

Ανοικτό Πανεπιστήμιο Κύπρου

Σχολή Θετικών και Εφαρμοσμένων Επιστημών

**Μεταπτυχιακό Πρόγραμμα Σπουδών *Ασφάλεια
Υπολογιστών και Δικτύων***

Μεταπτυχιακή Διατριβή



**Συλλογή και Προεπεξεργασία Δεδομένων από Ενεργητικές
Επιθέσεις Αναγνώρισης (Gathering & preprocessing data
from active footprinting attacks)**

Βασίλειος Ίβρος

**Επιβλέπων Καθηγητής
Ιωάννης Μαυρίδης**

Μάιος 2017

Ανοικτό Πανεπιστήμιο Κύπρου

Σχολή Θετικών και Εφαρμοσμένων Επιστημών

Μεταπτυχιακό Πρόγραμμα Σπουδών: Ασφάλεια

Υπολογιστών και Δικτύων

Μεταπτυχιακή Διατριβή

**Συλλογή και Προεπεξεργασία Δεδομένων από Ενεργητικές
Επιθέσεις Αναγνώρισης (Gathering & preprocessing data
from active footprinting attacks)**

Βασίλειος Ίβρος

**Επιβλέπων Καθηγητής
Ιωάννης Μαυρίδης**

Η παρούσα μεταπτυχιακή διατριβή υποβλήθηκε προς μερική εκπλήρωση των απαιτήσεων για απόκτηση μεταπτυχιακού τίτλου σπουδών στο μεταπτυχιακό πρόγραμμα Ασφάλειας Υπολογιστών και Δικτύων από τη Σχολή Θετικών και Εφαρμοσμένων Επιστημών του Ανοικτού Πανεπιστημίου Κύπρου.

Δεκέμβριος 2016

Περίληψη

Οι περισσότερες μεθοδολογίες Ελέγχου Τρωτότητας (Penetration Testing) περιλαμβάνουν στο πρώτο στάδιό τους, τεχνικές Footprinting. Οι ενεργητικές τεχνικές Footprinting γίνονται αντιληπτές από τα συστήματα Ελέγχου Παρέισφρησης (Intrusion Detection System) ή άλλους, παρόμοιου σκοπού μηχανισμούς ασφάλειας. Ωστόσο, ο όγκος των δεδομένων που συλλέγονται και οι διαστάσεις του διανύσματος χαρακτηριστικών που λαμβάνουμε, καθιστούν δύσκολη την άμεση επεξεργασία τους. Στην τρέχουσα εργασία μελετώνται οι ενεργητικές τεχνικές Footprinting με σκοπό την κατηγοριοποίηση των δεδομένων, που ανιχνεύονται από τους μηχανισμούς ασφάλειας. Τα δεδομένα αυτά πρόκειται να αναλυθούν ως προς την πολυπλοκότητά τους για να μπορούν να εξαχθούν συμπεράσματα για κάθε δικτυακό κόμβο στη βάση σχετικών χαρακτηριστικών.

Μια ακόμη σημαντική παράμετρος που επηρεάζει την ροή της έρευνας, είναι η διαρκής καταγραφή στοιχείων σε αρχεία μητρώου, που πραγματοποιείται από τους σύγχρονους μηχανισμούς ασφάλειας υπολογιστών και δικτύων, που όμως φέρνει τον Αναλυτή Ασφάλειας αντιμέτωπο με το πρόβλημα διαχείρισης μεγάλου όγκου δεδομένων. Η επιλογή των ελάχιστων, αλλά σημαντικότερων χαρακτηριστικών, θα οδηγήσει στην μείωση των διαστάσεων του διανύσματος χαρακτηριστικών. Με τον τρόπο αυτό η επεξεργασία των δεδομένων θα είναι λιγότερο απαιτητική σε υπολογιστικούς πόρους, ενώ η εξαγωγή συμπερασμάτων θα μπορεί να πραγματοποιηθεί σε πραγματικό χρόνο, βοηθώντας προς την κατεύθυνση του έγκαιρου εντοπισμού ηλεκτρονικών επιθέσεων.

Λέξεις κλειδιά: Εξόρυξη Δεδομένων, Footprinting, Επιλογή Χαρακτηριστικών, Ανίχνευση Παρεισφρήσεων, Σύνολα Δεδομένων, kdd 1999, Rstudio.

Summary

Most Penetration Testing methodologies include Footprinting techniques at their first stage. Active Footprinting techniques can be detected by Intrusion Detection Systems (IDS) or other security systems of similar purpose. Nevertheless, the amount of data collected and the dimensions of the features vector received make their immediate process quite difficult. In this dissertation, a study of active footprinting techniques will be studied in order to categorize the data collected by security systems. These collected data are going to be analysed according to their complexity, so that conclusions can be excluded for each network node according to relevant features.

An also important factor effecting the flow of the research is the constant recording of data in log files, which is carried out by up to date computer and network security systems, arousing for the Security Analyst the problem of having to handle a great volume of data. The selection of, the least in amount but the most important in significance, specifics will lead in the features vector's dimensions' reduction. By this way the data process will be less demanding as it comes to computational resources and conclusions could be reached in real time contributing with the on time detection of electronic attacks.

Keywords: Data Mining, Footprinting, Feature Selection, intrusion Detection, Datasets, KDD 1999, Rstudio

Ευχαριστίες

Πριν ξεκινήσω την παρουσίαση των ευρημάτων και την ανάλυση τους, θα ήθελα αρχικά να ευχαριστήσω τον επιβλέποντα καθηγητή μου κ. Ιωάννη Μαυρίδη, που μου πρότεινε το συγκεκριμένο θέμα μεταπτυχιακής διατριβής και με εμπιστεύτηκε ώστε να το υλοποιήσω καθώς και για την καθοδήγηση του κατά την διάρκεια της υλοποίησης της μεταπτυχιακής διατριβής.

Κλείνοντας θα ήθελα να ευχαριστήσω τους γονείς μου Γιάννη και Βασιλική καθώς επίσης και τα αδέρφια μου Δημήτρη και Έλενα που βρίσκονται δίπλα μου και με στηρίζουν σε κάθε μου βήμα.

Περιεχόμενα

1	Εισαγωγή	1
1.1	Δομή της Μεταπτυχιακής Διατριβής.....	1
1.2	Μεθοδολογία της Μεταπτυχιακής Διατριβής.....	2
1.3	Συνεισφορά της Μεταπτυχιακής Διατριβής.....	3
2	Θεωρητική θεμελίωση	4
2.1	Κατηγορίες Footprinting.....	4
2.1.1	Τύποι Ενεργητικού Footprinting.....	6
2.1.2	Ποιος υλοποιεί τεχνικές Ενεργητικού Footprinting και γιατί;.....	7
2.2	Ανίχνευση και Αποτροπή Εισβολών.....	9
2.2.1	Συστήματα Αποτροπής Παρεισφρήσεων.....	11
2.2.2	Συστήματα Ανίχνευσης Παρεισφρήσεων.....	11
2.3	Μεθοδολογίες Ανίχνευσης Εισβολής	11
2.3.1	Ανίχνευση Διαταραχών.....	13
2.3.2	Ανίχνευση Κακής Συμπεριφοράς.....	14
2.3.3	Συστήματα Ανίχνευσης Παρεισφρήσεων βασισμένα στην ανάλυση υπογραφών.....	15
2.3.4	Συστήματα Ανίχνευσης Παρεισφρήσεων βασισμένα στην ανάλυση ανωμαλιών.....	15
2.3.5	Συστήματα ανίχνευσης Παρεισφρήσεων Δικτύου.....	16
2.3.6	Τοπικά Συστήματα Ανίχνευσης Παρεισφρήσεων.....	20
3	Εξόρυξη Δεδομένων	24
3.1	Μέθοδοι Εξόρυξης Δεδομένων.....	25
3.1.1	Παλινδρόμηση.....	25
3.1.2	Σχέση κανόνα μάθησης.....	26
3.1.3	Κατηγοριοποίηση.....	27
3.1.4	Ομαδοποίηση.....	30
3.1.4.1	Αλγόριθμοι μη Ιεραρχική Ομαδοποίησης.....	33
4	Ανάλυση Συνόλου Δεδομένων	36
4.1	Πρωτόκολλα Επικοινωνίας.....	37
4.2	Ανάλυση Χαρακτηριστικών.....	39
5	Υλοποίηση Εφαρμογής	44
5.1	Εισαγωγή στο Rstudio.....	44
5.2	Βήματα Ανάλυσης Συνόλου Δεδομένων.....	44
5.2.1	Φόρτωση και Προεπεξεργασία Δεδομένων.....	45
5.2.2	Μοντέλο Κατηγοριοποίησης Naïve Bayes.....	49
5.2.3	Μοντέλο Κατηγοριοποίησης Δέντρων Απόφασης.....	51

5.2.4	Ανάλυση Παραγόμενου Συνόλου Δεδομένων.....	52
6	Επίλογος.....	55
	Βιβλιογραφία.....	56

Κεφάλαιο 1

Εισαγωγή

Η ραγδαία αύξηση της χρήσης των πληροφοριακών συστημάτων και των δικτύων σε όλο το εύρος της καθημερινότητας των χρηστών, των οργανισμών και των εταιριών έχει δημιουργήσει και την ανάλογη αύξηση των υποψήφιων εισβολέων. Οι τελευταίοι στοχεύουν να παρεισφρήσουν στο πληροφοριακό σύστημα στο οποίο περιλαμβάνονται και διακινούνται από χαμηλής σημαντικότητας έως άκρως εμπιστευτικά και απόρρητα δεδομένα με διάφορες μεθόδους που θα αναλυθούν στην παρούσα μεταπτυχιακή διατριβή. Στόχοι των εισβολέων είναι: η απόκτηση των δεδομένων, η τροποποίηση ή διαγραφή των δεδομένων και η άρνηση υπηρεσιών(denial of service) του πληροφοριακού συστήματος ή του δικτύου. Συνεπώς ήταν μονόδρομος για την προστασία των πληροφοριακών συστημάτων και δικτύων να αναπτυχθούν μέθοδοι και συστήματα ανίχνευσης παρεισφρήσεων, τα οποία έχουν ως στόχο τον εντοπισμό των εισβολών και μπορούν ακόμη και να εκτελέσουν σενάρια αναχαίτισης των εισβολών που εντοπίζονται. Οι μέθοδοι αυτές που αναπτύχθηκαν καθώς και τα συστήματα ανίχνευσης παρεισφρήσεων βρίσκονται σε μια συνεχή έρευνα του τρόπου λειτουργίας τους έτσι ώστε να βελτιώνονται και να προσαρμόζονται στα καινούργια δεδομένα.

1.1 Δομή της Μεταπτυχιακής Διατριβής

Η παρούσα μεταπτυχιακή διατριβή αποτελείται από έξι κεφάλαια. Επιγραμματικά το κάθε κεφάλαιο περιέχει τα εξής θέματα:

Στο κεφάλαιο 1 παρουσιάζεται το πρόβλημα που πραγματεύεται η μεταπτυχιακή διατριβή. Επίσης η μεθοδολογία που θα ακολουθηθεί στην παρούσα μεταπτυχιακής διατριβής καθώς και η συνεισφορά της.

Στο κεφάλαιο 2 παρουσιάζονται θεωρητικές έννοιες που σχετίζονται με την μεταπτυχιακή διατριβή όπως: Το Footprinting, τα Συστήματα Ανίχνευσης Παρεισφρήσεων και οι μέθοδοι Κατηγοριοποίησης Συλλογής Δεδομένων.

Στο κεφάλαιο 3 παρουσιάζονται θεωρητικές έννοιες που σχετίζονται άμεσα με την μεθοδολογία που ακολουθήθηκε στην μεταπτυχιακή διατριβή όπως οι μεθοδολογίες Εξόρυξης Δεδομένων(Data Mining).

Στο κεφάλαιο 4 παρουσιάζετε το σύνολο δεδομένων το οποίο θα αναλυθεί.

Στο κεφάλαιο 5 παρουσιάζετε το λογισμικό ανάλυσης που θα χρησιμοποιηθεί για την ανάλυση του συνόλου δεδομένων, καθώς επίσης και η μεθοδολογία ανάλυσης.

Στο κεφάλαιο 6 παρουσιάζονται τα συμπεράσματα που εξάγονται με την ολοκλήρωση της μεταπτυχιακής διατριβής.

1.2 Μεθοδολογία της Μεταπτυχιακής Διατριβής

Μια σημαντική παράμετρος που επηρεάζει την ροή της έρευνας, είναι η διαρκής καταγραφή στοιχείων σε αρχεία μητρώου, που πραγματοποιείται από τους σύγχρονους μηχανισμούς ασφάλειας Η/Υ και δικτύων, που όμως φέρνει τον Αναλυτή Ασφάλειας αντιμέτωπο με το πρόβλημα διαχείρισης μεγάλου όγκου δεδομένων. Η επιλογή των ελάχιστων, αλλά σημαντικότερων χαρακτηριστικών, θα οδηγήσει στην μείωση των διαστάσεων του διανύσματος χαρακτηριστικών. Με τον τρόπο αυτό η επεξεργασία των δεδομένων θα είναι λιγότερο απαιτητική σε υπολογιστικούς πόρους, και η εξαγωγή συμπερασμάτων θα μπορεί να πραγματοποιηθεί σε πραγματικό χρόνο, βοηθώντας προς την κατεύθυνση του έγκαιρου εντοπισμού ηλεκτρονικών επιθέσεων. Για το σκοπό αυτό, θα αξιοποιηθεί ένα έτοιμο ελεύθερο διαθέσιμο σύνολο δεδομένων(kdd 1999). Το σύνολο δεδομένων(kdd 1999) θα εισάγει τα δεδομένα τα οποία όταν θα τα επεξεργαστούμε με μία τεχνική επιλογής χαρακτηριστικών, χρησιμοποιώντας το ολοκληρωμένο περιβάλλον ανάπτυξης του λογισμικού ανάλυσης δεδομένων R, θα παραχθεί το τελικό παραγόμενο διάνυσμα χαρακτηριστικών, ώστε να γίνει η τελική παρουσίαση των αποτελεσμάτων.

1.3 Συνεισφορά της Μεταπτυχιακής Διατριβής

Η συγκεκριμένη έρευνα θα συνεισφέρει στη βελτίωση της λειτουργίας των συστημάτων ανίχνευσης παρεισφρήσεων, αναλύοντας τα διάφορα datasets, αξιολογώντας τις πτυχές τους και συνδυάζοντας μέσω της επιλογής, χαρακτηριστικών και των τεχνικών κατηγοριοποίησης τα επιστημονικά πεδία της Ασφάλειας Η/Υ και της Μηχανικής Μάθησης

Κεφάλαιο 2

Θεωρητική θεμελίωση

2.1 Κατηγορίες Footprinting

Ως Footprinting ονομάζεται η διαδικασία συλλογής πληροφοριών για το περιβάλλον ενός στόχου και το δίκτυο του. Συγκεκριμένα χωρίζετε στις παρακάτω κατηγορίες:

1. Το Ανώνυμο Footprinting (Anonymous Footprinting) που ονομάζεται η διαδικασία συλλογής πληροφοριών από πηγές στις οποίες δεν μπορεί να προσδιοριστεί ή να εντοπιστεί ο συντάκτης των πληροφοριών.
2. Το Ψευδώνυμο Footprinting(Pseudonymous Footprinting) που ονομάζεται η διαδικασία συλλογής πληροφοριών που δεν δημοσιεύονται με το κανονικό όνομα του συντάκτη άλλα με την χρήση ενός διαφορετικού ονόματος(ψευδώνυμο). Σε μια προσπάθεια του συντάκτη να διατηρήσει την ανωνυμία του.
3. Το Οργανωτικό ή Ιδιωτικό Footprinting(Organizational or private Footprinting) που ονομάζετε η διαδικασία συλλογής πληροφοριών από το διαδικτυακά βασισμένο(web-based) ημερολόγιο του οργανισμού και τις υπηρεσίες ηλεκτρονικού ταχυδρομείου(email).
4. Το Διαδικτυακό Footprinting (Internet Footprinting) που ονομάζετε η διαδικασία συλλογής πληροφοριών σχετικά με τον στόχο από το διαδίκτυο(Internet).
5. Το Παθητικό Footprinting (Passive ή Open Source Footprinting) που είναι η διαδικασία με την οποία ο κακόβουλος χρήστης χρησιμοποιεί διάφορα εργαλεία και τεχνικές, στην προσπάθεια του να συλλέξει όσο το δυνατόν περισσότερες πληροφορίες μπορεί για τον στόχο, ώστε να επιλέξει τον

κατάλληλο τρόπο για να του επιτεθεί. Βασική προϋπόθεση στην συλλογή των πληροφοριών αυτών είναι η άγνοια του θύματος. Αφού οι πληροφορίες που συλλέγονται από τον κακόβουλο χρήστη είναι δημόσιες ή ορατές από απλούς χρήστες (McGreevy, 2002). Ο κακόβουλος χρήστης μπορεί να κάνει μια αναζήτηση δεδομένων από ιστοσελίδες όπως ο www.whois.com (Sanghvi, & Dahiya, 2013), οι οποίες παρέχουν ελεύθερα πληροφορίες σχετικά με το όνομα τομέα του θύματος, τον διακομιστή του και το δίκτυο του επίσης χρησιμοποιώντας προγράμματα όπως το `dir`, το `nslookup` κ.α. ο κακόβουλος χρήστης μπορεί να αποκτήσει περισσότερες πληροφορίες που σχετίζονται με τους πίνακες ονομάτων τομέα (DNS tables). Επιπλέον ο κακόβουλος χρήστης πλοηγείτε και σε άλλες δημόσιες πληροφορίες είτε από το δημόσιο ιστότοπο είτε από ανώνυμους ιστότοπους. Συλλέγοντας περισσότερες πληροφορίες για το δίκτυο, τις διευθύνσεις IP, τα πρωτόκολλα δικτύου που χρησιμοποιούνται, τα εσωτερικά ονόματα τομέα, τα συστήματα ανίχνευσης παρεισφρήσεων που χρησιμοποιεί το θύμα, τους τηλεφωνικούς αριθμούς, πληροφορίες σχετικά με το μοντέλο ελέγχου πρόσβασης κ.α. Οι παραπάνω πληροφορίες μπορούν να χρησιμοποιηθούν και από τους κακόβουλους χρήστες για να επιτεθούν σε ένα στόχο και από το προσωπικό ασφαλείας για να ενισχύσει την ασφάλεια του (McGreevy, 2002).

6. Ενεργητικό Footprinting (Active Footprinting) ονομάζετε η προσπάθεια συλλογής πληροφοριών που αφορούν τον στόχο(θύμα), τις οποίες για να τις αποκτήσει ο κακόβουλος χρήστης χρειάζεται να ενεργήσει ο ίδιος για την συλλογή τους με κίνδυνο να εκτεθεί και να αποκαλυφθεί η ταυτότητα του, αφού δεν είναι δημόσιες ή προσβάσιμες από απλούς χρήστες. Αυτό επιτυγχάνετε με διάφορες μεθόδους και εργαλεία όπως: ψάχνοντας για χρήσιμες πληροφορίες στα σκουπίδια της εταιρίας, χρησιμοποιώντας συσκευές keyloggers για να υποκλέψει δεδομένα, κάνοντας χρήση εργαλείων που υλοποιούν τεχνικές packet sniffing, port scanning κ.α.

2.1.1 Τύποι Ενεργητικού Footprinting

Ορισμένοι τύποι Ενεργητικού Footprinting είναι οι εξής:

- Η Κοινωνική μηχανική (social engineering) που είναι η τεχνική η οποία χρησιμοποιεί την επιρροή και την πειθώ για να εξαπατάει ανθρώπους. Αυτό επιτυγχάνετε με δύο τρόπους, είτε προσποιούμενος αυτός που ασκεί κοινωνική μηχανική κάποιον άλλο με περισσότερα δικαιώματα είτε χειραγωγώντας άλλους ανθρώπους π.χ. υπαλλήλους για να του προσκομίσουν πληροφορίες. Αυτό έχει ως αποτέλεσμα αυτός που ασκεί κοινωνική μηχανική να έχει ένα πλεονέκτημα απέναντι σε ανθρώπους που κατέχουν πληροφορίες για να τους τις αποσπάσει χωρίς να κάνει απαραίτητα χρήση τεχνολογικών μέσων(Mitnick, Simon, & Wozniak, 2002).
- Η Υποκλοπή (eavesdropping) που είναι η μη εξουσιοδοτημένη παρακολούθηση μεθόδου επικοινωνίας άλλων ανθρώπων π.χ. τηλεφωνικών κλήσεων και ηλεκτρονικών μηνυμάτων. Σε περίπτωση όπου η υποκλοπή δεν επηρεάζει την ποιότητα επικοινωνίας μεταξύ αποστολέα και παραλήπτη είναι πολύ δύσκολο να γίνει αντιληπτή. Ο κακόβουλος χρήστης εκμεταλλεύεται την απαρχαιωμένη ή ανύπαρκτη κρυπτογράφηση των δεδομένων που μεταφέρονται, ή οποιοδήποτε κενό ασφαλείας σε ένα δίκτυο π.χ. μη ενημερωμένο πρόγραμμα αποτροπής ιών, ώστε να ξεκινήσει μια επίθεση με την οποία θα αποσπάσει χρήσιμες πληροφορίες π.χ. man in the middle.
- Το Shoulder surfing που αναφέρετε στην χρήση άμεσων τεχνικών παρατήρησης από τον κακόβουλο χρήστη. Όπως να παρατηρεί πάνω από τους ώμους κάποιου για να συλλέξει πληροφορίες π.χ κωδικούς πρόσβασης.
- Η Έρευνα στα Σκουπίδια(Dumpster Diving) που ερευνά τα σκουπίδια ενός στόχου ώστε να εντοπιστούν πολύτιμες πληροφορίες. Βασική πτυχή της ύπαρξης αυτής της τεχνικής είναι η απερισκεψία των ανθρώπων για τα απορρίμματα τους. Αφού αποθέτουν στα σκουπίδια έγγραφα

οικονομικών συναλλαγών, ιατρικά σκευάσματα, κωδικούς, λίστες μισθοδοσίας κ.α. Με τον καιρό ο «δύτης σκουπιδιών»(κακόβουλος χρήστης) θα αποκτήσει την ανάλογη εμπειρία για να ξεχωρίζει χρήσιμες σακούλες καθώς και να αποφεύγει σακούλες από τουαλέτες κλπ(Mitnick, Simon, & Wozniak, 2002).

- Το Piggybacking που είναι η τεχνική για την απόκτηση πρόσβασης σε μια υπηρεσία στην οποία δεν έχει πρόσβαση ο κακόβουλος χρήστης και το πετυχαίνει ακλουθώντας κάποιον υπάλληλο/υπεύθυνο που έχει πρόσβαση στην υπηρεσία όπως το άνοιγμα μιας πόρτας με κωδικό.
- Το Ψάρεμα(Phishing) που είναι μια τεχνική επίθεσης κατά την οποία ο κακόβουλος χρήστης προσπαθεί να αποσπάσει πολύτιμες πληροφορίες παραπλανώντας το θύμα, προσποιούμενος μια αξιόπιστη οντότητα και ποντάροντας στην έλλειψη παρατηρητικότητας του θύματος. Η τεχνική αυτή εμφανίστηκε για πρώτη φορά το 1995, αποστέλλοντας παραπλανητικά μηνύματα ηλεκτρονικού ταχυδρομείου προς πολλούς χρήστες. Τα μηνύματα περιείχαν ένα σύνδεσμο που παρέπεμπε σε πλαστή ιστοσελίδα που ζητούσε από τους χρήστες να συνδεθούν καταχωρώντας τα στοιχεία τους με σκοπό να τους τα αποσπάσει(Shi, & Saleem, 2012).

2.1.2 Ποιος υλοποιεί τεχνικές Ενεργητικού Footprinting και γιατί;

Οι τεχνικές Ενεργητικού Footprinting συνήθως υλοποιούνται από δύο μεγάλες κατηγορίες χρηστών:

1. Στην πρώτη κατηγορία κατατάσσονται οι χρήστες που έχουν την πρόθεση να επιτεθούν στα δεδομένα ενός στόχου όπως οι κακόβουλοι χρήστες, κυβερνοεγκληματίες και οι ανταγωνιστές οι οποίοι έχουν ως σκοπό τους να αποκτήσουν μια πιο βαθιά γνώση σχετικά με τον στόχο/θύμα, την οποία σπάνια την παρέχουν οι υπόλοιπες τεχνικές Footprinting ανάλογα βέβαια και το επίπεδο ασφαλείας του στόχου. Οι πληροφορίες αυτές σχετίζονται με τα συστήματα ασφαλείας(π.χ. IDS, IPS,

Antivirus, firewall), τις συσκευές και το δίκτυο του στόχου. Οι κακόβουλοι χρήστες χρησιμοποιούν το active footprinting για:

- Να βρουν τον πιο εύκολο τρόπο ώστε να διεισδύσουν σε ένα πληροφοριακό σύστημα.
- Όστε να δημιουργήσουν ένα πιο οργανωμένο σχέδιο επίθεσης.
- Να περιορίσουν την περιοχή επίθεσης.

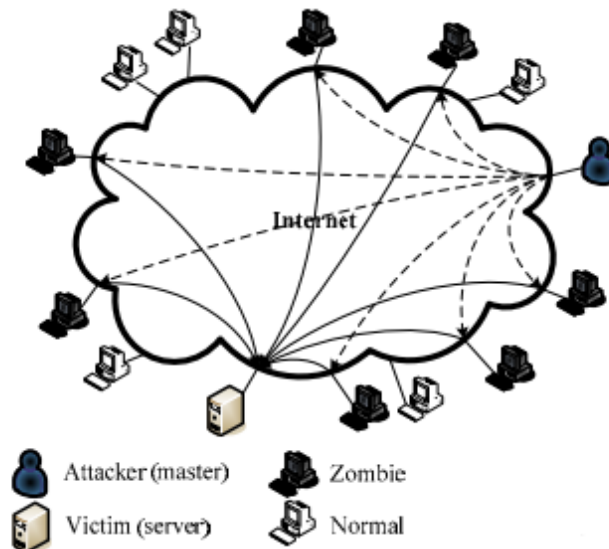
2. Στην δεύτερη κατηγορία κατατάσσεται το προσωπικό ασφαλείας μιας εταιρίας, ενός οργανισμού και γενικά ενός υποψήφιου στόχου στοχεύοντας στην ενίσχυση της ασφάλειας του στόχου(McGreeny, 2002). Το Footprinting έχει γίνει ένα αναγκαίο κακό για το προσωπικό ασφαλείας, αφού χρησιμοποιώντας το συλλέγει αρκετά δεδομένα που μπορεί να συλλέξει και ένας κακόβουλος χρήστης για τον στόχο, δίνοντας όμως την δυνατότητα στο προσωπικό ασφαλείας να δράση προληπτικά πριν την οποιαδήποτε επίθεση (Rechtin, & Maier, 2010).

Αν και σε αντίθεση με το παθητικό(Passive ή Open Source) Footprinting που έχει το πλεονέκτημα ότι δεν ενεργοποιεί κανέναν συναγερμό του στόχου, καθώς επίσης δεν θεωρείτε και παραβίαση ασφαλείας. Το ενεργητικό footprinting(active footprinting) όμως είναι παραβίαση ασφαλείας και σε περίπτωση που εντοπιστεί ενεργοποιεί συναγερμούς του στόχου. Αυτό άλλωστε είναι το βασικό του μειονέκτημα, αφού ο κακόβουλος χρήστης έρχεται σε επαφή με το δίκτυο του στόχου και με πληροφορίες τις οποίες δεν είναι εξουσιοδοτημένος να έχει πρόσβαση. Όμως το ενεργητικό footprinting έχει το πλεονέκτημα να μπορεί να αποκτήσει ο κακόβουλος χρήστης με την χρήση του πιο βαθιές πληροφορίες για το δίκτυο και το περιβάλλον του στόχου, όπως για παράδειγμα: ενεργά τερματικά(live hosts) με την χρήση σαρωτή δικτύου (network scanner) και ανοιχτές θύρες(open ports)με την χρήση σαρωτή θυρών(port scanner).

2.2 Ανίχνευση και Αποτροπή Εισβολών

Για να κατανοήσουμε τι ορίζουμε ως εισβολή σε ένα υπολογιστικό σύστημα θα αναλύσουμε τις τέσσερις κατηγορίες εισβολών:

1. Η άρνηση εκτέλεσης εφαρμογής(Denial of Service attack) είναι μια σημαντική επίθεση που προκαλεί πλήγμα στα δίκτυα υπολογιστών και τα υπολογιστικά συστήματα. Στοχεύοντας στην δημιουργία πλήγματος στην διαθεσιμότητα των πόρων όπως: στο εύρος ζώνης του καναλιού(link bandwidth), στην ενδιάμεση μνήμη(buffer) που αποθηκεύει προσωρινά δεδομένα συνδέσεων υπολογιστικών συσκευών καθώς επίσης και στην ενδιάμεση μνήμη εφαρμογών, στον επεξεργαστή κ.α. Έχοντας αποτελέσματα όπως το κατέβασμα μιας ιστοσελίδας, ή κατάρρευση καναλιού IRC(Internet Relay Chat) κ.α. Επειδή είναι δύσκολο για τον κακόβουλο χρήστη να υπερφορτώσει τους πόρους ενός πληροφοριακού συστήματος ή δικτύου κάνοντας χρήση ενός μόνο υπολογιστή. Για τον λόγο αυτό ο κακόβουλος χρήστης χρησιμοποιεί παγιδευμένες συσκευές(ονομάζονται zombie ή bots) στις οποίες έχει εγκαταστήσει κακόβουλο λογισμικό ώστε να καταφέρει να τις αναγκάσει να ξεκινήσουν παράλληλα επίθεση όλες μαζί στον στόχο, δημιουργώντας έναν τεράστιο όγκο δεδομένων τα οποία συναθροίζονται με δεδομένα που στέλνουν συσκευές χωρίς κακόβουλη πρόθεση στον στόχο κάνοντας έτσι πιο αποτελεσματική την επίθεση (Εικόνα 1). Η επίθεση άρνησης εκτέλεσης εφαρμογής με την χρήση και άλλων παγιδευμένων συσκευών από τον κακόβουλο χρήστη ονομάζεται κατανεμημένη άρνηση εκτέλεσης εφαρμογής(Distributed Denial of Service Attack) (Gu, & Liu, 2007).



Εικόνα 1. Επίθεση_κατανεμημένης άρνησης εκτέλεσης εφαρμογής(Distributed Denial of Service attack).

2. Η επίθεση απομακρυσμένου χρήστη(Remote to User attack, R2L) είναι μια επίθεση στην οποία ένας απομακρυσμένος χρήστης χρησιμοποιώντας μια απομακρυσμένη μηχανή π.χ. υπολογιστή στέλνει πακέτα στο υπολογιστικό σύστημα στόχο, προσπαθώντας να αποκτήσει μη εξουσιοδοτημένη πρόσβαση σε αυτό με διευρυμένα δικαιώματα χρήστη π.χ. admin ή root εκμεταλλευόμενος κάποια ευπάθεια του συστήματος (Chauhan, Mishra, & Kumar, 2011). Γνωστές επιθέσεις απομακρυσμένου χρήστη είναι οι: Spy, Phf, Multihop, Ftp_write, Imap, Warezmasters και Guess_passwd (Paliwal, & Gupta, 2012).

3. Η επίθεση χρήστη σε διαχειριστή(User to Root attack, U2R) είναι μια επίθεση όπου ένας τοπικός χρήστης με περιορισμένα δικαιώματα χρήστη προσπαθεί να εκμεταλλευτεί διάφορες ευπάθειες του υπολογιστικού συστήματος, ώστε να αποκτήσει μη εξουσιοδοτημένη πρόσβαση σε πληροφορίες που απαιτούν διευρυμένα δικαιώματα χρήστη π.χ. root ή admin(Paliwal, & Gupta, 2012). Ο ποιο διαδεδομένος τρόπος επίθεσης είναι η «υπερχείλιση μνήμης» (buffer overflow) που συνήθως προκαλούνται από προγραμματιστικά λάθη και υπάρχουν και άλλοι τύποι επιθέσεων όπως: το loadmodule, Perl, Rootkit, Httpunnel, Ps, Sqlattack και Xterm(Revathi, & Malathi, 2014).

4. Η επίθεση ανίχνευσης (Probing attack) είναι μια επίθεση κατά την διάρκεια της οποίας όπου ο κακόβουλος χρήστης σαρώνει το δίκτυο ή ένα υπολογιστικό σύστημα όπως, μια μηχανή ή μια συσκευή δικτύωσης, προκειμένου να συλλέξει πληροφορίες για ευπάθειες ή αδυναμίες τις οποίες θα εκμεταλλευτεί αργότερα να εισβάλει στον στόχο, ως τεχνική επίθεσης ανίχνευσης θεωρείτε και η κοινωνική μηχανική (social engineering) (Chauhan, Mishra, & Kumar, 2011). Επίσης τύποι επιθέσεων ανίχνευσης είναι οι: Nmap, Portswear, Ipsweep και Satan (Revathi, & Malathi, 2014).

2.2.1 Συστήματα Αποτροπής Παρεισφρήσεων

Τα Συστήματα Αποτροπής Παρεισφρήσεων (Intrusion Detection Systems, IPS) είναι μια επέκταση των Συστημάτων Ανίχνευσης Παρεισφρήσεων, επιπλέον όμως έχουν και ενεργητικό ρόλο αφού δεν αρκούνται μόνο στον εντοπισμό κακόβουλων ενεργειών, παραβιάσεων πολιτικών ασφαλείας και την δημιουργία αναφορών. Συγκεκριμένα όταν εντοπίσουν μια κακόβουλη ενέργεια πέραν της αναφοράς του γεγονότος που είναι η παθητική ενέργεια, μπορούν επίσης αυτόματα να απορρίψουν τα πακέτα που παρεκκλίνουν από τα καθορισμένους κανόνες και θεωρούνται ύποπτα και να διακόψουν τις συνδέσεις με τους απομακρυσμένους χρήστες που συμμετέχουν. Αξίζει να σημειωθεί ότι ορισμένα συστήματα αποτροπής Παρεισφρήσεων έχουν την ικανότητα να εφαρμόσουν διαφορετικές επιλογές προστασίας για διαφορετικά τμήματα του δικτύου. Χωρίζοντας με αυτόν τον τρόπο το δίκτυο σε τμήματα με διαφορετικούς κανόνες ασφαλείας με βάση τις ιδιαιτερότητες και τις ανάγκες του, αυξάνετε η ευελιξία του συστήματος αποτροπής Παρεισφρήσεων αφού γίνεται χρήσιμο για μεγάλα δίκτυα.

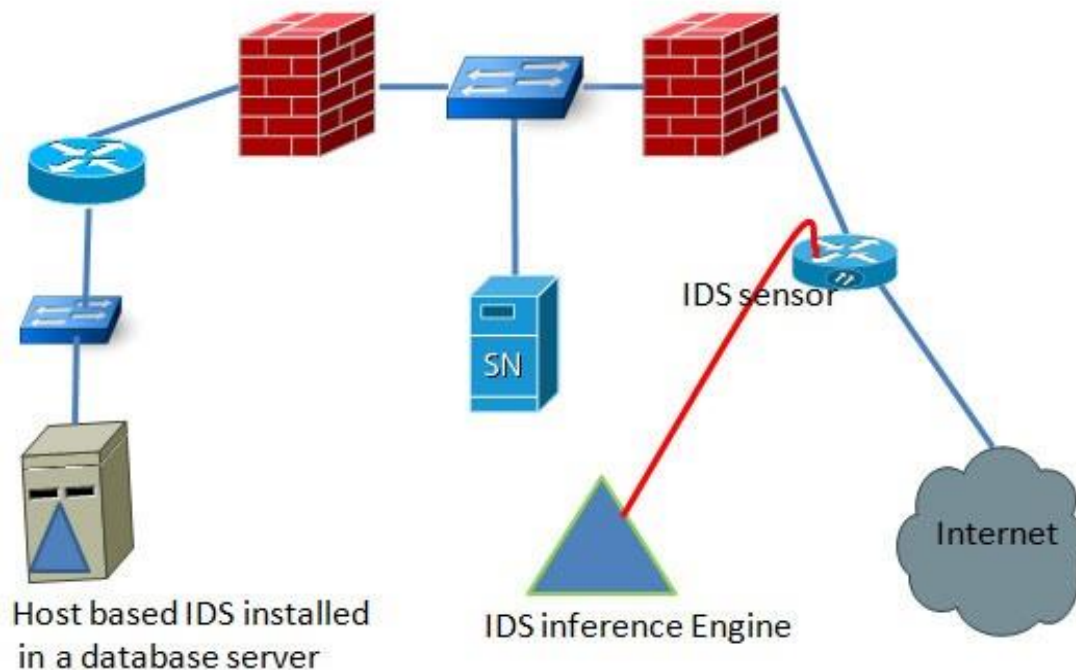
2.2.2 Συστήματα Ανίχνευσης παρεισφρήσεων

Τα Συστήματα Ανίχνευσης Παρεισφρήσεων (Intrusion Detection Systems) έχουν παθητικό ρόλο, αφού εντοπίζουν τις κακόβουλες ενέργειες, τις παραβιάσεις πολιτικών ασφαλείας και συντάσσουν αναφορές χωρίς να προβαίνουν σε καμία

ενέργεια αποτροπής. Συνήθως είναι μια συσκευή, ή εφαρμογή λογισμικού, ή ένας συνδυασμός υλικού και λογισμικού. Τα Συστήματα Ανίχνευσης Παρεισφρήσεων μπορούν να εποπτεύουν είτε ένα μεμονωμένο υπολογιστικό σύστημα (Host Based Intrusion Detection Systems, HIDS) είτε ένα ολόκληρο δίκτυο (Network Based Intrusion Detection Systems, NIDS). Κάθε φορά που το Σύστημα Ανίχνευσης Παρεισφρήσεων (Intrusion Detection Systems) ανιχνεύει μια παράξενη ή παράτυπη κίνηση ή παράτυπη ενέργεια στο υπολογιστικό σύστημα ή στο δίκτυο που εποπτεύει σε σχέση με τους κανόνες ασφαλείας που του έχουν καθορίσει, ετοιμάζει και μια σχετική αναφορά. Όλες αυτές οι αναφορές είναι διαθέσιμες στον υπεύθυνο ή στους υπεύθυνους διαχειριστές, οι οποίοι αφού τις αξιολογήσουν αποφασίζουν για το αν πρέπει να πάρουν μέτρα προς την ενίσχυση της ασφάλειας. Τα Συστήματα Ανίχνευσης Παρεισφρήσεων (Intrusion Detection Systems) κατηγοριοποιούνται σε δύο ομάδες:

1. Συστήματα Ανίχνευσης Παρεισφρήσεων Δικτύου (Network Based Intrusion Detection Systems, NIDSs).
2. Συστήματα Ανίχνευσης Παρεισφρήσεων υπολογιστικού συστήματος (Host Based Intrusion Detection Systems, HIDSs).

Όμως πολλές φορές σε ένα δίκτυο γίνεται χρήση και συνδυασμού τους (Εικόνα 2). Αξίζει να σημειωθεί ότι ορισμένα Συστήματα Ανίχνευσης Εισβολής έχουν την ικανότητα να ανιχνεύσουν απειλές έχοντας ως βάση διαφορετικούς κανόνες ασφαλείας για διαφορετικά τμήματα του δικτύου. Χωρίζοντας με αυτόν τον τρόπο το δίκτυο σε τμήματα με διαφορετικούς κανόνες ασφαλείας με βάση τις ιδιαιτερότητες και τις ανάγκες του, αυξάνετε η ευελιξία του συστήματος ανίχνευσης Παρεισφρήσεων αφού γίνεται χρήσιμο για μεγάλα δίκτυα.



Εικόνα 2. Συνδιασμός HIDS και NIDS ("Host Based IDS vs Network Based IDS | securitywing", 2012).

2.3 Μεθοδολογίες Ανίχνευσης Εισβολής

Οι μεθοδολογίες ανίχνευσης εισβολής είναι δύο κατηγοριών:

2.3.1 Ανίχνευση Διαταραχών

Η μεθοδολογία της ανίχνευσης διαταραχών (Anomaly Detection) έχει ως δεδομένη μια «φυσιολογική συμπεριφορά» του δικτύου, των υπολογιστικών συστημάτων και των χρηστών. Η φυσιολογική συμπεριφορά ορίζεται εντός κάποιων ορίων π.χ. μια βάση δεδομένων με τις συμπεριφορές χρηστών. Και όταν παρεκκλίνει από αυτή την συμπεριφορά λαμβάνοντας τιμές (μοτίβα) που βρίσκονται εκτός ορίων τα αποκαλεί ανωμαλίες, τις οποίες κατά προέκταση το Σύστημα Ανίχνευσης Παρεισφρήσεων (Intrusion Detection System) τις εκλαμβάνει ως εισβολή. Η παραπάνω μεθοδολογία είναι ιδανική για να εντοπίσει άγνωστες εισβολές αφού εντοπίζει τα συμπτώματα τους χωρίς να τις γνωρίζει, όμως υπάρχουν και περιπτώσεις που η μεθοδολογία έχει λάθος αποτελέσματα τέτοιες περιπτώσεις είναι (Chauhan, Mishra, & Kumar, 2011):

1.Διεργασίες που επιβαρύνουν το σύστημα(δίκτυο, υπολογιστή) και δεν χρησιμοποιούνται συχνά μπορεί σε συνδυασμό μεταξύ τους να επιφέρουν ακραίες τιμές και να εκληφθούν ως εισβολή (false positives).

2.Εισβολές που γίνονται χωρίς να διαταράσσουν την «φυσιολογική συμπεριφορά» δεν γίνονται αντιληπτές από το Σύστημα Ανίχνευσης Παρεισφρήσεων (false negatives).

Οπότε σημαντική παράμετρος στην σωστή λειτουργία ενός συστήματος ανίχνευσης εισβολής που χρησιμοποιεί anomaly detection είναι η επιλογή ορίων. Όστε να επιτυγχάνετε η βέλτιστη λειτουργία του Συστήματος Ανίχνευσης Παρεισφρήσεων(Intrusion Detection System).

2.3.2 Ανίχνευση Κακής Συμπεριφοράς

Η μεθοδολογία της ανίχνευσης κακής συμπεριφοράς (misuse detection) βασίζεται στη φιλοσοφία της αναγνώρισης προκαθορισμένων προτύπων γεγονότων τα οποία αντιστοιχούν σε κάποια γνωστή επίθεση. Αυτό επιτυγχάνετε κάνοντας χρήση της ανάλυσης υπογραφών και για αυτό τα πρότυπα αναφέρονται και ως signatures. Ένα πρότυπο μπορεί να περιγράψει μια εισβολή ή μια κατηγορία εισβολών. Οι βιβλιοθήκες των Συστημάτων Ανίχνευσης Παρεισφρήσεων(Intrusion Detection System) πρέπει να ενημερώνονται συνέχεια με καινούργια πρότυπα ώστε να είναι επικαιροποιημένες για να μην καθιστούν τα υπολογιστικά συστήματα ευάλωτα σε καινούργιες τεχνικές εισβολών. Η μεθοδολογία της ανίχνευσης κακής συμπεριφοράς (misuse detection) έχει το πλεονέκτημα ότι δεν παράγει πολλές λάθος διαγνώσεις εισβολών(false positives). Όμως έχει το μειονέκτημα ότι είναι αδύνατον να εντοπιστούν άγνωστες τεχνικές εισβολής καθώς και παραλλαγές γνωστών εισβολών.

Τα συστήματα ανίχνευσης Παρεισφρήσεων μπορούν να χρησιμοποιήσουν δύο τύπους μηχανισμών ώστε να αποφασίσουν αν υπάρχει εισβολή ή όχι, οι τύποι είναι οι ακόλουθοι:

- Ανάλυση με βάση τις υπογραφές(signature based intrusion detection system).
- Ανάλυση με βάση της ανωμαλίες (anomaly based intrusion detection system).

2.3.3 Συστήματα Ανίχνευσης Παρεισφρήσεων βασισμένα στην ανάλυση υπογραφών

Τα συστήματα ανίχνευσης Παρεισφρήσεων που είναι βασισμένα στην ανάλυση υπογραφών(signature based Intrusion Detection Systems) λειτουργούν με παρόμοια μεθοδολογία με τα προγράμματα αντιμετώπισης ιών(antivirus). Σε αυτά τα Συστήματα Ανίχνευσης Παρεισφρήσεων αρχικά ο κατασκευαστής αρχικά καθορίζει τους κανόνες ασφαλείας που τους αποκαλεί ως «υπογραφές» και στη συνέχεια τις ενημερώνει. Ως υπογραφές ορίζονται οι αλληλουχίες ενεργειών που μπορούν να εκληφθούν είτε ως επιθέσεις είτε ως πιθανές προσπάθειες εισβολής. Οι οποίες όταν εντοπίζονται δημιουργούν τις αντίστοιχες αναφορές ως προς τον υπεύθυνο π.χ. διαχειριστή. Η ομοιότητα των Συστημάτων Ανίχνευσης Παρεισφρήσεων που είναι βασισμένα στην ανάλυση με βάση τις υπογραφές σε σχέση με την λειτουργία των προγραμμάτων αντιμετώπισης ιών προϋποθέτει και το μειονέκτημα των προγραμμάτων αντιμετώπισης ιών, αφού ενώ είναι πολύ αποτελεσματικά στο να εντοπίζουν γνωστές τεχνικές επίθεσης παραμένουν όμως ευάλωτα σε άγνωστες τεχνικές επίθεσης(zero days). Επίσης όσο πιο μεγάλο όγκο έχει η βάση δεδομένων με τις υπογραφές τόσο αυξάνετε ο φόρτος εργασίας του επεξεργαστή για την ανάλυση των υπογραφών αυξάνοντας την πολυπλοκότητα του υπολογιστικού συστήματος και το κόστος του (Brox, 2002).

2.3.4 Συστήματα Ανίχνευσης Παρεισφρήσεων βασισμένα στην ανάλυση ανωμαλιών

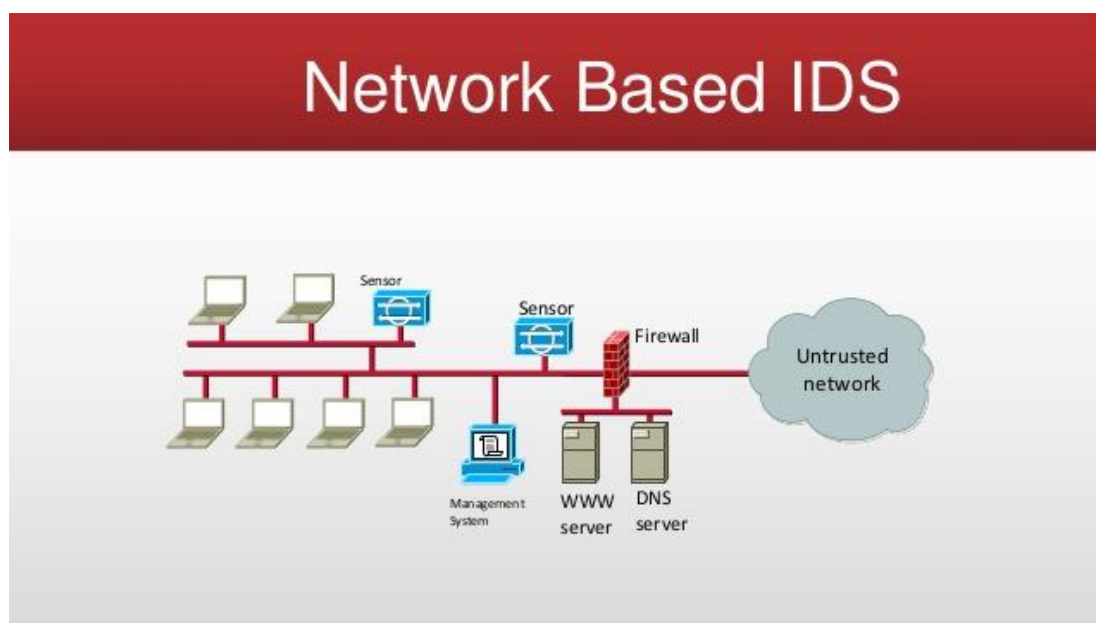
Τα Συστήματα Ανίχνευσης Παρεισφρήσεων βασισμένα στην ανάλυση ανωμαλιών(anomaly based Intrusion Detection Systems) έχουν την ικανότητα

να εντοπίζουν έως τώρα άγνωστες επιθέσεις. Ο τρόπος που το επιτυγχάνουν είναι δημιουργώντας ένα στατιστικό μοντέλο που περιγράφει την φυσιολογική λειτουργία(Συνήθης συμπεριφορά) του υπολογιστικού συστήματος και στην συνέχεια γίνονται έλεγχοι για τυχόν αποκλίσεις από το αρχικό στατιστικό μοντέλο που υπερβαίνουν τα όρια ανοχής τότε εκλαμβάνονται από το Σύστημα Ανίχνευσης Παρεισφρήσεων είτε ως εισβολές είτε ως προσπάθειες εισβολής. Στοιχεία που μπορεί να εμπεριέχονται στο στατιστικό μοντέλο είναι ο μέσος αριθμός πακέτων που δέχεται ένας υπολογιστής, η μέση διάρκεια μιας συνόδου, μέσος αριθμός συνδέσεων που επιχειρείτε κ.α. Ένα σημαντικό μειονέκτημα των Συστημάτων Ανίχνευσης Παρεισφρήσεων που είναι βασισμένα στην ανάλυση με βάση της ανωμαλίες, είναι ότι πολλές φορές εκλαμβάνουν ως επιθέσεις ενώ δεν είναι διεργασίες που επιβαρύνουν το σύστημα και δεν χρησιμοποιούνται συχνά ή σπάνιες συμπεριφορές χρηστών αφού ο συνδυασμός τους μπορεί να επιφέρει ακραίες τιμές(positive false) (Nascimento, & Correia, n.d.). Επιπλέον πέραν των ανθρώπινων πόρων που χρειάζονται για να φιλτράρουν τους λάθος συναγερμούς(false positives), χρειάζονται και περισσότεροι υλικοί πόροι αφού είναι μεγάλη η πολυπλοκότητα του υπολογιστικού συστήματος και απαιτούν πόρους όπως υψηλής απόδοσης επεξεργαστές. Τέλος οι διαχειριστές του συστήματος για την καλύτερη διαχείριση των υπολογιστικών πόρων θα πρέπει να ακολουθήσουν και διάφορες τεχνικές για να μειώσουν τον φόρτο διακίνησης δεδομένων στο δίκτυο(bandwidth), όπως την τοποθέτηση των αισθητήρων ανώμαλης λειτουργίας κοντά στον server και στο δίκτυο που παρακολουθείτε(Brox, 2002).

2.3.5 Συστήματα ανίχνευσης Παρεισφρήσεων δικτύου

Τα Συστήματα Ανίχνευσης Παρεισφρήσεων δικτύου (Network Based Intrusion Detection Systems, NIDS) συλλέγουν πληροφορίες για την κίνηση του δικτύου(network traffic) και όχι για την κάθε μονάδα(host) του δικτύου ξεχωριστά. Για να συλλέξουν τις πληροφορίες κάνουν χρήση αισθητήρων δικτύου, στη συνέχεια αποστέλλουν τις πληροφορίες στον σταθμό ελέγχου(management system) ο οποίος τις αναλύει για να εντοπίσει ίχνη εισβολής(Εικόνα 3). Οι αισθητήρες δικτύου δεν περιορίζονται μόνο στην δική τους IP διεύθυνση λαμβάνουν ολόκληρο τον φόρτο του δικτύου που διέρχεται

από εκείνο το σημείο, γι' αυτό τον λόγο έχει μεγάλη σημασία η στρατηγική τοποθέτησης των αισθητήρων για την αποτελεσματικότητα του Συστήματος Ανίχνευσης Παρεισφρήσεων δικτύου. Αξίζει να σημειωθεί ότι τα Συστήματα Ανίχνευσης Παρεισφρήσεων δεν μπορούν να ελέγξουν την κίνηση του δικτύου σε τηλεφωνικές γραμμές. Ο κάθε αισθητήρας έχει μια δικτυακή διεπαφή για να συνδέετε με ένα σημείο του δικτύου ώστε να παρακολουθεί την κίνηση για ασυνήθιστες συμπεριφορές τις οποίες μεταβιβάζει στον σταθμό ελέγχου ο οποίος τα αναλύει και στην συνέχεια τα παρουσιάζει με μορφή αναφοράς στον υπεύθυνο ασφαλείας που είτε διαχειρίζεται τον σταθμό ελέγχου είτε λαμβάνει αναφορές από αυτόν("Host- vs. Network-Based Intrusion Detection Systems", 2000). Ένα σημαντικό ελάττωμα των Συστημάτων Ανίχνευσης Παρεισφρήσεων δικτύου είναι η εμφάνιση αρκετών ψευδών συναγερμών(false positives) εξαιτίας της συνεχώς αυξανόμενης ταχύτητας μετάδοσης των δεδομένων και του τεράστιου όγκου των δεδομένων που διαχειρίζονται(Chen, n.d.).



Εικόνα 3. Συστήματα ανίχνευσης Παρεισφρήσεων δικτύου ("Intrusion detection and prevention system", n.d.).

Τα Πλεονεκτήματα των Συστημάτων Ανίχνευσης Παρεισφρήσεων Δικτύου (Network Based Intrusion Detection Systems, NIDS) είναι τα εξής:

1. Τα Συστήματα Ανίχνευσης Παρεισφρήσεων Δικτύου είναι ιδανικά για να εντοπίζουν στο ξεκίνημα τους επιθέσεις που στοχεύουν στο δίκτυο όπως

συμβαίνει με τις επιθέσεις που στοχεύουν στην κατανάλωση του εύρους ζώνης (bandwidth) με ιούς τύπου worm σε αντίθεση με τα Συστήματα Ανίχνευσης Παρεισφρήσεων Τοπικών Υπολογιστικών Συστημάτων (Host Based Intrusion Detection Systems, NIDS)(Chen, n.d.).

2. Τα Συστήματα Ανίχνευσης Παρεισφρήσεων Δικτύου έχουν πιο εύκολη εγκατάσταση σε σχέση τα Συστήματα Ανίχνευσης Παρεισφρήσεων Τοπικών Υπολογιστικών Συστημάτων (Host Based Intrusion Detection Systems, NIDS). Ενώ επίσης έχουν χαμηλές απαιτήσεις ενημέρωσης από τον διαχειριστή του συστήματος ασφαλείας αφού λειτουργούν συνήθως με αυτοματοποιημένες διαδικασίες που έχουν σχεδιαστεί από τον κατασκευαστή τους. Δυσχεραίνοντας με αυτόν τον τρόπο πιθανές προσπάθειες πρόσβασης του κακόβουλου χρήστη αποκτώντας δικαιώματα διαχειριστή ώστε: να απενεργοποιήσει το Σύστημα Ανίχνευσης Παρεισφρήσεων Δικτύου, να το καταστήσει ακίνδυνο για τις κακόβουλες ενέργειες του ή ακόμη και να «διαβάσει» τον τρόπο λειτουργίας του Συστήματος Ανίχνευσης Παρεισφρήσεων Δικτύου ώστε να σχεδιάσει καλύτερα τις επόμενες επιθέσεις του σε όμοια συστήματα.
3. Τα Συστήματα Ανίχνευσης Παρεισφρήσεων Δικτύου παρουσιάζουν υψηλή προσαρμοστικότητα σε τυχόν διαφοροποιήσεις του δικτύου π.χ. αλλαγή τοπολογίας, αλλαγές στην βαθμονόμηση της σημαντικότητας των τοπικών υπολογιστικών συστημάτων του δικτύου. Αφού πολλές φορές το μόνο που απαιτητέ είναι η αλλαγή στις ρυθμίσεις των καρτών δικτύου των συσκευών του δικτύου("Host- vs. Network-Based Intrusion Detection Systems", 2000).
4. Τα Συστήματα Ανίχνευσης Παρεισφρήσεων σε περίπτωση που σταματήσουν για οποιονδήποτε λόγο να λειτουργούν δεν διακόπτουν την απρόσκοπτη λειτουργία του δικτύου και δεν επηρεάζει ούτε γίνετε αντιληπτό από του υπαλλήλους που χρησιμοποιούν τοπικά υπολογιστικά συστήματα του δικτύου. Γίνετε αντιληπτό μόνο από τον διαχειριστή του συστήματος.

Τα Μειονεκτήματα των συστημάτων ανίχνευσης Παρεισφρήσεων δικτύου είναι τα εξής:

1. Οι αισθητήρες των Συστημάτων Ανίχνευσης Παρεισφρήσεων Δικτύου λαμβάνουν δεδομένα για την κίνηση του δικτύου (network traffic) μόνο από εκείνο το σημείο που είναι τοποθετημένοι. Συνεπώς εάν η επίθεση γίνεται σε συγκεκριμένο κομμάτι του δικτύου που δεν υπάρχει αισθητήρας δεν θα γίνει αντιληπτή. Η λύση σε αυτό το πρόβλημα είναι η τοποθέτηση αισθητήρων σε κάθε κομμάτι του δικτύου, αυξάνοντας σημαντικά το κόστος του εξοπλισμού και καθιστώντας το σύστημα πολυδαίδαλο.
2. Τα Συστήματα Ανίχνευσης Παρεισφρήσεων Δικτύου τις περισσότερες φορές χρησιμοποιούν ανάλυση υπογραφών (signatures) για τον εντοπισμό εισβολών. Αυτό όμως εμπεριέχει τον κίνδυνο να είναι απόλυτος αδύνατον τα Συστήματα Ανίχνευσης Παρεισφρήσεων δικτύου να εντοπίσουν άγνωστες τεχνικές εισβολής("Host- vs. Network-Based Intrusion Detection Systems", 2000).
3. Τα Συστήματα Ανίχνευσης Παρεισφρήσεων Δικτύου πολλές φορές διακινούν αρκετά ογκώδεις ποσότητες δεδομένων προς τον σταθμό ελέγχου. Για να αντιμετωπιστεί αυτό το πρόβλημα χρησιμοποιούνται αλγόριθμοι φιλτραρίσματος των δεδομένων. Όμως τότε παρουσιάζετε το πρόβλημα της αποσιώπησης πολύτιμων δεδομένων με αποτέλεσμα να μην εντοπίζονται εισβολές.
4. Τα Συστήματα Ανίχνευσης Παρεισφρήσεων Δικτύου αντιμετωπίζουν προβλήματα εντοπισμού εισβολών όταν η κίνηση στο δίκτυο είναι κρυπτογραφημένη. Αφού σε αυτή την περίπτωση δεν μπορεί να πραγματοποιηθεί σάρωση (scan) των πρωτοκόλλων ή του περιεχόμενου των πακέτων. Επίσης η φύση των διακοπών (switch) δυσχεραίνει επίσης την παρακολούθηση του δικτύου, αφού σε αντίθεση με τους

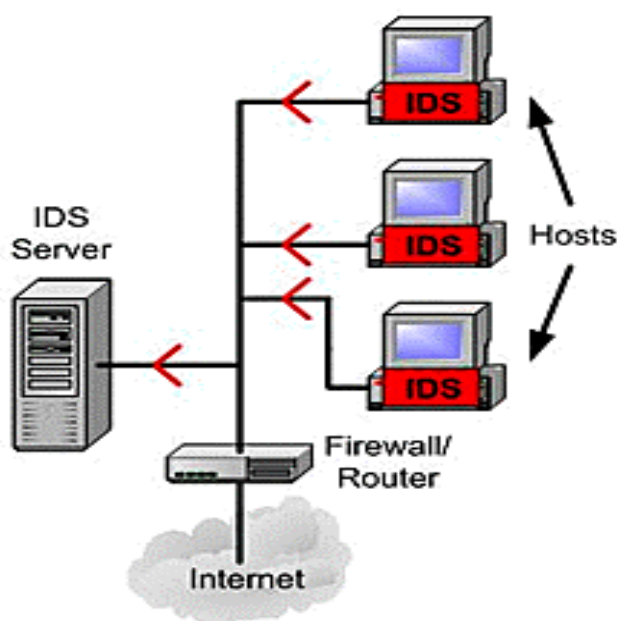
διανομείς(hubs) απομονώνει τα τοπικά υπολογιστικά συστήματα ώστε να έχουν πρόσβαση στις πληροφορίες που απευθύνονται μόνο σε αυτά. Γι' αυτό το λόγο μερικοί σύγχρονοι διακόπτες(switch) έχουν μία θύρα που επιτρέπει την παρακολούθηση και την σάρωση(Das, & Sarkar, 2014).

5. Τα Συστήματα Ανίχνευσης Παρεισφρήσεων Δικτύου είναι πολύ συγκεντρωτικά αφού ο έλεγχος και η καταγραφή γίνετε στον Διακομιστή (Server), δεν γίνετε ξεχωριστά όπως στα Τοπικά Συστήματα Ανίχνευσης Παρεισφρήσεων(Host Based Intrusion Detection Systems, HIDS) όπου ο έλεγχος και η καταγραφή γίνετε στον κάθε υπολογιστή στον οποίον είναι εγκατεστημένο το Τοπικό Σύστημα Ανίχνευσης Παρεισφρήσεων. Οπότε σε περίπτωση που υπάρξει πρόβλημα στην λειτουργία ή στις ρυθμίσεις του Server καταρρέει όλο το δίκτυο ασφαλείας.

2.3.6 Τοπικά Συστήματα Ανίχνευσης Παρεισφρήσεων

Τα Τοπικά Συστήματα Ανίχνευσης Παρεισφρήσεων (Host Based Intrusion Detection Systems, HIDS) εστιάζουν στην παρακολούθηση και ανάλυση δραστηριοτήτων που λαμβάνουν χώρα στα τοπικά υπολογιστικά συστήματα(host) ή στα μεμονωμένα συστήματα για ίχνη εισβολής(Das, & Sarkar, 2014). Τα Τοπικά Συστήματα Ανίχνευσης Παρεισφρήσεων πολλές φορές αποκαλούνται και ως αισθητήρες και συνήθως εγκαθίστανται σε υπολογιστές ή συσκευές που θεωρούνται πιθανοί στόχοι επίθεσης(Εικόνα 4). Οπότε για κάθε μεμονωμένο υπολογιστικό σύστημα(host) ή συσκευή που θέλουμε να παρακολουθήσουμε χρειαζόμαστε και έναν αισθητήρα που θα έχει εγκατασταθεί και προσαρμοστεί σε αυτό. Τα δεδομένα που συλλέγουν τα Τοπικά Συστήματα Ανίχνευσης Παρεισφρήσεων αποκαλούνται ίχνη ελέγχου. Τα οποία συλλέγονται αναλύοντας τα αρχεία καταγραφής συστήματος(system logs), καταγραφές που αναφέρονται σε διεργασίες του λειτουργικού συστήματος καθώς και από το περιεχόμενο αντικειμένων που δεν ανταποκρίνεται στο πρότυπο ασφαλείας που έχει καθοριστεί("Host- vs. Network-Based Intrusion Detection Systems", 2000). Επίσης παρέχουν πιο ακριβή πληροφορία για την παρουσία ή όχι κάποιας εισβολής. Αφού παρακολουθούν πιθανές παραβιάσεις εμπιστευτικότητας, ακεραιότητας και

διαθεσιμότητας π.χ. εάν κάποιος χρήστης ή εφαρμογή αποκτά πρόσβαση σε πόρους και εάν ξεκινάει τροποποίηση των πόρων χωρίς να έχει δικαίωμα. Είναι πιο αποτελεσματικά έναντι των προγραμμάτων αντιμετώπισης ιών(antivirus) αφού σε αντίθεση με αυτά, τα Τοπικά Συστήματα Ανίχνευσης Παρεισφρήσεων ελέγχουν το υπολογιστικό σύστημα και για επιθέσεις υπερχειλίσης της μνήμης (buffer overflow) (Das, & Sarkar, 2014).



Εικόνα 4. Τοπικά Συστήματα Ανίχνευσης Παρεισφρήσεων(Host Based IDS).

Τα πλεονεκτήματα των Τοπικών Συστημάτων Ανίχνευσης Παρεισφρήσεων(Host Based Intrusion Detection Systems, HIDS) είναι τα εξής:

1. Τα Τοπικά Συστήματα Ανίχνευσης Παρεισφρήσεων πραγματοποιούν μια πιο ολοκληρωμένη καταγραφή των εισβολών. Αφού συνήθως καταγράφουν με ακρίβεια τις κινήσεις του κακόβουλου χρήστη π.χ. εντολές που εκτέλεσε, κενά ασφαλείας που προσπάθησε να εκμεταλλευτεί κ.α.
2. Τα Τοπικά Συστήματα Ανίχνευσης επειδή είναι πολύ συγκεκριμένα στο τι ελέγχουν πολύ σπάνια προβαίνουν σε λάθος συναγερμούς(false positives).
3. Τα Τοπικά Συστήματα Ανίχνευσης Παρεισφρήσεων είναι ιδανικά για πληροφοριακά συστήματα που το εύρος ζώνης του δικτύου δεν επαρκεί

για την επικοινωνία του αισθητήρα με τον σταθμό ελέγχου. Καθώς επίσης και σε περιπτώσεις που οι αισθητήρες παράγουν μεγάλο όγκο δεδομένων σε βαθμό που καθιστούν τα Συστήματα Ανίχνευσης Εισβολής Δικτύου (Network Based Intrusion Detection Systems, NIDS) αναξιόπιστα.

4. Τα Τοπικά Συστήματα Ανίχνευσης Παρεισφρήσεων έχουν χαμηλές απαιτήσεις ενημέρωσης από τον διαχειριστή αφού λειτουργούν συνήθως με αυτοματοποιημένες διαδικασίες που έχουν σχεδιαστεί από τον κατασκευαστή τους. Σταματώντας με αυτόν τον τρόπο τις προσπάθειες εξουδετέρωσης από έναν κακόβουλο χρήστη που έχει αποκτήσει δικαιώματα διαχειριστή.

Τα μειονεκτήματα των Τοπικών Συστημάτων Ανίχνευσης Παρεισφρήσεων είναι τα εξής:

1. Τα Τοπικά Συστήματα Ανίχνευσης Παρεισφρήσεων προϋποθέτουν εγκατάσταση σε κάθε υπολογιστική συσκευή του δικτύου την οποία θέλουμε να παρακολουθήσουμε και κατά προέκταση να προστατεύουμε. Αυξάνοντας με αυτόν τον τρόπο το κόστος και τον χρόνο εγκατάστασης. Οπότε πολλές φορές για να περιοριστεί κυρίως το κόστος πολλοί υπεύθυνοι ασφαλείας(διαχειριστές) αφήνουν μια ή και περισσότερες υπολογιστικές συσκευές από ένα δίκτυο χωρίς Τοπικό Σύστημα Ανίχνευσης Παρεισφρήσεων. Σε αυτήν την περίπτωση δημιουργούν εν γνώσει τους μια ολοφάνερη ευπάθεια, η οποία σε περίπτωση όπου γίνει κάποια παραβίαση σε αυτήν την υπολογιστική συσκευή δεν θα γίνει αντιληπτή από τον υπεύθυνο ασφαλείας. Επίσης από την στιγμή που εγκαθιστούν σε υπολογιστές χρηστών προϋποθέτει την ύπαρξη χωρητικότητας σε αυτούς.
2. Τα Τοπικά Συστήματα Ανίχνευσης Παρεισφρήσεων έχουν κάποιες ελάχιστες απαιτήσεις συστήματος για να εγκατασταθούν σε μια υπολογιστική συσκευή που θέλουμε να προστατεύσουμε. Οπότε σε

περίπτωση που δεν τηρούνται από μια ή και περισσότερες υπολογιστικές συσκευές οι ελάχιστες απαιτήσεις όπως: σε απόδοση επεξεργαστή (CPU), μνήμης τυχαίας προσπέλασης (RAM) και χωρητικότητας στον σκληρό δίσκο δεν γίνεται να παρακολουθείτε η κίνηση στις συγκεκριμένες υπολογιστικές συσκευές. Και σε περίπτωση που θελήσει ο υπεύθυνος ασφαλείας να εγκαταστήσει ένα Τοπικό Σύστημα Ανίχνευσης Παρεισφρήσεων (Host Based Intrusion Detection Systems, HIDS) σε κάθε μια από αυτές τις υπολογιστικές συσκευές θα αναγκαστεί να τις αναβαθμίσει αυξάνοντας έτσι το χρηματικό κόστος.

3. Τα Τοπικά Συστήματα Ανίχνευσης Παρεισφρήσεων αναφέραμε ότι για να καλύπτουν με ασφάλεια ένα δίκτυο πρέπει να είναι εγκατεστημένα σε όσο το δυνατόν περισσότερες υπολογιστικές συσκευές. Αυτό όμως δημιουργεί ένα άλλο πρόβλημα που έχει να κάνει με την ανάλυση δεδομένων στον Server, αφού του αυξάνει το φόρτο εργασίας και δημιουργεί καθυστέρηση στα αποτελέσματα της ανάλυσης π.χ. ένα δίκτυο με 254 host που έχουν εγκατεστημένο IDS οι 250 σε σχέση με ένα με 6 host που έχουν όλοι εγκατεστημένο IDS το πρώτο έχει ολοφάνερα μεγαλύτερη καθυστέρηση στην ανάλυση των δεδομένων εάν έχουν παρόμοιο Server και δίκτυο.
4. Τα Τοπικά Συστήματα Ανίχνευσης Παρεισφρήσεων τις περισσότερες φορές χρησιμοποιούν ανάλυση υπογραφών (signatures) για τον εντοπισμό εισβολών. Αυτό όμως εμπεριέχει τον κίνδυνο να είναι απόλυτος αδύνατον τα Συστήματα Ανίχνευσης Παρεισφρήσεων δικτύου να εντοπίσουν άγνωστες τεχνικές εισβολής ("Host- vs. Network-Based Intrusion Detection Systems", 2000).

Κεφάλαιο 3

Εξόρυξη Δεδομένων

Με τον όρο Μηχανική Μάθηση(Machine Learning) περιγράφουμε μια μορφή τεχνικής νοημοσύνης που αναλύει δεδομένα με σκοπό την εξόρυξη άγνωστης γνώσης. Είναι ένας συνεχώς αναπτυσσόμενος κλάδος που βρίσκει εφαρμογή σε ζητήματα όπως: η αναγνώριση προσώπων, η αναγνώριση φωνής, η ρομποτική κ.α.(Alpaydm, 2013).

Η εφαρμογή της Μηχανικής Μάθησης όταν πραγματοποιείτε αναλύοντας βάσης δεδομένων με σκοπό την ανακάλυψη πληροφοριών αποκαλείτε εξόρυξη δεδομένων(Data Mining)(Alpaydm, 2013). Ανάλυση των δεδομένων αυτών επιτυγχάνετε με την χρήση αλγορίθμων. Σήμερα η Εξόρυξη Δεδομένων είναι ένα χρήσιμο εργαλείο που δίνει την δυνατότητα αναλύοντας συνήθως τεράστιο όγκο δεδομένων να εξάγουμε όλες τις κρυμμένες πληροφορίες. Είναι ένας συνεχώς αναπτυσσόμενος κλάδος αφού βρίσκει εφαρμογή σε πεδία όπως η δικανική υπολογιστών(computer forensics), το εμπόριο (marketing), η επιστημονική έρευνα κ.α.

Η Εξόρυξη Δεδομένων(Data Mining) διακρίνεται σε τρεις διαδικασίες:

1. Την Εξερεύνηση(exploration) η οποία περιλαμβάνει την προετοιμασία των δεδομένων.
2. Την Δημιουργία και Επικύρωση του Μοντέλου (Model Building and Validation) η οποία περιλαμβάνει την επιλογή μεθόδου Εξόρυξης Δεδομένων(Data Mining).

3. Την Ανάπτυξη(Deployment) η οποία περιλαμβάνει την χρήση των δεδομένων που εξάχθηκαν από την μέθοδο Εξόρυξης Δεδομένων που επιλέχτηκε, ώστε να παραχθεί αποτέλεσμα(Uack, 2013).

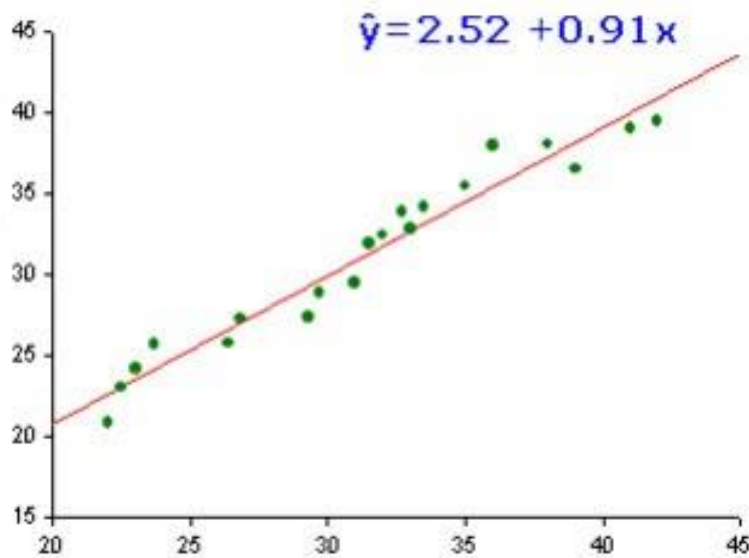
3.1 Μέθοδοι Εξόρυξης Δεδομένων

Η Εξόρυξη Δεδομένων(Data Mining) συνήθως χωρίζετε σε τέσσερεις μεθόδους συγκεκριμένα είναι οι εξής(Chauhan, Mishra, & Kumar, 2011):

1. Η Ομαδοποίηση (Clustering)
2. Η Κατηγοριοποίηση (Classification)
3. Η Παλινδρόμηση (Regression)
4. Η Σχέση Κανόνα Μάθησης (Association rule learning)

3.1.1 Παλινδρόμηση

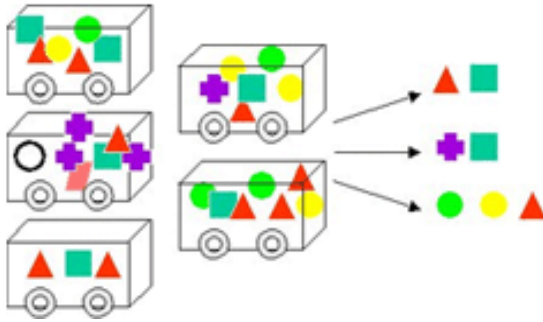
Η Παλινδρόμηση είναι μια τεχνική ταξινόμησης δεδομένων η οποία περιγράφει την εφαρμογή μιας συνάρτησης σε ένα σύνολο δεδομένων, ώστε να ταξινομήσει τα δεδομένα με το ελάχιστο σφάλμα(Chauhan, Mishra, & Kumar, 2011). Η πιο γνωστή συνάρτηση παλινδρόμησης είναι η γραμμική συνάρτηση $y=mx+b$ όπου m και b είναι σταθερές τιμές κάθε φορά και x η τιμή που αντλείτε από το σύνολο δεδομένων ώστε να προβλεφθεί με το ελάχιστο σφάλμα το y . Στην Εικόνα 5 φαίνετε η γραμμική παλινδρόμηση για $y=0,91*x+2,52$. Επίσης βρίσκει εφαρμογή σε προβλέψεις τιμών μετοχών, ισοτιμίας νομισμάτων κ.α.



Εικόνα 5. Παλινδρόμηση (Regression) (Calbimonte, 2014).

3.1.2 Σχέση κανόνα μάθησης

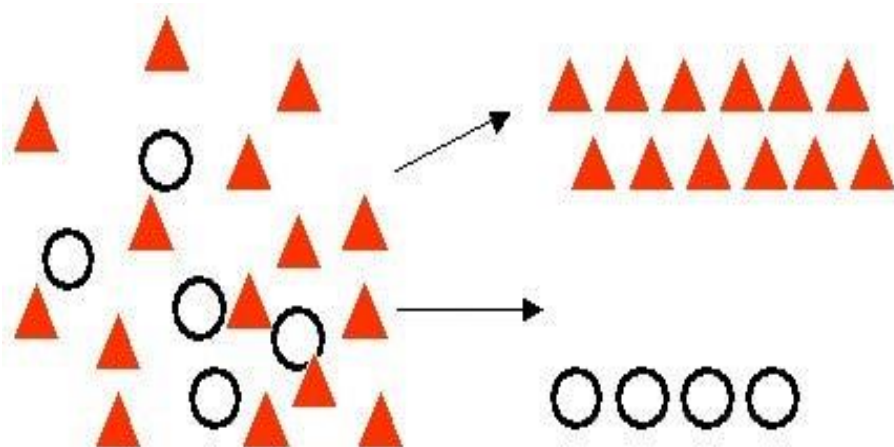
Η Σχέση κανόνα μάθησης (Association rule learning) είναι μια τεχνική η οποία περιγράφει την διαδικασία ανεύρεσης σχέσεων μεταξύ των μεταβλητών των δεδομένων. Χρησιμοποιείται συνήθως στο εμπόριο συλλέγοντας δεδομένα για τις αγοραστικές συνήθειες των πελατών. Για παράδειγμα παρατηρείτε ότι σε μια υπεραγορά(supermarket) οι καταναλωτές που αγοράζουν αψρό ξυρίσματος, έχουν πιθανότητα 70% να αγοράσουν και ξυραφάκια. Αυτό βρίσκει πρακτική εφαρμογή σε πολλούς τομείς του εμπορίου ενδεικτικά σε ένα ηλεκτρονικό κατάστημα όπου ανάλογα με το προϊόν που βλέπουμε ή αγοράζουμε μας προτείνει προϊόντα που κοίταξαν ή αγόρασαν άλλοι πελάτες μετά από αυτό το προϊόν, είναι ένα αποτέλεσμα συλλογής δεδομένων με χρήση σχέση κανόνα μάθησης. Στην Εικόνα 6 φαίνετε το φιλτράρισμα και η ταξινόμηση του δείγματος με την χρήση σχέση κανόνα μάθησης. Αφού το τρίγωνο συνδυάζεται με το τετράγωνο, ο σταυρός με το τετράγωνο και το ο κίτρινος κύκλος με τον πράσινο και το τρίγωνο.



Εικόνα 6. Σχέση κανόνα μάθησης (Association rule learning) (Calbimonte, 2014).

3.1.3 Κατηγοριοποίηση

Η Κατηγοριοποίηση(classification) είναι μια τεχνική ταξινόμησης δεδομένων, με την οποία τα δεδομένα του συνόλου δεδομένων, ταξινομούνται με βάση τις καινούργιες πληροφορίες που εξάγονται από τον συνδυασμό των υφιστάμενων πληροφοριών που εμπεριέχονται αρχικά σε αυτά. Το αποτέλεσμα της κατηγοριοποίησης σε αντίθεση με την ομαδοποίηση είναι δύο ομάδες π.χ. ένα ηλεκτρονικό μήνυμα θα κατηγοριοποιηθεί ως ανεπιθύμητη αλληλογραφία ναι ή όχι(δύο ομάδες). Στο παράδειγμα μας το αποτέλεσμα μπορεί να εξαρτηθεί από συνδυασμό πολλών παραμέτρων όπως συχνότητα λήψης ηλεκτρονικού μηνύματος, περιεχόμενο ηλεκτρονικού μηνύματος κ.α. (Chauhan, Mishra, & Kumar, 2011). Στην Εικόνα 7 φαίνεται το φιλτράρισμα και η ταξινόμηση του δείγματος σε 2 ομάδες.

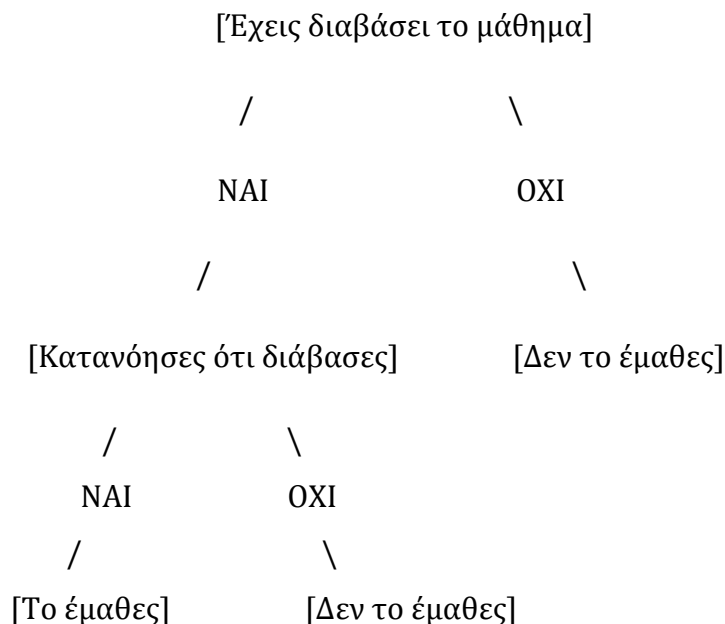


Εικόνα. 7 Κατηγοριοποίηση (Classification)("Oracle Advanced Analytics Data Mining Algorithms and Functions SQL API", n.d.).

Οι Τεχνικές Κατηγοριοποίησης(classification) είναι οι εξής(Chauhan, Mishra, & Kumar, 2011):

1. Τα Δέντρα Απόφασης (decision trees) τα οποία ταξινομούν τα δεδομένα σε διάφορα επίπεδα (levels) αποφάσεων ώστε να φτάσουν στην τελική απόφαση. Έτσι ώστε η δομή του δέντρου με τους κόμβους «ρίζας» και τα «φύλλα» να παριστάνουν επίπεδα απόφασης. Ένα δέντρο απόφασης που περιλαμβάνει διακριτές ετικέτες αποκαλείτε δέντρο ταξινόμησης, ενώ ένα δέντρο απόφασης το οποίο περιλαμβάνει συνεχής ετικέτες αποκαλείτε δέντρο παλινδρόμησης. Τα δέντρα απόφασης κάνουν χρήση των αλγορίθμων ID3 και J48(Chauhan, Mishra, & Kumar, 2011).

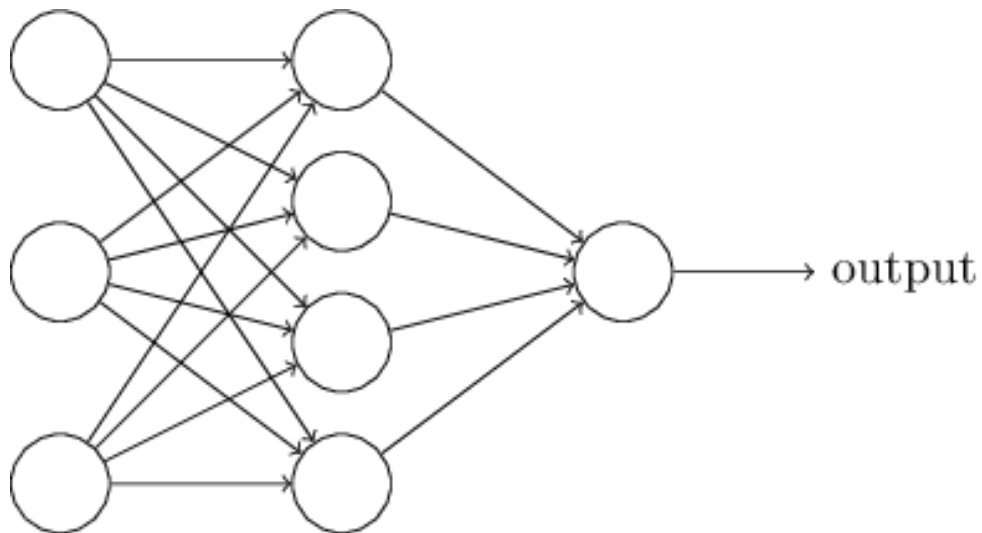
Ένα απλό παράδειγμα δέντρου απόφασης είναι το εξής:



2. Οι Μηχανές Υποστήριξης Διανυσμάτων(Support Vector Machines, SVM) αναπτύχθηκαν από τον Vladimir Vapnik και σχεδιάστηκαν για δυαδική ταξινόμηση. Οι Μηχανές Υποστήριξης Διανυσμάτων(Support Vector Machines) εντάσσουν τα δεδομένα σε μια υπέρ-κατηγορία (υπέρ-διάσταση) και στην συνέχεια τα διαχωρίζουν σε υποκατηγορίες(διαστάσεις).
3. Η Ασαφής λογική (Fuzzy Logic) επεξεργάζεται τα δεδομένα του δικτύου και περιγράφει μέτρα για την ανίχνευση σημαντικών ανωμαλιών στο σύστημα. Για να το επιτύχει αυτό διαβαθμίζει τα δεδομένα του δικτύου

σε μια κλίμακα αληθείας από 0 έως 1 (συμπεριλαμβάνοντας και τις ενδιάμεσες τιμές) (Chauhan, Mishra, & Kumar, 2011).

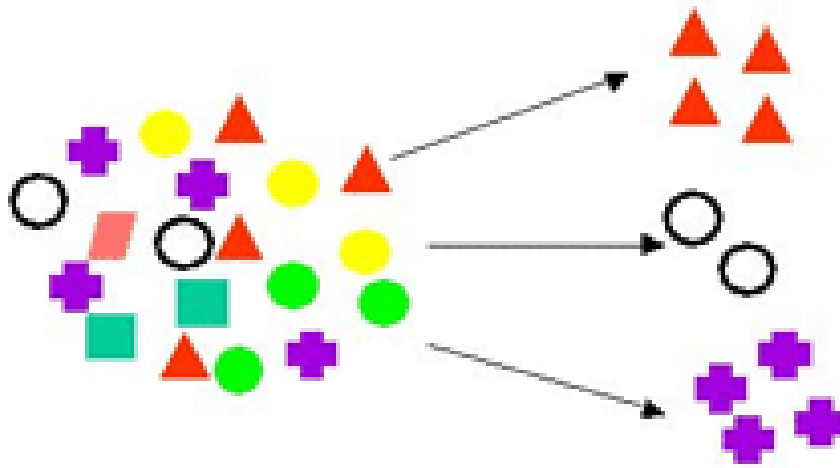
4. Οι τεχνικές Naive Bayes ταξινόμησης χρησιμοποιούν μια αναβάθμιση του θεωρήματος Bayes, το οποίο βασίζεται σε στατιστικές εξαρτίσεις και σχέσεις μεταξύ μεταβλητών. Οι Bayesian ταξινομητές είναι λιγότερο επιρρεπείς σε λάθη. Αφού θεωρούν ότι η επίδραση της τιμής ενός γνωρίσματος σε μια κατηγορία δεν επηρεάζει τα υπόλοιπα γνωρίσματα. Θωρακίζοντας με αυτόν τον τρόπο την κατά συνθήκη ανεξαρτησία της κατηγορίας. Οπότε σε ένα δείγμα δεδομένων $X=(x_1,x_2,...,x_n)$ τα όποια ανήκει σε m κατηγορίες ο αλγόριθμος Naive Bayes αναζητεί την μέγιστη πιθανότητα να ανήκει ένα x σε μια κατηγορία C . Αυτό εκφράζεται με τον τύπο $P(X|C_i)P(C_i) > P(X|C_j)P(C_j)$ για κάθε i, j να ανήκουν στο διάστημα $(1,m)$ (Joshi, 2012).
5. Τα Νευρωνικά δίκτυα (Neural Networks) είναι συστήματα που σχεδιάζονται με βάση τον τρόπο λειτουργίας του ανθρώπινου εγκεφάλου. Άλλωστε η ονομασία τους προέρχεται από τους νευρώνες του ανθρώπινου εγκεφάλου ο οποίος αποτελείται από εκατομμύρια νευρώνες που συνδέονται μεταξύ τους. Των ρόλο των νευρώνων στο Νευρωνικό δίκτυο των αναπαριστά ένα σύνολο από μονάδες εισόδου και εξόδου που συνδέονται μεταξύ τους (Joshi, 2012). Η φιλοσοφία των νευρωνικών δικτύων (Neural Networks) είναι ότι εκπαιδεύουμε το σύστημα ώστε να προβλέπει την επόμενη κίνηση ενός χρήστη σε μια αλληλουχία κινήσεων (εντολές, ενέργειες), με βάση το ιστορικό κινήσεων του. Όστε ανάλογα με τον βαθμό απόκλισης των επόμενων κινήσεων να ενεργοποιούνται συναγερμοί (alerts). Μπορούμε να το προσαρμόσουμε σε διάφορους χρήστες. Όμως είναι μια μέθοδος που σίγουρα αρκετές φορές εκλαμβάνει ενέργειες ως επιθέσεις ενώ δεν υφίστανται (false positives). Καθώς επίσης δεν εντοπίζει επιθέσεις που υφίστανται (false negatives). Ένα παράδειγμα Νευρωνικού Δικτύου παρουσιάζετε στην Εικόνα 8.



Εικόνα 8. Αναπαράσταση Νευρωνικού Δικτύου(Neural Network) (Nielsen, 2015).

3.1.4 Ομαδοποίηση

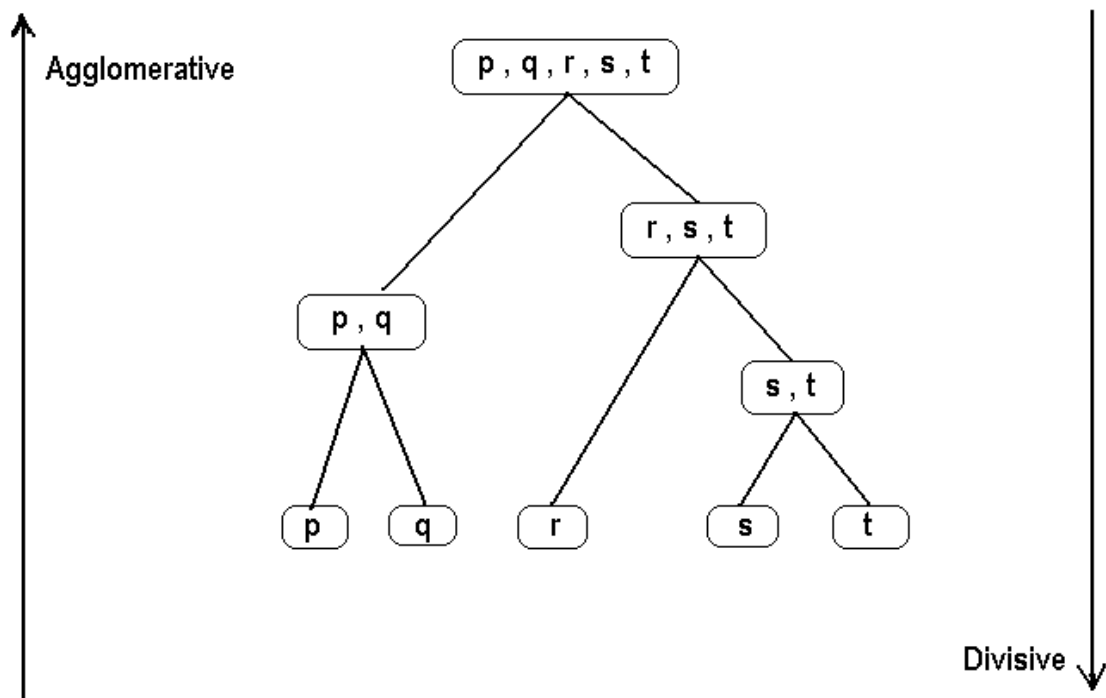
Η Ομαδοποίηση(clustering) είναι μια τεχνική ταξινόμησης δεδομένων η οποία περιγράφει την διαδικασία αναζήτησης ομάδων παρόμοιων δεδομένων μέσα στο σύνολο δεδομένων. Όστε να τα ταξινομήσει με κριτήριο κάποιες βασικές τους ομοιότητες σε ομάδες(clusters)π.χ. ηλικιακές ομάδες ατόμων. Με αυτή την μέθοδο επιτυγχάνετε η απλοποίηση της διαδικασίας προσπέλασης και ελέγχου των δεδομένων. Η ομαδοποίηση διαδραματίζει σημαντικό ρόλο και έχει εφαρμογή σε κλάδους όπως η έρευνα επιστημονικών δεδομένων, η ανάκτηση πληροφοριών στην ασφάλεια υπολογιστικών συστημάτων(computer security), το μάρκετινγκ, η Ιατρική, η βιολογία κ.α. Στην Εικόνα 9 φαίνεται το φιλτράρισμα και η ταξινόμηση του δείγματος σε τρεις ομάδες.



Εικόνα 9. Ομαδοποίηση (Clustering)("Oracle Advanced Analytics Data Mining Algorithms and Functions SQL API", n.d.).

Ορισμένες κατηγορίες τεχνικών ομαδοποίησης είναι οι εξής:

1. Η Ιεραρχική μέθοδος (Hierarchical method) ομαδοποίησης είναι μια μέθοδος ομαδοποίησης η οποία είτε διαιρώντας σύνολα δεδομένων είτε αθροίζοντας σύνολα δεδομένων φτάνουν στην τελική μορφή ομαδοποίησης. Υπάρχουν δύο κατηγορίες ιεραρχικών μεθόδων η συσσωρευτική μέθοδος(agglomerative method) και η διαιρετική μέθοδος(divisive method). Στην συσσωρευτική μέθοδος(agglomerative method)τα αντικείμενα ξεκινούν από μικρές ομάδες και μεμονωμένα αντικείμενα ώστε να συγκροτήσουν την μεγάλη ομάδα με βάση έναν πίνακα που καταγράφει αποστάσεις ομοιότητας(Εικόνα 10). Αντίθετα στην διαιρετική μέθοδος(divisive method) ξεκινάει όλα τα δεδομένα από μία ομάδα και τα υποδιαίρει μέχρι το κάθε αντικείμενο να βρίσκεται μόνο του με βάση έναν πίνακα που καταγράφει αποστάσεις ομοιότητας(Εικόνα10).



Εικόνα 10. Ιεραρχική μέθοδος ομαδοποίησης προς τα πάνω είναι συσσωρευτική και προς τα κάτω είναι διαιρετική("Hierarchical clustering", n.d).

2. Η μη Ιεραρχική μέθοδος(Non Hierarchical method) ομαδοποίησης διαίρει ένα σύνολο δεδομένων από N αντικείμενα σε M ομάδες(clusters), με επικάλυψη(overlap) ή χωρίς επικάλυψη. Όπου τα αντικείμενα κατατάσσονται στις ομάδες με κριτήριο την ομοιότητα τους μετά από διαδοχικές προσεγγίσεις. Ο αριθμός M των ομάδων(clusters) μπορούν να ορίζονται είτε κατά την διάρκεια της διαδικασίας ομαδοποίησης(clustering) είτε εκ των προτέρων, αυτή είναι μια βασική διάφορα με την Ιεραρχική μέθοδο ομαδοποίησης. Ουσιαστικά οι μη Ιεραρχικές μέθοδοι ομαδοποίησης χρησιμοποιούν μεθόδους Διαχωρισμού(partitioning methods)οι οποίες καταλήγουν σε ένα σύνολο ομάδων M . Οι οποίες περιέχουν αντικείμενα που ανήκουν σε μία μόνο ομάδα το κάθε ένα. Κάθε ομάδα μπορεί να εκπροσωπηθεί από ένα κεντρικό αντικείμενο που προσδιορίζει και τα υπόλοιπα ή ένα σύμπλεγμα αντικειμένων. Εάν είναι διαθέσιμες οι αριθμητικές τιμές των αντικειμένων, τότε ο μέσος όρος των τιμών των αντικειμένων είναι ο κατάλληλος εκπρόσωπος της ομάδας(Joshi, 2012).

3.1.4.1 Αλγόριθμοι μη Ιεραρχικής ομαδοποίησης

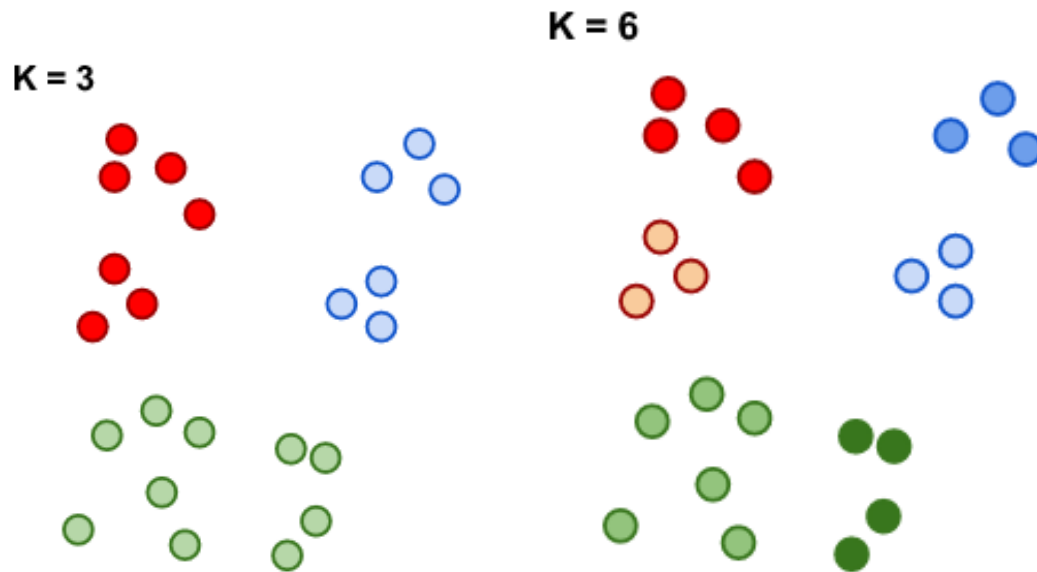
Ορισμένοι αλγόριθμοι μη Ιεραρχικής ομαδοποίησης(Non hierarchical Method) είναι οι εξής:

1. Ο αλγόριθμος K-means ομαδοποίησης(clustering) που αναπτύχθηκε από τον MacQueen το 1967 (MacQueen, 1967). Είναι μια μέθοδος ανάλυσης ομαδοποίησης που έχει ως στόχο να στεγανοποιήσει τις παρατηρήσεις σε k ομάδες.

Αποτελείτε από τα εξής 4 βήματα

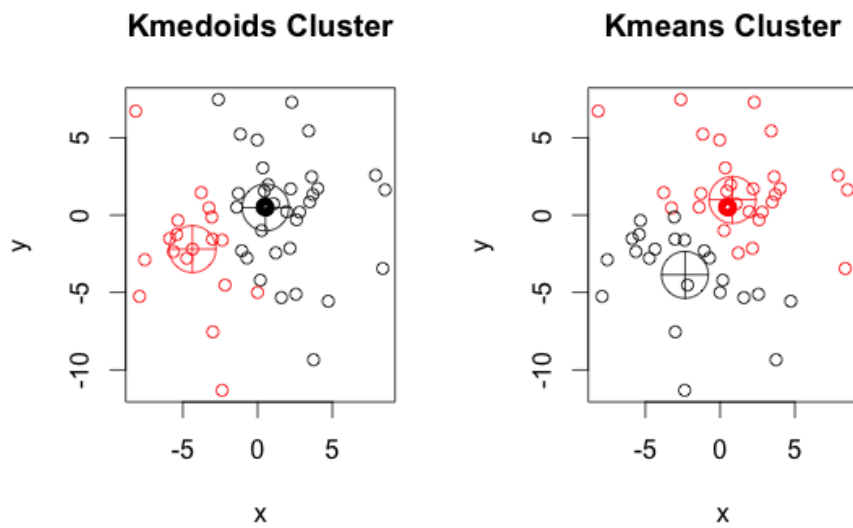
1. Χωρίζει τυχαία τα αντικείμενα σε k ομάδες.
2. Υπολογίζει την μέση τιμή(κέντρο βάρους) της κάθε ομάδας.
3. Ομαδοποιεί το κάθε αντικείμενο στην ομάδα με την πλησιέστερη μέση τιμή(Joshi, 2012).
4. Επαναλαμβάνει τα βήματα από το βήμα 2 μέχρι να μην υπάρχει καμία νέα τροποποίηση στην ομαδοποίηση.

Επειδή δημιουργεί αντικείμενα που ανήκουν σε μια και μόνο ομάδα και όχι σε περισσότερες, έχει πλεονέκτημα όταν τα δεδομένα είναι πάρα πολλά και η επεξεργασία τους με ιεραρχικές κυρίως μεθόδους ομαδοποίησης είναι είτε περίπλοκη είτε αδύνατη. Ένα όμως βασικό μειονέκτημα της μεθόδου K-means ομαδοποίησης είναι ο μεταβαλλόμενος τρόπος επιλογής του αριθμού των ομάδων(clusters). Ενδεικτικά στην εικόνα 11 παρουσιάζετε στο σύνολο δεδομένων η διάφορα της ομαδοποίησης με όταν ο αριθμός ομάδων είναι τρεις(K=3) και όταν ο αριθμός ομάδων είναι έξι(κ=6).



Εικόνα 11. K-means για $k=3$ και $k=6$ ("Non-hierarchical cluster analysis", n.d.).

2. Ο αλγόριθμος K-mediod ομαδοποίησης (clustering) είναι μια μέθοδος ανάλυσης ομαδοποίησης η οποία συγκεντρώνει το σύνολο δεδομένων των n αντικειμένων σε k ομάδες που είναι γνωστές εκ των προτέρων. Η συγκεκριμένη μέθοδος είναι πιο ισχυρή στο θόρυβο και στις ακραίες τιμές σε σύγκριση με την k -means (Joshi, 2012). Αφού η ομάδα αντιπροσωπεύεται από το πιο κεντρικό αντικείμενο (medoid) της ομάδας (cluster).



Εικόνα 12. Βλέπουμε την διαφορά της K-medoids με K-means ("difference between k means and k medoid", 2015).

Στον πίνακα 1 παρουσιάζονται συγκεντρωτικά σε κάθε στήλη οι τεχνικές ομαδοποίησης και ταξινόμησης:

Τεχνικές ομαδοποίησης(clustering)	Τεχνικές ταξινόμησης (classification)
Ιεραρχική μέθοδος(Hierarchical method)	Δέντρα Απόφασης (decision trees)
Μη Ιεραρχική μέθοδος (Non-Hierarchical method)	Μηχανές Υποστήριξης Διανυσμάτων (Support vector Machines, SVM)
	Ασαφής λογική (Fuzzy Logic)
	Naive Bayes
	Νευρωνικά δίκτυα (Neural Networks)

Πίνακας 1. Τεχνικές Ομαδοποίησης(clustering) και Τεχνικές Ταξινόμησης (classification).

Κεφάλαιο 4

Ανάλυση Συνόλου Δεδομένων

Το σύνολο δεδομένων KDD 1999 χρησιμοποιήθηκε στον 3ο Διεθνή Διαγωνισμό Εργαλείων Αναγνώρισης Γνώσης και Εξόρυξης Δεδομένων [(Third international Knowledge Discovery and Data Mining (KDD) Tools Competition)], ο οποίος διεξήχθη με την συνδρομή της 5ης Διεθνής Διάσκεψης για την Ανακάλυψη Γνώσης και Εξόρυξης Δεδομένων (The Fifth International Conference on Knowledge Discovery and Data Mining). Ο στόχος του διαγωνισμού ήταν να κατασκευαστεί ένας ανιχνευτής παρεισφρήσεων δικτύου που θα προστατεύει το δίκτυο υπολογιστών από μη εξουσιοδοτημένους χρήστες εντός και εκτός δικτύου. Ο ανιχνευτής παρεισφρήσεων δικτύου θα είχε ένα πρόγραμμα αξιολόγησης και ανίχνευσης παρεισφρήσεων το οποίο θα διέκρινε τις κακές συνδέσεις (εισβολές) από τις καλές συνδέσεις (κανονικές συνδέσεις). Ως σύνδεση που την αξιολογούμε αναλόγως είτε σε καλή σύνδεση (κανονική σύνδεση) είτε σε κακή σύνδεση (εισβολή) ορίζουμε μια ακολουθία πακέτων που αρχίζουν και τελειώνουν σε αυστηρά καθορισμένους χρόνους. Όπου τα δεδομένα που μεταφέρονται μπορούν να ταξιδεύουν αμφίδρομα και από τα δύο άκρα της σύνδεσης (διευθύνσεις IP τερματικών) με την χρήση πρωτοκόλλων μεταφοράς. Τέλος κάθε μεμονωμένη σύνδεση μεταφέρει περίπου 100 bytes (Stolfo, Fan, Lee, Prodromidis, & Chan, 2000). Ουσιαστικά το σύνολο δεδομένων KDD 1999 χρησιμοποιήθηκε για να ερευνηθεί και να αναπτυχθεί περαιτέρω η ανίχνευση εισβολών μέσω από την συνεχή αξιολόγηση Συστημάτων Ανίχνευσης Παρεισφρήσεων. Συγκεκριμένα το σύνολο δεδομένων KDD 1999 χρησιμοποιεί μια έκδοση του συνόλου δεδομένων DARPA 98. Το σύνολο δεδομένων DARPA 98 κατασκευαστικά και διαχειριστικά από την MIT Lincoln Labs υπό την αιγίδα των

Προηγμένων Αμυντικών Ερευνητικών Προγραμμάτων(Defense Advanced Research Projects) και το Εργαστήριο Έρευνας της Αεροπορίας(Air Force Research Laboratory). Το σύνολο δεδομένων DARPA 98 περιλάμβανε μια ποικιλία εισβολών που προσαρμόστηκαν σε ένα περιβάλλον στρατιωτικού δικτύου. Για να δημιουργηθεί το σύνολο δεδομένων χρειάστηκαν εννέα εβδομάδες συλλογής ακατέργαστων δεδομένων από ένα τοπικό δίκτυο που δεχόταν επιθέσεις, το οποίο δίκτυο το είχαν προσομοιώσει με ένα τυπικό τοπικό δίκτυο των αεροπορικών δυνάμεων των Η.Π.Α. Συνολικά συλλέχθηκαν 4 gigabytes ακατέργαστα δεδομένα(Stolfo, Fan, Lee, Prodromidis, & Chan, 2000).

Συγκεκριμένα οι τύποι εισβολών που υπάρχουν στο σύνολο δεδομένων KDD 1999 χωρίζονται σε τέσσερις κύριες κατηγορίες και αυτές αντίστοιχα σε 22 υποκατηγορίες όπως παρατίθενται και στον πίνακα 2(Stolfo, Fan, Lee, Prodromidis, & Chan, 2000).

Κατηγορίες εισβολών	Υποκατηγορίες εισβολών
Denial of Service attack	back, land, neptune, pod, smurf, teardrop
Remote to User attack (R2L)	ftp-write, guess-passwd, imap, multihop, phf, spy, warezclient, warezmaster
User to Root attack (U2R)	U2R buffer-overflow, perl, loadmodule, rootkit
Probing attack	ipsweep, nmap, portsweep, satan

Πίνακας 2. Διάφοροι τύποι επιθέσεων στο σύνολο δεδομένων KDD 1999(Stolfo, Fan, Lee, Prodromidis, & Chan, 2000).

4.1 Πρωτόκολλα Επικοινωνίας

Τα πρωτόκολλα επικοινωνίας που συμπεριλαμβάνονται στο σύνολο δεδομένων KDD 1999 είναι τα εξής:

1. Το Πρωτόκολλο Ελέγχου Μεταφοράς (Transmission Control Protocol, TCP) είναι ένα από τα σημαντικότερα πρωτόκολλα της σουίτας των πρωτοκόλλων Διαδικτύου (Internet Protocol) στο επίπεδο μεταφοράς. Το πρωτόκολλο Ελέγχου Μεταφοράς είναι το 4ο επίπεδο στο μοντέλο αναφοράς Ανοικτής Διασύνδεσης Συστημάτων (μοντέλο OSI). Το συγκεκριμένο πρωτόκολλο είναι υπεύθυνο για την αξιόπιστη μεταφορά των δεδομένων, αφού είναι προσανατολισμένης σύνδεσης (connection oriented). Που συνεπάγεται ότι δεδομένα αποστέλλονται από την μια πλευρά στην άλλη με την σωστή σειρά. Ο τρόπος που το επιτυγχάνει είναι χωρίζοντας τα δεδομένα σε πακέτα με ετικέτες και τα αποστέλλει μέσω του δικτύου. Το πρωτόκολλο Ελέγχου Μεταφοράς χρησιμοποιείται από πολλές υπηρεσίες όπως το Πρωτόκολλο Μεταφοράς Υπερκειμένου (Hypertext Transfer Protocol, HTTP) και το SMTP (Simple Mail Transfer Protocol) (Khubeb, Naahid, & Naahid, 2013).
2. Το Πρωτόκολλο User Datagram Protocol (UDP) αν και έχει παρομοία συμπεριφορά με το Πρωτόκολλο Ελέγχου Μεταφοράς διαφοροποιείται ως προς την αξιοπιστία του, αφού κάνει μεταφορά πακέτων δεδομένων χωρίς σύνδεση (connection-less). Δεδομένου ότι τα δεδομένα ταξιδεύουν σε αναξιόπιστα μέσα, ενδέχεται να μην φτάσουν στην σωστή σειρά, επίσης μπορεί να υπάρχει απώλεια πακέτων, καθώς και αλληλοεπικάλυψη (duplication) πακέτων. Το συγκεκριμένο πρωτόκολλο είναι χρήσιμο όταν μας ενδιαφέρει περισσότερο η παράδοση των πακέτων σε συγκεκριμένο χρόνο έναντι της απώλειας δεδομένων (Khubeb, Naahid, & Naahid, 2013).
3. Το Πρωτόκολλο Ελέγχου Μηνυμάτων Διαδικτύου (Internet Control Message Protocol, ICMP) χρησιμοποιείται κυρίως για την επικοινωνία δύο υπολογιστών που βρίσκονται σε σύνδεση. Ο κύριος σκοπός του Πρωτοκόλλου Ελέγχου Μηνυμάτων Διαδικτύου είναι να στέλνει μηνύματα μεταξύ υπολογιστών που βρίσκονται συνδεδεμένοι στο δίκτυο σχετικά με το δίκτυο. Συγκεκριμένα ανακατευθύνει τα μηνύματα ανάλογα με τις τροποποιήσεις του δικτύου και χρησιμοποιείται από τους δρομολογητές (routers) ώστε να ενημερώνει με πληροφορίες δρομολόγησης τα τερματικά (hosts) του δικτύου, τα οποία στα πρώτα

στάδια της σύνδεσης τους στο δίκτυο γνωρίζουν ελάχιστες πληροφορίες δρομολόγησης και στην συνέχεια χρειάζονται ενημέρωση για της τροποποίησης του δικτύου (Khubeb, Naahid, & Naahid, 2013).

Οι υποκατηγορίες εισβολών κατανέμονται ανάλογα με το κάθε πρωτόκολλο επικοινωνίας το οποίο προσβάλλουν όπως φαίνετε στον Πίνακα 3.

Τύπος Πρωτοκόλλου	Όνομα εισβολής
Πρωτόκολλο Ελέγχου Μεταφοράς [Transmission Control Protocol, TCP]	Teardrop, satan, nmap, rootkit
Πρωτόκολλο User Datagram Protocol (UDP)	Neptune, guess_passwd, land, portsweep, buffer_overflow, phf, warezmaster, ipsweep, multihop, wwarezclient, perl, back, ftp_write, loadmodule, satan, spy, imap, rootkit
Πρωτόκολλο Ελέγχου Μηνυμάτων Διαδικτύου(Internet Control Message Protocol, ICMP)	Portsweep, ipsweep, smurf, satan, pod, nmap

Πίνακας 3. Υποκατηγορίες εισβολών κατανεμημένες ανάλογα με το πρωτόκολλο το οποίο προσβάλλουν(Khubeb, Naahid, & Naahid, 2013).

4.2 Ανάλυση Χαρακτηριστικών

Το σύνολο δεδομένων KDD 1999 αποτελείτε από 4.900.000 μεμονωμένους φορείς συνδέσεων, καθένας από τους οποίους περιέχει 41 χαρακτηριστικά τα οποία χαρακτηρίζουν μια σύνδεση είτε ως καλή σύνδεση(κανονική σύνδεση) είτε ως κακή σύνδεση(εισβολή). Κάνοντας αντιστοίχιση με έναν τύπο επίθεσης στην περίπτωση που είναι κακή σύνδεση (εισβολή) (Campos, Oliveira, & Roisenberg, 2012). Τα 41 χαρακτηριστικά που υπάρχουν στο σύνολο δεδομένων KDD 1999, τα οποία χαρακτηρίζουν μια σύνδεση είτε ως καλή σύνδεση(κανονική σύνδεση) είτε ως κακή σύνδεση(εισβολή), κάνοντας αντιστοίχιση με έναν τύπο επίθεσης στην περίπτωση της εισβολής. Παρουσιάζονται στους παρακάτω πίνακες (Πίνακας 4, Πίνακας 5, Πίνακας 6) με

το όνομα τους ,την περιγραφή τους και τον τύπο τους(Stolfo, Fan, Lee, Prodromidis, & Chan, 2000).

Όνομα	Περιγραφή	Τύπος
Duration	Διάρκεια (δευτερολέπτων) της σύνδεσης.	Συνεχής
Protocol_type	Τύπος πρωτοκόλλου.	Διακριτή
service	Υπηρεσίες δικτύου προορισμού.	Διακριτή
src_bytes	Αριθμός δεδομένων bytes από την πηγή προς τον προορισμό.	Συνεχής
dst_bytes	Αριθμός δεδομένων bytes από τον προορισμό προς την πηγή.	Συνεχής
flag	Φυσιολογική ή εσφαλμένη κατάσταση της σύνδεσης.	Διακριτή
land	1 εάν η σύνδεση είναι από τον ίδιο host/port. Αλλιώς 0.	Διακριτή
wrong_fragment	Αριθμός των λάθος fragments.	Συνεχής
urgent	Αριθμός επειγουσών πακέτων.	Διακριτή

Πίνακας 4. Χαρακτηριστικά που σχετίζονται με συνδέσεις στο πρωτόκολλο TCP(Stolfo, Fan, Lee, Prodromidis, & Chan, 2000).

Όνομα	Περιγραφή	Τύπος
hot	Αριθμός των «καυτών» δεικτών.	Συνεχής
num_failed_logins	Αριθμός των αποτυχημένων προσπαθειών σύνδεσης(login).	Συνεχής
logged_in	1 ένα είναι επιτυχημένη προσπάθεια σύνδεσης. Αλλιώς 0.	Διακριτή
num_compromised	Αριθμός των συνθηκών που βρίσκονται σε κίνδυνο.	Συνεχής
root_shell	1 εάν λαμβάνετε το κέλυφος root. Αλλιώς 0.	Διακριτή
su_attempted	1 εάν γίνετε αποπήρα εντολής sudo (root) Αλλιώς 0	Διακριτή
num_root	0 αριθμός των προσβάσεων ως root	Συνεχής

	(διαχειριστής).	
num_file_creations	Τον αριθμό των λειτουργιών που δημιουργούν αρχείο.	Συνεχής
num_shells	Ο αριθμός των shell παρακινήσεων.	Συνεχής
num_access_files	Ο αριθμός των λειτουργιών σε αρχεία πρόσβασης.	Συνεχής
num_outbound_cmds	αριθμός των εξερχόμενων εντολών σε μια συνεδρία FTP.	Συνεχής
is_hot_login	1 εάν η σύνδεση ανήκει στην καυτή λίστα. Αλλιώς 0.	Διακριτή
is_guest_login	1 εάν η σύνδεση είναι τύπου "guest" login. Αλλιώς 0.	Διακριτή

Πίνακας 5. Χαρακτηριστικά που σχετίζονται με συνδέσεις που είναι γνωστό το domain(Stolfo, Fan, Lee, Prodromidis, & Chan, 2000).

Όνομα	Περιγραφή	Τύπος
count	αριθμό των συνδέσεων που γίνονται στον host όπως συμβαίνει στην τρέχουσα σύνδεση κατά τα τελευταία δύο δευτερόλεπτα.	Συνεχής
<i>Σημείωση: Τα ακόλουθα χαρακτηριστικά αναφέρονται σε συνδέσεις στον ίδιο host.</i>		
serror_rate	Το ποσοστό επί τοις % των συνδέσεων που έχουν σφάλματα «SYN».	Συνεχής
rerror_rate	Το ποσοστό επί τοις % των συνδέσεων που έχουν σφάλματα «REJ».	Συνεχής
same_srv_rate	Το ποσοστό επί τοις % των συνδέσεων που προέρχονται από τις ίδιες λειτουργίες (service).	Συνεχής
diff_srv_rate	Το ποσοστό επί τοις % των συνδέσεων που προέρχονται από διαφορετικές λειτουργίες (service).	Συνεχής
srv_count	Τον αριθμό των συνδέσεων που προέρχονται από την ίδια λειτουργία (service) με την τρέχουσα	Συνεχής

	σύνδεση κατά τα τελευταία 2 δευτερόλεπτα.	
<i>Σημείωση: Τα ακόλουθα χαρακτηριστικά αναφέρονται σε συνδέσεις με ίδια λειτουργία.</i>		
srv_serror_rate	Το ποσοστό επί τοις % των συνδέσεων που έχουν σφάλματα «SYN».	Συνεχής
srv_rerror_rate	Το ποσοστό επί τοις % των συνδέσεων που έχουν σφάλματα «REJ».	Συνεχής
srv_diff_host_rate	Το ποσοστό επί τοις % των συνδέσεων που γίνονται σε διαφορετικούς host.	Συνεχής
dst_host_count	Μετράει τις συνδέσεις που έχουν τον ίδιο host προορισμού.	Συνεχής
dst_host_srv_count	Μετράει τις συνδέσεις που έχουν τον ίδιο host προορισμού και προέρχονται από τις ίδιες λειτουργίες (service).	Συνεχής
dst_host_same_srv_rate	Το ποσοστό επί τοις % των συνδέσεων που έχουν τον ίδιο host προορισμού και προέρχονται από τις ίδιες λειτουργίες (service).	Συνεχής
dst_host_diff_srv_rate	Το ποσοστό επί τοις % των συνδέσεων που έχουν τον ίδιο host προορισμού και προέρχονται από διαφορετικές λειτουργίες (service).	Συνεχής
dst_host_same_src_port_rate	Το ποσοστό επί τοις % των συνδέσεων που έχουν τον ίδιο host προορισμού και την ίδια src port.	Συνεχής
dst_host_srv_diff_host_rate	Το ποσοστό επί τοις % των συνδέσεων που προέρχονται από τις ίδιες λειτουργίες (service) άλλα από διαφορετικό host.	Συνεχής
dst_host_serror_rate	Το ποσοστό επί τοις % των συνδέσεων που προέρχονται από τον ίδιο host και έχουν το S0 σφάλμα.	Συνεχής
dst_host_srv_serror_rate	Το ποσοστό επί τοις % των συνδέσεων που προέρχονται από τον ίδιο host και καθορίζουν την	Συνεχής

	λειτουργία που έχει το So σφάλμα.	
dst_host_rerror_rate	Το ποσοστό επί τοις % των συνδέσεων που προέρχονται από τον ίδιο host και έχουν το RST σφάλμα.	Συνεχής
dst_host_srv_rerror_rate	Το ποσοστό επί τοις % των συνδέσεων που προέρχονται από τον ίδιο host και καθορίζουν την λειτουργία που έχει το RST σφάλμα.	Συνεχής

Πίνακας 6. Χαρακτηρίστηκα κίνησης στο δίκτυο που καταγράφονται χρησιμοποιώντας χρονικό παράθυρο δύο δευτερόλεπτων (Stolfo, Fan, Lee, Prodromidis, & Chan, 2000).

Κεφάλαιο 5

Υλοποίηση Εφαρμογής

5.1 Εισαγωγή στο Rstudio

Για την υλοποίηση του πρακτικού μέρους της μεταπτυχιακής διατριβής χρησιμοποιήθηκε το R Studio, το οποίο είναι ένα ολοκληρωμένο περιβάλλον ανάπτυξης (Integrated Development Environment, IDE) κώδικα R, στο οποίο επεξεργάστηκαν και αναλύθηκαν τα δεδομένα. Το Rstudio και οι παραγόμενες ρουτίνες του χρησιμοποιούν τη γλώσσα προγραμματισμού C++ σε επίπεδο συστήματος προκειμένου να εκτελέσουν ταχύτατα τους αλγόριθμους τοπικά ή απομακρυσμένα σε εγκατάσταση εξυπηρέτη (server). Είναι ιδανικό για στατιστικούς υπολογισμούς και γραφικές αναπαραστάσεις, αφού παρέχει τη δυνατότητα χρήσης μιας πληθώρας στατιστικών εργαλείων και τεχνικών όπως είναι: η γραμμική μοντελοποίηση (linear modeling), η μη γραμμική μοντελοποίηση (non-linear modeling), η κατηγοριοποίηση (classification), η ομαδοποίηση (clustering) κ.α. Τέλος είναι λογισμικό ανοιχτού κώδικα με δυνατότητα παραμετροποίησης και επέκτασης.

5.2 Βήματα Ανάλυσης Συνόλου Δεδομένων

Γνωρίζουμε ήδη ότι το σύνολο δεδομένων KDD 1999 έχει έναν μεγάλο αριθμό εγγραφών όπως αναλύθηκε και στο κεφάλαιο 4. Το μεγάλο πλήθος εγγραφών προς επεξεργασία αυξάνει ανάλογα και το υπολογιστικό κόστος. Στην προσπάθεια μείωσης του υπολογιστικού κόστους δημιουργήθηκε το σύνολο δεδομένων KDD 10% το οποίο είναι ένα υποσύνολο που περιέχει μόνο το 10% των δεδομένων εκπαίδευσης, τα οποία ελήφθησαν τυχαία από το αρχικό σύνολο δεδομένων. Εκτός όμως από το αρχικό σύνολο δεδομένων KDD 1999 και

το σύνολο δεδομένων KDD 10% υπάρχει και ένα τρίτο σύνολο δεδομένων που αποκαλείτε διορθωμένο KDD(corrected KDD). Αυτό δεν έχει την ίδια κατανομή πιθανότητας επιθέσεων με τα άλλα δύο σύνολα δεδομένων αφού περιέχει 14 νέα είδη εισβολών (Al-mamory, & Jassim, 2013).

Στα πλαίσια της διάκρισης των «κακών» και καλών συνδέσεων θα πραγματοποιήσω ανάλυση δεδομένων χωρίς υψηλό υπολογιστικό κόστος, οπότε θα χρησιμοποιήσω το σύνολο δεδομένων 10%, το οποίο θα διαχειριστώ με τη χρήση του λογισμικού ανάλυσης δεδομένων Rstudio ώστε να εξάγω αποτελέσματα.

5.2.1 Φόρτωση και Προεπεξεργασία Δεδομένων

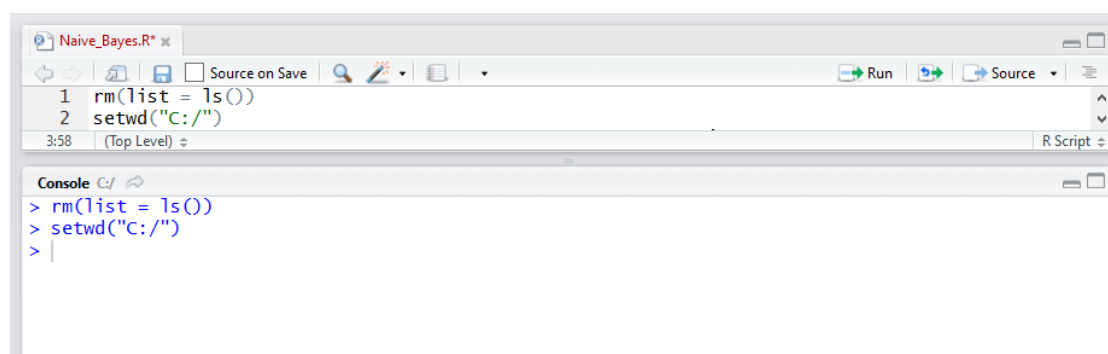
Στο παρακάτω κείμενο το σενάριο εντολών που εκτελέστηκαν εμφανίζεται με **έντονη** γραμματοσειρά Cambria.

Αρχικά καθάρισα τον περιβάλλον εργασίας με την εντολή:

```
rm( list = ls() )
```

Στη συνέχεια αρχικοποίησα το περιβάλλον εργασίας:

```
setwd("C:/")
```



The screenshot shows the RStudio interface. The top pane displays the R script with two lines of code: `1 rm(list = ls())` and `2 setwd("C:/")`. The bottom pane shows the console output, which is empty, indicating that the commands have been executed successfully. The console prompt is `>`.

Εικόνα 13. Στιγμιότυπο αρχικοποίησης περιβάλλοντος εργασίας.

Στην συνέχεια φορτώνεται το 10% dataset του KDD 1999 και ακολουθεί προεπεξεργασία των δεδομένων στην κονσόλα του Rstudio την εξής διαδικασία:

Αρχικά φορτώνεται το 10% σύνολο δεδομένων του KDD με την εντολή `read.csv` και την παράμετρο `stringsAsFactors = FALSE` η οποία απενεργοποιεί την αυτόματη μετατροπή των συμβολοσειρών χαρακτήρων(`strings`) σε παράγοντες. Φορτώνεται το αρχείο `names` στο οποίο αναγράφονται τα ονόματα των χαρακτηριστικών του συνόλου δεδομένων `kdd 1999` προσθέτοντας το χαρακτηριστικό `label`. Αναλυτικά τα βήματα:

```
train_raw <- read.csv ("kddcup.data_10_percent_corrected.csv" ,
stringsAsFactors = FALSE)
```

```
colnames <- read.table("names", skip = 1, sep = - ":")
```

```
names(train_raw) <- colnames$V1
```

```
d <- dim(train_raw)
```

προσθέτω στο αρχείο την ετικέτα `label` στο τέλος:

```
names(train_raw)[d[2]] <- "label"
```

Εμφανίζω τα ονόματα των χαρακτηριστικών :

```
names(train_raw)
```

Παρακάτω το στιγμιότυπο (Εικόνα 14):

```
[1] "duration"           "protocol_type"      "service"
[4] "flag"               "src_bytes"          "dst_bytes"
[7] "land"               "wrong_fragment"    "urgent"
[10] "hot"                "num_failed_logins" "logged_in"
[13] "num_compromised"    "root_shell"        "su_attempted"
[16] "num_root"           "num_file_creations" "num_shells"
[19] "num_access_files"   "num_outbound_cmds" "is_host_login"
[22] "is_guest_login"     "count"              "srv_count"
[25] "serror_rate"        "srv_error_rate"     "error_rate"
[28] "srv_rerror_rate"    "same_srv_rate"     "diff_srv_rate"
[31] "srv_diff_host_rate" "dst_host_count"     "dst_host_srv_count"
[34] "dst_host_same_srv_rate" "dst_host_diff_srv_rate" "dst_host_same_src_port_rate"
[37] "dst_host_srv_diff_host_rate" "dst_host_serror_rate" "dst_host_srv_serror_rate"
[40] "dst_host_rerror_rate" "dst_host_srv_rerror_rate" "label"
> |
```

Εικόνα 14.

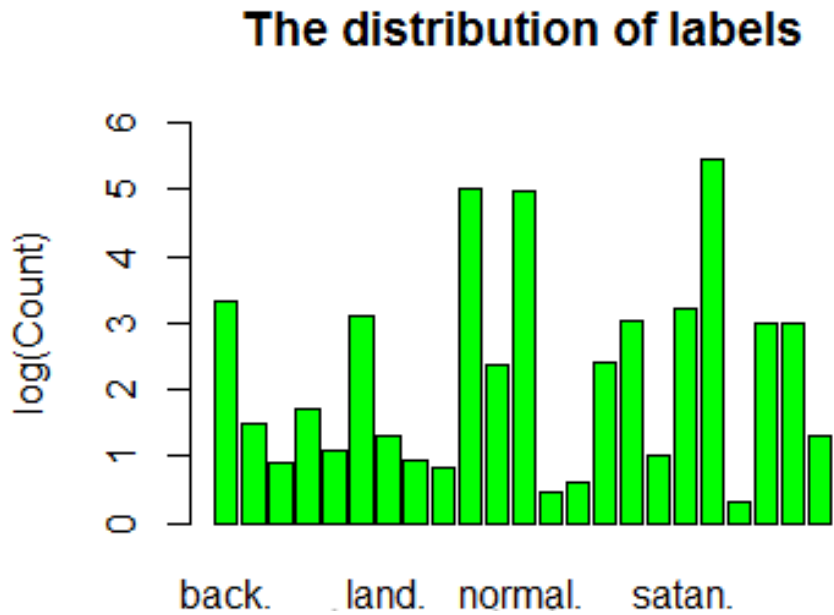
Παρατηρώ ότι στο 10% KDD σύνολο δεδομένων αποτελείτε από 494020 εγγραφές (Εικόνα 15) που έχουν 41 χαρακτηριστικά(Εικόνα 14) και το τελευταίο το 42^ο είναι η ετικέτα(label).

Data	
colnames	40 obs. of 1 variable
train_raw	494020 obs. of 42 variables
Values	
d	int [1:2] 494020 42

Εικόνα 15.

Στη συνέχεια παρουσιάζονται οι κατανομές των ετικετών σε γράφημα με τις υποκατηγορίες των επιθέσεων(Εικόνα 16).

```
sum_label <- aggregate(rep(1, d[1]),  
                        by = list(train_raw$label), FUN = sum)  
names(sum_label) <- c("label", "count")  
barplot(beside = TRUE, log10(sum_label$count),  
        Ορίζω τον άξονα y.  
        names.arg = sum_label$label, ylim = c(0,6),  
        xlab = "Label", ylab = "log(Count)",  
        col = "Blue", main = "The distribution of labels")
```



Εικόνα 16.

Εγκαθίστανται και φορτώνονται τα απαραίτητα πακέτα (βιβλιοθήκες)
 Φορτώνεται η βιβλιοθήκη caret η οποία θα χρησιμοποιήσει το πακέτο lattice
 και το πακέτο ggplot2.

library(caret)

Έπειτα, στην προσπάθεια μας να αφαιρέσουμε δεδομένα τα οποία δεν προσθέτουν σημαντική πληροφορία για την κατηγοριοποίηση των συνδέσεων επιλέγουμε και απομακρύνουμε τις εγγραφές χωρίς τιμή (NA) και τις εγγραφές με μηδενική διακύμανση

```
l <- train_raw$label
```

```
sum(is.na(l))
```

Το αποτέλεσμα είναι [1] 0

```
>
> # select the features
> library(caret)
> l <- train_raw$label
> sum(is.na(l))
[1] 0
> |
```


Εικόνα 17.

Στην εικόνα 17 φαίνεται ότι δεν υπάρχουν χαρακτηριστικά με κενές τιμές (με τιμή NA). Στην συνέχεια απομακρύνω τις τιμές που παρουσιάζουν μηδενική διακύμανση:

```
nzvc <- nearZeroVar(train_raw)
train_raw <- train_raw[, -nzvc]
```

Μετατρέπεται η ετικέτα σε χαρακτηριστικό τύπου «παράγοντας»

```
training <- train_raw
training$label <- factor(training$label)
d <- dim(training)
```

Αφού ολοκληρώθηκε η προεπεξεργασία των δεδομένων του συνόλου δεδομένων 10% KDD το παραγόμενο σύνολο δεδομένων μετά την προεπεξεργασία αποτελείται από 494020 παρατηρήσεις και 19 χαρακτηριστικά όπως παρουσιάζετε και στην Εικόνα 18, ενδεικτικά βλέπουμε ένα μέρος των χαρακτηριστικών και των παρατηρήσεων στην Εικόνα 19.

train_raw	494020 obs. of 19 variables
training	494020 obs. of 19 variables

Εικόνα 18.

	service	src_bytes	dst_bytes	flag	num_compromised	serror_rate	r
209	tcp	finger	SF	9	0	1	
213	udp	domain_u	SF	33	0	1	
214	udp	domain_u	SF	30	0	1	
217	udp	domain_u	SF	30	0	1	

< Showing 1 to 5 of 494,020 entries

Εικόνα 19.

5.2.2 Κατηγοριοποίηση με μοντέλο Naïve Bayes

Στο 10% KDD σύνολο δεδομένων θα εφαρμοστεί η τεχνική κατηγοριοποίησης Naïve Bayes.

Φορτώνεται η βιβλιοθήκη e1071:

```
library(e1071)
```

Πραγματοποιείται διαχωρισμός των δεδομένων (training set, test set) και δημιουργία του μοντέλου Naïve Bayes:

```
label_result = training[,d[2]]
```

```
training_data = training[,1:(d[2]-1)]
```

```
navie_bayes_tree_model = naiveBayes(as.factor(label_result)~.,  
                                   training_data)
```

Ακολουθεί πρόβλεψη στο σύνολο testing_data:

```
testing_data = testing[, 1: (d[2]-1)]
```

```
navie_bayes_pred = predict(navie_bayes_tree_model, testing_data)
```

```
golden_answer = testing[, d[2]]
```

```
navie_bayes_pred = factor(navie_bayes_pred, levels  
=levels(golden_answer))
```

Εξάγεται η ακρίβεια πρόβλεψης:

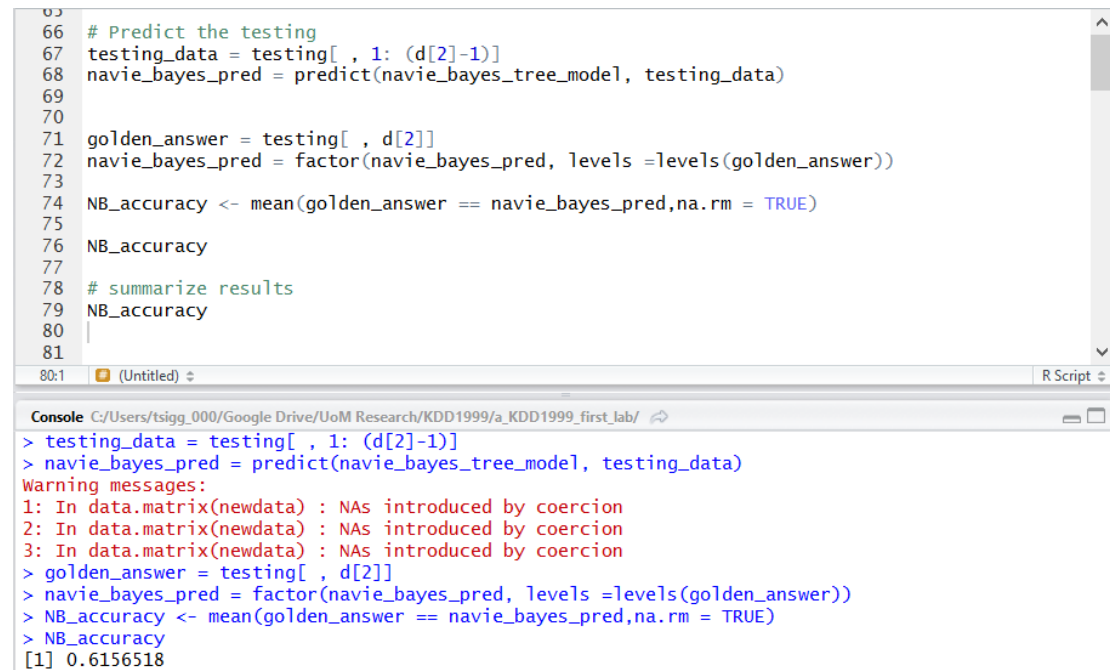
```
golden_answer = testing[, d[2]]
```

```
navie_bayes_pred = factor(navie_bayes_pred, levels
=levels(golden_answer))
```

```
NB_accuracy <- mean(golden_answer == navie_bayes_pred,na.rm = TRUE)
```

```
NB_accuracy
```

Η πρόβλεψη ακρίβειας του αλγόριθμου Naïve Bayes είναι 0,6156518(Εικόνα 20).



```
66 # Predict the testing
67 testing_data = testing[, 1: (d[2]-1)]
68 navie_bayes_pred = predict(navie_bayes_tree_model, testing_data)
69
70
71 golden_answer = testing[, d[2]]
72 navie_bayes_pred = factor(navie_bayes_pred, levels =levels(golden_answer))
73
74 NB_accuracy <- mean(golden_answer == navie_bayes_pred,na.rm = TRUE)
75
76 NB_accuracy
77
78 # summarize results
79 NB_accuracy
80
81
```

```
> testing_data = testing[, 1: (d[2]-1)]
> navie_bayes_pred = predict(navie_bayes_tree_model, testing_data)
Warning messages:
1: In data.matrix(newdata) : NAs introduced by coercion
2: In data.matrix(newdata) : NAs introduced by coercion
3: In data.matrix(newdata) : NAs introduced by coercion
> golden_answer = testing[, d[2]]
> navie_bayes_pred = factor(navie_bayes_pred, levels =levels(golden_answer))
> NB_accuracy <- mean(golden_answer == navie_bayes_pred,na.rm = TRUE)
> NB_accuracy
[1] 0.6156518
```

Εικόνα 20.

5.2.3 Μοντέλο Κατηγοριοποίησης Δέντρων Απόφασης

Χτίζω ένα νέο μοντέλο κατηγοριοποίησης χρησιμοποιώντας τεχνική κατηγοριοποίησης δέντρων απόφασης.

Εγκαθιστώ και φορτώνω απαραίτητες βιβλιοθήκες:

```
library(caret)
```

```
library(e1071)
```

```
library(randomForest)
```

```
library(doParallel)
```

Θέτω ένα τυχαίο αριθμό το λεγόμενο random seed:

```
set.seed(100)
```

Χρησιμοποιώ την συνάρτηση DoParallel:

```
registerDoParallel(makeCluster(detectCores()))
```

```
model_rf <- train(label ~ ., method = "rf", data = training)
```

Φορτώνω την βιβλιοθήκη `rpart`

```
library(rpart)
```

Κάνω τεχνική ταξινόμηση του Decision Tree:

```
decision_tree_model <- rpart(label ~ ., data = training, method = "class")
```

Πρόβλεψη του Δέντρου απόφασης:

```
decision_tree_pred <- predict(decision_tree_model, testing_data, type =  
"class")
```

Η παρακάτω εντολή θα μου εμφανίσει το Δέντρο απόφασης:

```
rpart.plot(decision_tree_model, main = "Classification Tree",  
  
extra = 102, under = TRUE, faclen = 0)
```

Με την παρακάτω εντολή ελέγχω τα αποτελέσματα των υπό επεξεργασία δεδομένων του συνόλου δεδομένων:

```
confusionMatrix(prediction1, subTesting$classe)
```

Παρατηρώ ότι τρέχει ατέρμονα χωρίς σταματημό άρα το μοντέλο των δέντρων απόφασης είναι αναποτελεσματικό στην προκειμένη περίπτωση. Αφού σε υπολογιστικά συστήματα 32bit υπερφορτώνετε η χρήση μνήμης και καταρρέει το υπολογιστικό σύστημα.

5.2.4 Ανάλυση Παραγόμενου Συνόλου Δεδομένων

Μετά την χρήση τεχνικών Εξόρυξης Δεδομένων θα κάνω ανάλυση του παραγόμενου συνόλου δεδομένων:

Με τις παρακάτω εντολές εμφανίζω τις περιττές εγγραφές στο σύνολο εξάσκησης αφού έχω εκτελέσει την προεπεξεργασία των δεδομένων(ενότητα 5.2.1):

```
d_train <- dim(train_raw)
```

```
d_unique_train <- dim(unique(train_raw))
```

```
d_train_percent <- (d_train[1] - d_unique_train[1]) / d_train[1]
```

Όπως παρουσιάζετε και στην εικόνα 21 στο 10% σύνολο δεδομένων εξάσκησης μετά την αφαίρεση των κενών εγγραφών υπάρχουν 494020 παρατηρήσεις. Όμως το νούμερο των μοναδικών παρατηρήσεων είναι 110192, ενώ επίσης έχουμε 77,6% (0.7769483) περιττές εγγραφές.

values	
d_train	int [1:2] 494020 19
d_train_perc...	0.776948301688191
d_unique_train	int [1:2] 110192 19

Εικόνα 21.

Με τις παρακάτω εντολές εμφανίζω τις περιττές εγγραφές στο σύνολο εξάσκησης μετά την εφαρμογή του αλγόριθμου Naïve Bayes(ενότητα 5.2.2):

```
d_test <- dim(test_raw)
```

```
d_unique_test <- dim(unique(test_raw))
```

```
d_test_percent <- (d_test[1] - d_unique_test[1]) / d_test[1]
```

Όπως παρουσιάζετε και στην εικόνα 22 στο 10% σύνολο δεδομένων εξάσκησης μετά την εφαρμογή του αλγόριθμου Naïve Bayes υπάρχουν 494020 παρατηρήσεις. Όμως το νούμερο των μοναδικών παρατηρήσεων είναι 145585. Ενώ επίσης έχουμε 70,5%(0.7053054532) περιττές παρατηρήσεις.

values	
d_test	int [1:2] 494020 42
d_test_perce...	0.705305453220517
d_train	int [1:2] 494020 19
d_train_perc...	0.776948301688191
d_unique_test	int [1:2] 145585 42
d_unique_train	int [1:2] 110192 19

Εικόνα 22.

Κεφάλαιο 6

Επίλογος

Αναλύοντας τα αποτελέσματα που προκύπτουν από την επεξεργασία του 10% KDD 1999 συνόλου δεδομένων και διαβάζοντας μελέτες πάνω στο συγκεκριμένο σύνολό δεδομένων δημιουργούνται δύο ζητήματα:

1. Υπάρχει ένα πλήθος περιττών εγγραφών όπως το 78% του συνόλου δεδομένων. Το συγκεκριμένο πλήθος δημιουργεί μια προκατάληψη στους αλγορίθμους να στρέφουν την προσοχή τους ως προς τις πιο συχνά εμφανιζόμενες εγγραφές, οπότε στις εκτιμήσεις τους μεροληπτούν υπέρ των εγγραφών που περιγράφουν πιο συχνά επιθέσεις όπως είναι η άρνηση εκτέλεσης εφαρμογής (Denial of Service attack). Αφού λοιπόν δεν εστιάζουν στις πιο σπάνια εμφανιζόμενες εγγραφές οι οποίες είναι εγγραφές που περιγράφουν συνήθως επιθέσεις όπως η επίθεση απομακρυσμένου χρήστη (Remote to User attack, R2L) και η επίθεση χρήστη σε διαχειριστή (User to Root attack, U2R), δημιουργούν μια ξεκάθαρη ευπάθεια στα Συστήματα Ανίχνευσης Παρειασφρήσεων να τις εντοπίσουν. Για να επιλυθεί το συγκεκριμένο πρόβλημα πρέπει να αφαιρέσουμε όλες τις περιττές εγγραφές και να διατηρήσουμε μόνο ένα μοναδικό αντίγραφο από τις περιττές εγγραφές (Tavallaee, Baghe, Lu, & A. Ghorbani, 2009).
2. Η πρόβλεψη ακρίβειας του μοντέλου κατηγοριοποίησης Naïve Bayes κυμαίνεται σε φυσιολογικά επίπεδα. Γεγονός που μας φανερώνει ότι λειτουργεί ικανοποιητικά και υπολογίζει είτε ελάχιστα είτε καθόλου ψευδείς συναγερούς είτε θετικούς (false positives) είτε αρνητικούς (false negatives).

Βιβλιογραφία

Al-mamory, S., & Jassim, F. (2013). Evaluation of Different Data Mining Algorithms with KDD CUP 99 Data Set. *Journal Of Babylon University/Pure And Applied Sciences*, 4.

Alpaydın, E. (2013). Introduction to Machine Learning 2nd ed. *Massachusetts Institute of Technology*, 2-3.

Brox, A. (2002). Signature-Based or Anomaly-Based Intrusion Detection: The Practice and Pitfalls. *scmagazine*. <https://www.scmagazine.com/signature-based-or-anomaly-based-intrusion-detection-the-practice-and-pitfalls/article/548733/> [Πρόσβαση: 26 Φεβρουαρίου 2017]

Calbimonte, D. (2014). Data Mining Introduction Part 9: Microsoft Linear Regression. <http://www.sqlservercentral.com>.
<http://www.sqlservercentral.com/articles/Data+Mining/110116/> [Πρόσβαση: 18 Δεκεμβρίου 2016]

Campos, L., Oliveira, R., & Roisenberg, M. (2012). Network Intrusion Detection System Using Data Mining. <Http://www.Campuscastanh.Ufpa.Br/>, 106.
https://www.researchgate.net/profile/Lidio_Campos/publication/232707594_Network_Intrusion_Detection_System_Using_Data_Mining/links/0deec521235cf10d91000000/Network-Intrusion-Detection-System-Using-Data-Mining.pdf
[Πρόσβαση: 07 Δεκεμβρίου 2016]

Chauhan, A., Mishra, G., & Kumar, G. (2011). Survey on Data Mining Techniques in Intrusion Detection. *International Journal Of Scientific & Engineering Research*, 1-3.

- Chen, T. Intrusion Detection for Viruses and Worms.
<https://pdfs.semanticscholar.org/7bdf/8ccb2b34ad93eab0bdea9b4af45c04eb9b4b.pdf>
[Πρόσβαση: 04 Δεκεμβρίου 2016]
- Das, N., & Sarkar, T. (2014). Survey on Host and Network Based Intrusion Detection System. *Int. J. Advanced Networking And Applications*, 2266-2269.
- Gu, Q., & Liu, P. (2007). Denial of Service Attacks, 4-7.
<https://s2.ist.psu.edu/paper/ddos-chap-gu-june-07.pdf> [Πρόσβαση: 12 Απριλίου 2017]
- Joshi, M. (2012). Classification, Clustering And Intrusion Detection. *System. International Journal Of Engineering Research And Applications (IJERA)*, 2.
- Khubeb, M., Naahid, S., & Naahid, S. (2013). Analysis of KDD CUP 99 Dataset using Clustering based Data Mining. *International Journal Of Database Theory And Application*, 24-28.
- MacQueen, J. (1967). Some Methods for classification and Analysis of Multivariate Observations. . *University Of California Press*, 281-297.
- McGreevy, J. (2002). Footprinting: What Is It, Who Should Do It, and Why?. *SANS Institute 2002*, 3-4.
- Mitnick, K., Simon, W., & Wozniak, S. (2002). The Art Of Deception. *Steave Wozniak/John wiley & sons*, 2.
- Nascimento, G., & Correia, M. Anomaly-based Intrusion Detection in Software as a Service.
<http://www.covert.io/researchapers/security/Anomalybased%20intrusion%20detection%20in%20software%20as%20a%20service.pdf> [Πρόσβαση: 15 Δεκεμβρίου 2016]
- Nielsen, M. (2015). *neural networks and deep learning*.
<http://neuralnetworksanddeeplearning.com/chap1.html> [Πρόσβαση: 04 Ιανουαρίου 2017]

Olusola., A., S.Oladele., A., & Abosede, D. (2010). Analysis of KDD '99 Intrusion Detection Dataset for Selection of Relevance Features. *Proceedings Of The World Congress On Engineering And Computer Science, Vol I WCECS 2010*, 7.

Paliwal, S., & Gupta, R. (2012). Denial-of-Service, Probing & Remote to User (R2L) Attack Detection using Genetic Algorithm. *International Journal Of Computer Applications*, 2.

Rechtin, E., & Maier, M. (2010). *The Art of Systems Architecting*, (2nd ed). CRC Press, 254.

Revathi, S., & Malathi, A. (2014). Detecting User-To-Root (U2R) Attacks Based on Various Machine Learning Techniques. *International Journal Of Advanced Research In Computer And Communication Engineering*, 2.

Sanghvi, H., & Dahiya, M. (2013). Cyber Reconnaissance: An Alarm before Cyber Attack. *International Journal Of Computer Application*.

Shib, J., & Saleem, S. (2012). *Computer Security Research Reports*.

Stolfo, S., Fan, W., Lee, W., Prodromidis, A., & Chan, P. (2000). Cost-based Modeling for Fraud and Intrusion Detection Results from *the JAM Project*. <https://kdd.ics.uci.edu/databases/kddcup99/task.html> [Πρόσβαση: 19 Ιανουαρίου 2017]

Tavallaee, M., Baghe, E., Lu, W., & A. Ghorbani, A. (2009). A Detailed Analysis of the KDD CUP 99 Data Set. *Proceedings Of The 2009 IEEE Symposium On Computational Intelligence In Security And Defense Applications*, 3-4. <http://www.ee.ryerson.ca/%7Ebagheri/papers/cisda.pdf> [Πρόσβαση: 12 Απριλίου 2017]

Uack, J. (2013). Data Mining. www.the-data-mine.com. <http://www.the-data-mine.com/Misc/DataMining> [Πρόσβαση: 19 Φεβρουαρίου 2017]

difference between k means and k medoid. (2015). *Stats.stackexchange.com*.
<https://stats.stackexchange.com/questions/156210/difference-between-k-means-and-k-medoid> [Πρόσβαση: 02 Απριλίου 2017]

Hierarchical clustering. *http://www.solver.com*.
<http://www.solver.com/xlminer/help/hierarchical-clustering-intro> [Πρόσβαση: 23 Δεκεμβρίου 2016]

Host Based IDS vs Network Based IDS | securitywing. (2012). *Securitywing.com*.
<http://securitywing.com/host-based-ids-vs-network-based-ids> [Πρόσβαση: 20 Δεκεμβρίου 2016]

Host- vs. Network-Based Intrusion Detection Systems. (2000). *SANS Institute*, 3-8.

IDS Introduction :: Chapter 16. Intrusion-Detection System :: Part VII: Detecting and Preventing Attacks :: Router firewall security :: Networking :: *eTutorials.org*.
Etutorials.org.
<http://etutorials.org/Networking/Router+firewall+security/Part+VII+Detecting+and+Preventing+Attacks/Chapter+16.+IntrusionDetection+System/IDS+Introduction/>[Πρόσβαση: 04 Απριλίου 2017]

Intrusion detection and prevention system. *www.slideshare.net*.
<https://image.slidesharecdn.com/intrusiondetectionandpreventionsystem-130518035126-phpapp02/95/intrusion-detection-and-prevention-system-18-638.jpg?cb=1369196205> [Πρόσβαση: 02 Ιανουαρίου 2017]

Non-hierarchical cluster analysis. *Mb3is.megx.net*.
<https://mb3is.megx.net/gustame/dissimilarity-based-methods/cluster-analysis/non-hierarchical-cluster-analysis>[Πρόσβαση: 24 Μαρτίου 2017]

Oracle Advanced Analytics Data Mining Algorithms and Functions SQL API.
http://www.oracle.com. <http://www.oracle.com/technetwork/database/enterprise-edition/odm-techniques-algorithms-097163.html> [Πρόσβαση: 30 Ιανουαρίου 2017]