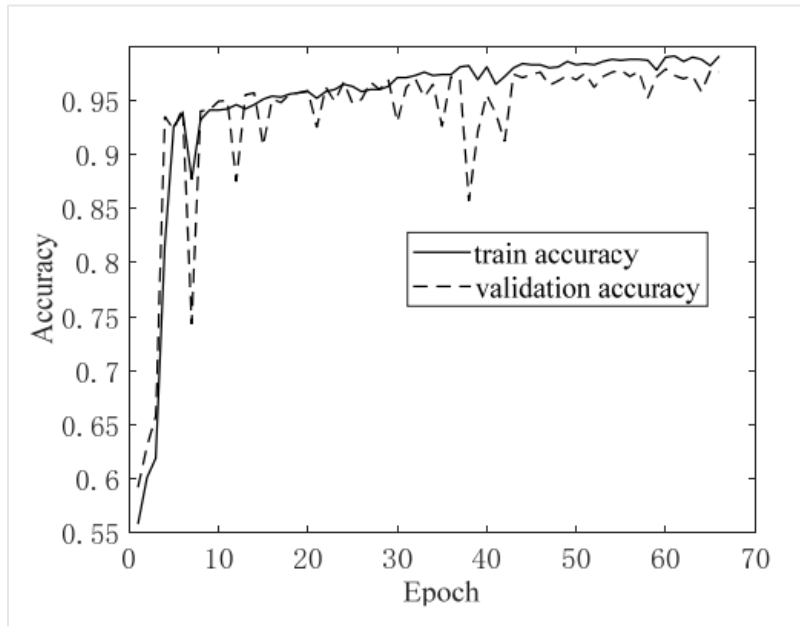
Διάγραμμα 23 : Καμπύλη απωλειών εκπαίδευσης και επικύρωσης

Πηγή: (Zhang & Lei, 2021)

Φαίνεται στο μοντέλο των Zhang & Lei (2021) ότι η συνολική τιμή των απωλειών παρουσιάζει πτωτική τάση (Διάγραμμα 23, ενώ η συνολική ακρίβεια παρουσιάζει ανοδική τάση (Διάγραμμα 24) . Το εύρος είναι μεγάλο στην αρχή της εκπαίδευσης, υποδεικνύοντας ότι ο ρυθμός μάθησης είναι κατάλληλος. Μετά τη σταθεροποίηση της τιμής απώλειας και της ακρίβειας, η τελική ακρίβεια εκπαίδευσης και επικύρωσης έφτασε σε ποσοστό άνω του 95% (Διάγραμμα 24).



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ



Διάγραμμα 24: Καμπύλη ακρίβειας εκπαίδευσης και επικύρωσης

Πηγή: (Zhang & Lei, 2021)

Στη στατιστική ανάλυση δυαδικών ή πολλαπλών ταξινομήσεων, η ακρίβεια, η ανάκληση και το σκορ F1 είναι μέτρα της ακρίβειας των αποτελεσμάτων της ταξινόμησης. Το τελευταίο (F1) είναι ο αρμονικός μέσος όρος των δύο πρώτων. Η ανάκληση σχετίζεται με το δείγμα, δηλαδή πόσες θετικές τιμές του δείγματος έχουν προβλεφθεί σωστά. Για παράδειγμα στα δείγματα τύπου ελλειπόντων τιμών υπάρχουν συνολικά 603 δείγματα. Εάν 602 προβλέπονται σωστά, η ανάκληση είναι $602/603 = 99,83\%$. Η ακρίβεια είναι σχετική του αποτελέσματος της πρόβλεψης και δείχνει πόσα από τα δείγματα των οποίων οι προβλέψεις είναι θετικές είναι σωστά. Λαμβάνοντας ως παράδειγμα τα δείγματα φυσιολογικού τύπου, συνολικά 2590 δείγματα προβλέπεται ότι είναι φυσιολογικοί τύποι. Εάν 2542 προβλέψεις είναι σωστές, η ακρίβεια είναι $2542/2590 = 98,15\%$. Οι δείκτες ανάκλησης και ακρίβειας είναι μερικές φορές αντιφατικοί. Εάν χρησιμοποιείται ένας ολοκληρωμένος δείκτης για να εκφράσει τα αποτελέσματα της ανάκλησης και της ακρίβειας, η πιο συνηθισμένη μέθοδος θα πρέπει να είναι η βαθμολογία F1:

$$F_1 = 2 \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} \times 100\%$$



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

(1.18)

Πίνακας 3: Αποτελέσματα της πρόβλεψης του συνόλου δοκιμής

Πηγή: (Zhang & Lei, 2021)

Πρόβλεψη μοτίβων δεδομένων										
		1	2	3	4	5	6	7	Σύνολο	Ανάκληση (%)
Μοτίβο πραγματικών δεδομένων	1-κανονικά	2542	0	60	72	14	0	0	2688	94,57
	2-ελλείποντα	0	602	1	0	0	0	0	603	99,83
	3- δευτερεύοντα	26	0	326	8	0	0	0	360	90,56
	4-ακραία	21	0	1	84	0	0	0	106	79,24
	5-τετράγωνα	1	0	1	2	612	0	0	616	99,35
	6-τάσης	0	0	2	0	0	1110	35	1147	96,77
	7-παρέκλισης	0	0	0	0	0	13	123	136	90,44
	Σύνολο	2590	602	391	166	626	1123	158	5656	95,45
	Ακρίβεια (%)	98.15	100.0	83.38	50.60	97.76	98.84	77.85		
	F1 score	0.96	1.00	0.87	0.62	0.99	0.98	0.84		

Από τα αποτελέσματα διαπιστώνεται ότι η προτεινόμενη μέθοδος μπορεί να εντοπίσει αποτελεσματικά διάφορα μοτίβα δεδομένων. Η ανάκληση των κατηγοριών κανονικών, ελλειπών, δευτερευουσών, τετραγωνικών, τάσης και παρέκκλισης μπορεί να φτάσει πάνω από το 90%. Εκτός από το χαμηλό σκορ F1 των ακραίων δεδομένων και των δεδομένων παρέκκλισης, οι άλλοι τύποι είναι όλοι υψηλοί. Ένας μικρός αριθμός δευτερευόντων δεδομένων ταξινομείται στην κατηγορία των κανονικών δεδομένων. Ορισμένα δείγματα ακραίων τιμών ταξινομούνται στην κανονική κατηγορία και μερικά στην κατηγορία δευτερευόντων. Το ακραίο δείγμα μπορεί να έχει μόνο λίγες κορυφές και τα περισσότερα χαρακτηριστικά του ακραίου δείγματος είναι πολύ παρόμοια με το κανονικό δείγμα και το χαρακτηριστικό που είναι πολύ μικρό θα χαθεί



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

στη διαδικασία συνέλιξης. Η τάση και η παρέκκλιση συγχέονται εν μέρει, πιθανώς επειδή και οι δύο έχουν κεκλιμένα χαρακτηριστικά.

Κεφάλαιο 7^ο: Συμπεράσματα

Καθεμία από τις πολυάριθμες τεχνικές ανίχνευσης ανωμαλιών που συζητήθηκαν προηγουμένως έχει τα μοναδικά της πλεονεκτήματα και αδυναμίες. Είναι σημαντική η επιλογή της κατάλληλης τεχνικής ανίχνευσης ανωμαλιών για ένα δεδομένο πρόβλημα ανίχνευσης ανωμαλιών. Δεδομένης της πολυπλοκότητας του χώρου των προβλημάτων, δεν είναι εφικτό να παρέχεται μια τέτοια κατανόηση για κάθε πρόβλημα ανίχνευσης ανωμαλιών.

Η υπολογιστική πολυπλοκότητα μιας τεχνικής ανίχνευσης ανωμαλιών είναι μια βασική πτυχή, ειδικά όταν η τεχνική εφαρμόζεται σε έναν πραγματικό τομέα. Ενώ οι τεχνικές που βασίζονται στην ταξινόμηση, την ομαδοποίηση και τις στατιστικές τεχνικές έχουν ακριβούς χρόνους εκπαίδευσης, η δοκιμή είναι συνήθως γρήγορη. Συχνά αυτό είναι αποδεκτό, δεδομένου ότι τα μοντέλα μπορούν να εκπαιδευτούν εκτός σύνδεσης, ενώ η δοκιμή απαιτείται να γίνεται σε πραγματικό χρόνο. Αντίθετα, τεχνικές όπως οι τεχνικές που βασίζονται στον πλησιέστερο γείτονα, οι θεωρητικές της πληροφορίας, οι οποίες δεν έχουν φάση εκπαίδευσης, έχουν ακριβή φάση δοκιμής, η οποία μπορεί να αποτελέσει περιορισμό σε ένα πραγματικό περιβάλλον.

Οι τεχνικές ανίχνευσης ανωμαλιών συνήθως υποθέτουν ότι οι ανωμαλίες στα δεδομένα είναι σπάνιες σε σύγκριση με τις κανονικές περιπτώσεις. Αν και η υπόθεση αυτή είναι γενικά αληθής, οι ανωμαλίες δεν είναι πάντα σπάνιες. Για παράδειγμα, όταν το θέμα είναι η ανίχνευση σκουληκιών σε δίκτυα υπολογιστών, η ανώμαλη κίνηση (σκουλήκι) είναι στην πραγματικότητα πιο συχνή από την κανονική κίνηση. Οι μη επιβλεπόμενες τεχνικές δεν είναι κατάλληλες για τέτοια μαζική ανίχνευση ανωμαλιών. Για την ανίχνευση μαζικών ανωμαλιών μπορούν να εφαρμοστούν τεχνικές που λειτουργούν με εποπτεία ή ημι-υπο εποπτεία (Sun et al., 2007).

Για πιο σύνθετα σύνολα δεδομένων, διαφορετικοί τύποι τεχνικών αντιμετωπίζουν διαφορετικές προκλήσεις. Οι τεχνικές που βασίζονται στον πλησιέστερο γείτονα και στην ομαδοποίηση υποφέρουν όταν ο αριθμός των διαστάσεων είναι μεγάλος, επειδή τα μέτρα απόστασης σε μεγάλο αριθμό διαστάσεων δεν είναι σε θέση να



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

διαφοροποιήσουν μεταξύ κανονικών και ανώμαλων περιπτώσεων. Οι τεχνικές που βασίζονται στην ταξινόμηση μπορούν να αποτελέσουν καλύτερη επιλογή για ένα σύνθεστο σύνολο δεδομένων, αλλά για να πιο αποτελεσματικές, οι τεχνικές που βασίζονται στην ταξινόμηση απαιτούν ετικέτες τόσο για τις κανονικές όσο και για τις ανώμαλες περιπτώσεις, οι οποίες συχνά δεν είναι διαθέσιμες.

Σε αυτή την επισκόπηση συζητήθηκαν διάφοροι τρόποι με τους οποίους το πρόβλημα της ανίχνευσης ανωμαλιών έχει διατυπωθεί στη βιβλιογραφία και πραγματοποιήθηκε μια προσπάθεια επισκόπησης του μεγάλου όγκου βιβλιογραφίας σχετικά με τις διάφορες τεχνικές. Για κάθε κατηγορία τεχνικών ανίχνευσης ανωμαλιών, προσδιορίστηκε μια μοναδική υπόθεση σχετικά με την έννοια των κανονικών και των ανώμαλων δεδομένων. Κατά την εφαρμογή μιας δεδομένης τεχνικής σε έναν συγκεκριμένο τομέα, οι υποθέσεις αυτές μπορούν να χρησιμοποιηθούν ως κατευθυντήριες γραμμές για την αξιολόγηση της αποτελεσματικότητας της τεχνικής στον εν λόγω τομέα. Ιδανικά, μια ολοκληρωμένη έρευνα σχετικά με την ανίχνευση ανωμαλιών θα πρέπει να επιτρέπει στον αναγνώστη όχι μόνο να κατανοήσει τα κίνητρα πίσω από τη χρήση μιας συγκεκριμένης τεχνικής ανίχνευσης ανωμαλιών, αλλά και να παρέχει μια συγκριτική ανάλυση των διαφόρων τεχνικών.

Εκτός όμως από την θεωρητική ανασκόπηση των διαφόρων τεχνικών για την ανίχνευση των ανωμαλιών, στην παρούσα εργασία συμπεριλήφθηκαν και δύο περιπτώσεις ειδικότερου ενδιαφέροντος. Μελετήθηκαν δύο εφαρμογές τεχνικών βαθιάς μάθησης σε δύο πολύ συγκεκριμένα προβλήματα, που ανακύπτουν σε πραγματικές συνθήκες. Η πρώτη εφαρμογή που μελετήθηκε αφορούσε το ζήτημα της κυκλοφορικής ροής, όπου μια τεχνική ανίχνευσης ανωμαλιών μπορεί να παράσχει δεδομένα με υψηλή χρησιμότητα σχετικά με τις ώρες εντός της ημέρας άλλα και τις περιόδους όπου κορυφώνεται η κίνηση ή μεταβάλλονται τα μοτίβα του κυκλοφοριακού δικτύου (Mendez et al., 2022). Η δεύτερη εφαρμογή, ενέπιπτε στον κλάδο της πολιτικής μηχανικής και αφορούσε το πολύ καίριο ζήτημα της παρακολούθησης των διάφορων δεδομένων που καταγράφονται από τους αισθητήρες στις γέφυρες (Zhang & Lei, 2021). Το συγκεκριμένο ζήτημα είναι καίριας σημασίας διότι πρέπει να υπάρχει η δυνατότητα καθορισμού του είδους των ανώμαλων δεδομένων τα οποία δημιουργούν



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

πιθανά προβλήματα στατικότητας σε σχέση με τα υπόλοιπα ανώμαλα δεδομένα που μπορεί να παρουσιάζονται λόγω προβλημάτων στους αισθητήρες (πχ φθορά).

Στην μελέτη των Mendez et al. (2022) τα δεδομένα προήλθαν από έναν σταθμό κυκλοφορίας στην πόλη της Μαδρίτης και είχαν διάρκεια από τον Ιανουάριο του 2018 έως τον Ιούλιο του 2020. Στην εν λόγω μελέτη χρησιμοποιήθηκαν δύο διαφορετικές τεχνικές βαθιάς μάθησης η CNN και η BiLSTM, που εφαρμόστηκαν σε δύο διαφορετικά σενάρια. Στο πρώτο σενάριο που ήταν το «βασικό» έγινε χρήση όλων των δεδομένων για την εκπαίδευση των μοντέλων (κανονικά και μη κανονικά δεδομένα), ενώ στο δεύτερο που ήταν «καθοδηγούμενο» το μοντέλο εκπαιδεύτηκε μόνο στα κανονικά δεδομένα. Τα κανονικά δεδομένα του «καθοδηγούμενου» σεναρίου ήταν αυτά που βρίσκονταν κάτω από ένα συγκεκριμένο τεταρτημόριο ντιλο. Η επιτυχία του «καθοδηγούμενου» σεναρίου για αμφότερες τις τεχνικές ήταν περιορισμένη. Η τεχνική CNN αναγνώρισε ανωμαλίες όλων των τύπων (παγκόσμιες, πλαισίου και συλλογικές) αλλά δεν μπορούσε να αναγνωρίσει έναν ικανοποιητικό αριθμό από αυτές, ενώ η BiLSTM αναγνώρισε μόνο ανωμαλίες παγκόσμιου τύπου. Στον αντίποδα στο «βασικό» σενάριο τα αποτελέσματα ήταν πολύ καλύτερα και οι δύο τεχνικές αναγνώρισαν όλων των τύπων τις ανωμαλίες.

Για την μελέτη Zhang & Lei (2021) χρησιμοποιήθηκαν δομικά δεδομένα ενός μήνα, από επιταχυνσιόμετρα, ανεμόμετρα, μετρητές τάσης κ.λπ., που αφορούσαν μια γέφυρα μεγάλου ανοίγματος στην Κίνα. Στην εργασία αυτή χρησιμοποιήθηκε επίσης η μέθοδος βαθιάς μάθησης CNN για τον εντοπισμό των ανωμαλιών στις χρονοσειρές των δεδομένων. Αυτό που διαφοροποιεί τις δύο μελέτες που εξετάστηκαν, αν και πρόκειται για την ίδια τεχνική CNN, είναι ότι στην μελέτη των Zhang & Lei (2021) ακολούθησαν διαφορετική προ-επεξεργασία των δεδομένων. Ειδικότερα, τα δεδομένα επεκτάθηκαν και μειώθηκε η δειγματοληψία, δημιουργώντας έτσι νέα δείγματα. Η επέκταση των δεδομένων αφορούσε ένα μικρό αριθμό δειγμάτων, τα δεδομένα ακραίων τιμών και τα δεδομένα παρέκκλισης. Τα πρώτα είναι οι ακραίες τιμές που εμφανίζονται στην χρονική απόκριση και τα δεύτερα είναι μη στάσιμα δεδομένα που προκύπτουν από την παρέκκλιση της σειράς. Στα μεν η επέκταση έγινε μέσα από την μεγέθυνση των μεμονωμένων σημείων στα δε με την μέθοδο της συμμετρικής αναστροφής και της



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

προσθήκης θορύβου. Επίσης, έγινε μείωση του δείγματος με συμμετρική εξαγωγή της μέγιστης και της ελάχιστης τιμής.

Τα αποτελέσματα της έρευνας έδειξαν γενικά ότι η μέθοδος μπορεί να εντοπίσει με αποτελεσματικό τρόπο τα διάφορα μοτίβα των δεδομένων, που είναι και το ζητούμενο καθώς ο στόχος της έρευνας ήταν το πρόβλημα της ανίχνευσης των ανωμαλιών να μετατραπεί σε πρόβλημα κατηγοριοποίησης των σειρών. Η ανάκληση ήταν αρκετά υψηλή σε >90% σχεδόν σε όλες τις κατηγορίες και τα σκορ F1 ήταν χαμηλά στην περίπτωση των ακραίων τιμών και των δεδομένων παρέκκλισης. Αυτό δείχνει πως για ανώμαλα δεδομένα με πολύ δυσδιάκριτα χαρακτηριστικά, όπως τα ακραία δεδομένα, υπάρχουν ακόμη πολλά περιθώρια βελτίωσης της ακρίβειας αναγνώρισης.

Γενικά μπορεί να εξαχθεί στο συμπέρασμα ότι οι μέθοδοι αυτόματου κωδικοποιητή λειτουργούν καλά για την ανίχνευση ανωμαλιών σε χρονοσειρές. Επιπλέον, λειτουργούν καλά με όλους τους τύπους ανώμαλων δεδομένων, επειδή βασίζονται τη λειτουργία τους στη διαφορά μεταξύ της εισόδου και της ανακατασκευασμένης εξόδου. Ωστόσο, οι τυπικές μέθοδοι βαθιάς μάθησης δεν λειτουργούν το ίδιο καλά, επειδή βασίζονται τη λειτουργία τους στην εύρεση παρόμοιων χαρακτηριστικών μεταξύ των δεδομένων, αλλά οι διαφορετικοί τύποι ανωμαλιών δεν είναι παρόμοιοι μεταξύ τους.

Βιβλιογραφία

Abraham, B. and Box, G. E. P. (1979). Bayesian analysis of some outlier problems in time series. *Biometrika* 66, 2, 229–236.

Aggarwal, C. (2005). On abnormality detection in spuriously populated data streams. In *Proceedings of 5th SIAM Data Mining*. 80.

Agliari, E., Barra, A., Barra, O. A., Fachechi, A., Franceschi Vento, L., & Moretti, L. (2020). Detecting cardiac pathologies via machine learning on heart-rate variability time series and related markers. *Scientific Reports*, 10(1). <https://doi.org/10.1038/s41598-020-64083-4>

Albrecht, S., Busch, J., Kloppenburg, M., Metze, F., and Tavan, P. (2000). Generalized radial basis function networks for classification and novelty detection: self-organization of optional bayesian decision. *Neural Networks* 13, 10, 1075–1093.



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

- Anscombe, F. J. and Guttman, I. (1960). Rejection of outliers. *Technometrics* 2, 2, 123–147.
- Baker, D., Hofmann, T., McCallum, A., and Yang, Y. (1999). A hierarchical probabilistic model for novelty detection in text. In *Proceedings of International Conference on Machine Learning*
- Basu, S., Bilenko, M., and Mooney, R. J. (2004). A probabilistic framework for semi-supervised clustering. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM Press, New York, NY, USA, 59–68.
- Bay, S. D. and Schwabacher, M. (2003). Mining distance-based outliers in near linear time with randomization and a simple pruning rule. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM Press, 29–38.
- Bentley, J. L. (1975). Multidimensional binary search trees used for associative searching. *Communications of the ACM* 18, 9, 509–517
- Bishop, C. (1994). Novelty detection and neural network validation. In *Proceedings of IEEE Vision, Image and Signal Processing*. Vol. 141. 217–222.
- Bolton, R. and Hand, D. (1999). Unsupervised profiling methods for fraud detection. In *Credit Scoring and Credit Control VII*
- Box, G. E., Jenkins, G. M., Reinsel, G. C., and Ljung, G. M. (2015). *Time series analysis: forecasting and control*. John Wiley and Sons.
- Brause, R., Langsdorf, T., and Hepp, M. (1999). Neural data mining for credit card fraud detection. In *Proceedings of IEEE International Conference on Tools with Artificial Intelligence*. 103–106.
- Breunig, M. M., Kriegel, H.-P., Ng, R. T., and Sander, J. (2000). Lof: identifying density-based local outliers. In *Proceedings of 2000 ACM SIGMOD International Conference on Management of Data*. ACM Press, 93–104.
- Byunggho, H. and Sungzoon, C. (1999). Characteristics of autoassociative mlp as a novelty detector. In *Proceedings of IEEE International Joint Conference on Neural Networks*. Vol. 5. 3086–3091.
- Chandola, V., Banerjee, A., & Kumar, V. (2009). Anomaly detection. *ACM Computing Surveys*, 41(3), 1–58. <https://doi.org/10.1145/1541880.1541882>
- Chandola, V., Boriah, S., and Kumar, V. (2008). Understanding categorical similarity measures for outlier detection. Tech. Rep. 08-008, University of Minnesota. Mar.



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

- Chatfield, C., amp; Xing, H. (2019) The analysis of time series: an introduction with R. CRC Press. Chatfield, C., amp;
- Chen, D., Shao, X., Hu, B., and Su, Q. (2005). Simultaneous wavelength selection and outlier detection in multivariate regression of near-infrared spectra. *Analytical Sciences* 21, 2, 161–167.
- Chen, P., Niu, A., Liu, D., Jiang, W., & Ma, B. (2018). Time Series Forecasting of Temperatures Using Sarima: An example from Nanjing. *IOP Conference Series: Materials Science and Engineering*, 394, 052024.
<https://doi.org/10.1088/1757-899x/394/5/052024>
- Chow, C. and Yeung, D.-Y. (2002). Parzen-window network intrusion detectors. In *Proceedings of the 16th International Conference on Pattern Recognition*. Vol. 4. IEEE Computer Society, Washington, DC, USA, 40385.
- Crook, P. A., Marsland, S., Hayes, G., and Nehmzow, U. (2002). A tale of two filters - on-line novelty detection. In *Proceedings of International Conference on Robotics and Automation*. 3894–3899.
- Crook, P. and Hayes, G. (2001). A robot implementation of a biologically inspired method for novelty detection. In *Proceedings of Towards Intelligent Mobile Robots Conference*. Manchester, UK.
- Cun, Y. L., Boser, B., Denker, J. S., Howard, R. E., Hubbard, W., Jackel, L. D., and Henderson, D. (1990). Handwritten digit recognition with a back-propagation network. *Advances in neural information processing systems*, 396–404.
- Das, K. and Schneider, J. (2007). Detecting anomalous records in categorical datasets. In *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM Press
- Davy, M. and Godsill, S. (2002). Detection of abrupt spectral changes using support vector machines. an application to audio signal segmentation. In *Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing*. Orlando, USA
- Desforges, M., Jacob, P., and Cooper, J. (1998). Applications of probability density estimation to the detection of abnormal conditions in engineering. In *Proceedings of Institute of Mechanical Engineers*. Vol. 212. 687–703
- Diehl, C. and Hampshire, J. (2002). Real-time object classification and novelty detection for collaborative video surveillance. In *Proceedings of IEEE International Joint Conference on Neural Networks*. IEEE, Honolulu, HI



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

- Ertoz, L., Eilertson, E., Lazarevic, A., Tan, P.-N., Kumar, V., Srivastava, J., and Dokas, P. (2004). MINDS - Minnesota Intrusion Detection System. In *Data Mining - Next Generation Challenges and Future Directions*. MIT Press.
- Eskin, E. (2000). Anomaly detection over noisy data using learned probability distributions. In *Proceedings of the Seventeenth International Conference on Machine Learning*. Morgan Kaufmann Publishers Inc., 255–262.
- Eskin, E., Arnold, A., Prerau, M., Portnoy, L., and Stolfo, S. (2002). A geometric framework for unsupervised anomaly detection. In *Proceedings of Applications of Data Mining in Computer Security*. Kluwer Academics, 78–100.
- Fan, W., Miller, M., Stolfo, S. J., Lee, W., and Chan, P. K. (2001). Using artificial anomalies to detect unknown and known network intrusions. In *Proceedings of the 2001 IEEE International Conference on Data Mining*. IEEE Computer Society, 123–130.
- Galeano, P., Pea, D., and Tsay, R. S. (2004). Outlier detection in multivariate time series via projection pursuit. *Statistics and Econometrics Working Papers ws044211*, Universidad Carlos III, Departamento de Estadística y Econometría. Sep.
- Goldberger, A. L., Amaral, L. A. N., Glass, L., Hausdorff, J. M., Ivanov, P. C., Mark, R. G., Mietus, J. E., Moody, G. B., Peng, C.-K., and Stanley, H. E. (2000). PhysioBank, PhysioToolkit, and PhysioNet: Components of a new research resource for complex physiologic signals. *Circulation* 101, 23, e215–e220. *Circulation Electronic Pages*:
<http://circ.ahajournals.org/cgi/content/full/101/23/e215>.
- Goldberger, A. L., Amaral, L. A., Glass, L., Hausdorff, J. M., Ivanov, P. C., Mark, R. G., Mietus, J. E., Moody, G. B., Peng, C.-K., & Stanley, H. E. (2000). PhysioBank, PhysioToolkit, and PhysioNet. *Circulation*, 101(23).
<https://doi.org/10.1161/01.cir.101.23.e215>
- Grzegorzek, J. (2023). Time-series analysis in sales forecasting — Super Business Manager. <https://www.superbusinessmanager.com/time-series-analysis-in-sales-forecasting/>
- Guttormsson, S., II, R. M., and El-Sharkawi, M. (1999). Elliptical novelty grouping for on-line short-turn detection of excited running rotors. *IEEE Transactions on Energy Conversion* 14, 1 (March).
- Gwadera, R., Atallah, M. J., and Szpankowski, W. (2005). Reliable detection of episodes in event sequences. *Knowledge and Information Systems* 7, 4, 415–437.



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

- Hawkins, S., He, H., Williams, G. J., and Baxter, R. A. (2002). Outlier detection using replicator neural networks. In Proceedings of the 4th International Conference on Data Warehousing and Knowledge Discovery. Springer-Verlag, 170–180
- He, Z., Deng, S., and Xu, X. (2002). Outlier detection integrating semantic knowledge. In Proceedings of the Third International Conference on Advances in Web-Age Information Management. Springer-Verlag, London, UK, 126–131.
- He, Z., Xu, X., and Deng, S. (2003). Discovering cluster-based local outliers. Pattern Recognition Letters 24, 9-10, 1641–1650.
- He, Z., Xu, X., Huang, J. Z., and Deng, S. (2004). A frequent pattern discovery method for outlier detection. 726–732.
- He, Z., Xu, X., Huang, J. Z., and Deng, S. (2004). Mining class outliers: Concepts, algorithms and applications. 588–589.
- Heller, K. A., Svore, K. M., Keromytis, A. D., and Stolfo, S. J. (2003). One class support vector machines for detecting anomalous windows registry accesses. In Proceedings of the Workshop on Data Mining for Computer Security
- Ho, T. V. and Rouat, J. (1998). Novelty detection based on relaxation time of a network of integrate-and-fire neurons. In Proceedings of Second IEEE World Congress on Computational Intelligence. Anchorage, AK, 1524–1529.
- Ihler, A., Hutchins, J., and Smyth, P. (2006). Adaptive event detection with time-varying poisson processes. In Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining. ACM Press, New York, NY, USA, 207–216.
- Jose, J. (2022). *Introduction to Time Series and Its Applications*.
https://www.researchgate.net/publication/362389180_INTRODUCTION_TO_TIME_SERIES_ANALYSIS_AND_ITS_APPLICATIONS
- Kearns, M. J. (1990). Computational Complexity of Machine Learning. MIT Press, Cambridge, MA, USA.
- King, S., King, D., P. Anuzis, K. A., Tarassenko, L., Hayton, P., and Utete, S. (2002). The use of novelty detection techniques for monitoring high-integrity plant. In Proceedings of the 2002 International Conference on Control Applications. Vol. 1. Cancun, Mexico, 221–226.
- Knorr, E. M., Ng, R. T., and Tucakov, V. (2000). Distance-based outliers: algorithms and applications. The VLDB Journal 8, 3-4, 237–253.



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

- Kojima, K. and Ito, K. (1999). Autonomous learning of novel patterns by utilizing chaotic dynamics. In IEEE International Conference on Systems, Man, and Cybernetics. Vol. 1. IEEE, Tokyo, Japan, 284–289.
- Kou, Y., Lu, C.-T., & Chen, D. (2006). Spatial weighted outlier detection. Proceedings of the 2006 SIAM International Conference on Data Mining. <https://doi.org/10.1137/1.9781611972764.71>
- Kou, Y., Lu, C.-T., and Chen, D. (2006). Spatial weighted outlier detection. In Proceedings of SIAM Conference on Data Mining.
- Kumar, V. (2005). Parallel and distributed computing for Cybersecurity. IEEE Distributed Systems Online, 6(10). <https://doi.org/10.1109/mdso.2005.53>
- Kumar, V. (2005). Parallel and distributed computing for cybersecurity. Distributed Systems Online, IEEE 6, 10.
- LAYEK, G. C. (2016). *Introduction to dynamical systems and Chaos*. SPRINGER.
- Lee, W., Stolfo, S. J., and Mok, K. W. (2000). Adaptive intrusion detection: A data mining approach. Artificial Intelligence Review 14, 6, 533–567.
- Lu, C.-T., Chen, D., and Kou, Y. (2003). Algorithms for spatial outlier detection. In Proceedings of 3rd International Conference on Data Mining. 597–600
- Ma, J. and Perkins, S. (2003). Online novelty detection on temporal sequences. In Proceedings of the 9th ACM SIGKDD international conference on Knowledge discovery and data mining. ACM Press, New York, NY, USA, 613–618.
- Mahoney, M. V., Chan, P. K., and Arshad, M. H. (2003). A machine learning approach to anomaly detection. Tech. Rep. CS–2003–06, Department of Computer Science, Florida Institute of Technology Melbourne FL 32901. march.
- Markou, M. and Singh, S. (2003). Novelty detection: a review-part 2: neural network based approaches. Signal Processing 83, 12, 2499–2521.
- Markou, M., & Singh, S. (2003). Novelty detection: A review—part 1: Statistical approaches. Signal Processing, 83(12), 2481–2497. <https://doi.org/10.1016/j.sigpro.2003.07.018>
- Martinez, D. (1998). Neural tree density estimation for novelty detection. IEEE Transactions on Neural Networks 9, 2, 330–338.
- Méndez, M., Ibbias, A., & Núñez, M. (2022). Using deep learning to detect anomalies in traffic flow. *Intelligent Information and Database Systems*, 299–312. https://doi.org/10.1007/978-3-031-21743-2_24



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

- Moya, M., Koch, M., and Hostetler, L. (1993). One-class classifier networks for target recognition applications. In Proceedings on World Congress on Neural Networks, International Neural Network Society. Portland, OR, 797–801.
- Parzen, E. (1962). On the estimation of a probability density function and mode. *Annals of Mathematical Statistics* 33, 1065–1076.
- Platt, J. (2000). Probabilistic outputs for support vector machines and comparison to regularized likelihood methods. A. Smola, P. Bartlett, B. Schoelkopf, and D. Schuurmans, Eds. 61–74
- Rajeswari, A. M., Yalini, S. K., Janani, R., Rajeswari, N., & Deisy, C. (2018). A comparative evaluation of supervised and unsupervised methods for detecting outliers. *2018 Second International Conference on Inventive Communication and Computational Technologies (ICICCT)*. <https://doi.org/10.1109/icicct.2018.8473123>
- Ramaswamy, S., Rastogi, R., and Shim, K. (2000). Efficient algorithms for mining outliers from large data sets. In Proceedings of the 2000 ACM SIGMOD international conference on Management of data. ACM Press, 427–438.
- Rousseeuw, P. J. and Leroy, A. M. (1987). Robust regression and outlier detection. John Wiley & Sons, Inc., New York, NY, USA.
- Rousseeuw, P., & Leroy, A. M. (2003). Robust regression and outlier detection. Wiley.
- Roussopoulos, N., Kelley, S., and Vincent, F. (1995). Nearest neighbor queries. In Proceedings of ACM-SIGMOD International Conference on Management of Data.
- Salvador, S. and Chan, P. (2003). Learning states and rules for time-series anomaly detection. Tech. Rep. CS–2003–05, Department of Computer Science, Florida Institute of Technology Melbourne FL 32901. March
- Scholkopf, B., Platt, J. C., Shawe-Taylor, J. C., Smola, A. J., and Williamson, R. C. (2001). Estimating the support of a high-dimensional distribution. *Neural Comput.* 13, 7, 1443–1471.
- Scott, S. L. 2001. Detecting network intrusion using a markov modulated nonhomogeneous poisson process. Submitted to the Journal of the American Statistical Association.
- Sebyala, A. A., Olukemi, T., and Sacks, L. (2002). Active platform security through intrusion detection using naive bayesian network for anomaly detection. In Proceedings of the 2002 London Communications Symposium.



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

- Sheikholeslami, G., Chatterjee, S., and Zhang, A. (1998). Wavecluster: A multi-resolution clustering approach for very large spatial databases. In Proceedings of the 24rd International Conference on Very Large Data Bases. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 428–439.
- Song, X., Wu, M., Jermaine, C., and Ranka, S. (2007). Conditional anomaly detection. *IEEE Transactions on Knowledge and Data Engineering* 19, 5, 631–645.
- Spence, C., Parra, L., & Sajda, P. (n.d.). Detection, synthesis and compression in mammographic image analysis with a hierarchical image probability model. *Proceedings IEEE Workshop on Mathematical Methods in Biomedical Image Analysis (MMBIA 2001)*. <https://doi.org/10.1109/mmbia.2001.991693>
- Spence, C., Parra, L., and Sajda, P. (2001). Detection, synthesis and compression in mammographic image analysis with a hierarchical image probability model. In *Proceedings of the IEEE Workshop on Mathematical Methods in Biomedical Image Analysis*. IEEE Computer Society, Washington, DC, USA, 3.
- Stefano, C., Sansone, C., and Vento, M. (2000). To reject or not to reject: that is the question— an answer in case of neural classifiers. *IEEE Transactions on Systems, Management and Cybernetics* 30, 1, 84–94.
- Sun, J., Qu, H., Chakrabarti, D., and Faloutsos, C. (2005). Neighborhood formation and anomaly detection in bipartite graphs. In *Proceedings of the 5th IEEE International Conference on Data Mining*. IEEE Computer Society, Washington, DC, USA, 418–425.
- Tan, P.-N., Steinbach, M., and Kumar, V. (2005). *Introduction to Data Mining*. Addison-Wesley.
- Tan, P.-N., Steinbach, M., Karpatne, A., & Kumar, V. (2020). *Introduction to data mining*. Pearson Education Limited.
- Tarassenko, L. (1995). Novelty detection for the identification of masses in mammograms. In *Proceedings of the 4th IEEE International Conference on Artificial Neural Networks*. Vol. 4. Cambridge, UK, 442–447
- Teng, H. S., Chen, K., & Lu, S. C. (1990). Adaptive real-time anomaly detection using inductively generated sequential patterns. *Proceedings. 1990 IEEE Computer Society Symposium on Research in Security and Privacy*. <https://doi.org/10.1109/risp.1990.63857>
- Teng, H., Chen, K., and Lu, S. (1990). Adaptive real-time anomaly detection using inductively generated sequential patterns. In *Proceedings of IEEE Computer Society Symposium on Research in Security and Privacy*. IEEE Computer Society Press, 27



ΠΑΝΕΠΙΣΤΗΜΙΟ ΘΕΣΣΑΛΙΑΣ
ΠΟΛΥΤΕΧΝΙΚΗ ΣΧΟΛΗ
ΤΜΗΜΑ ΠΟΛΙΤΙΚΩΝ ΜΗΧΑΝΙΚΩΝ

- Torr, P. and Murray, D. (1993). Outlier detection and motion segmentation. In Proceedings of SPIE, Sensor Fusion VI, Paul S. Schenker; Ed. Vol. 2059. 432–443.
- Tsay, R. S., Pea, D., and Pankratz, A. E. 2000. Outliers in multivariate time series. *Biometrika* 87, 4, 789–804
- Vinueza, A. and Grudic, G. (2004). Unsupervised outlier detection and semi-supervised learning. Tech. Rep. CU-CS-976-04, Univ. of Colorado at Boulder. May
- Xiao, D., & Su, J. (2022). Research on stock price time series prediction based on Deep Learning and autoregressive integrated moving average. *Scientific Programming*, 2022, 1–12. <https://doi.org/10.1155/2022/4758698>
- Xing, H.(2019). The analysis of time series: an introduction with R.
- Yairi, T., Kato, Y., and Hori, K. (2001). Fault detection by mining association rules from housekeeping data. In In Proceedings of International Symposium on Artificial Intelligence, Robotics and Automation in Space.
- Zhang, Y., & Lei, Y. (2021). Data anomaly detection of bridge structures using convolutional neural network based on structural vibration signals. *Symmetry*, 13(7), 1186. <https://doi.org/10.3390/sym13071186>