

Εξόρυξη Δεδομένων και Αλγόριθμοι Μάθησης

Project

2024-2025

Μαύρα Πολυδώρου ΑΜ: 1064885

Παραδοχές:

- ➡ Για την συγγραφή του κώδικα χρησιμοποιήθηκε το περιβάλλον jupyter notebook που προσφέρεται ως extension του vscode για αποφυγή πολλών αρχείων κώδικα αφού έτσι μπορούσε να γραφεί συνολικά και το κάθε ερώτημα να εκτελείται ανεξάρτητα.
- ➡ Χρησιμοποιήθηκαν οι βιβλιοθήκες Matplotlib (pyplot)-seaborn, Pandas, NumPy για οπτικοποίηση των δεδομένων, επεξεργασία και ανάλυση αυτών και για αριθμητικους υπολογισμούς αντίστοιχα.
- ➡ Χρησιμοποιήθηκαν επίσης οι PCA (Principal Component Analysis) για την μείωση διαστάσεων στα γραφήματα, Train-Test Split για τον διαχωρισμό των συνόλων εκπαίδευσης/επικύρωσης στο τελευταίο ερώτημα, StandardScaler για κανονικοποίηση των δεδομένων στο ίδιο ερώτημα.
- ➡ Τέλος χρησιμοποιήθηκαν επίσης DBSCAN, Birch, KMeans για τους αντίστοιχους αλγόριθμους που χρησιμοποιήθηκαν στα ερωτήματα.

Ερώτημα 1

Ακολουθώντας τις οδηγίες της εκφώνησης έγινε λήψη του dataset (το αντίστοιχο αρχείο βρίσκεται στον παρών φάκελο με όνομα **data.csv**). Στη συνέχεια έγινε μια προ-επεξεργασία στα δεδομένα. Πιο συγκεκριμένα αντικαταστάθηκαν οι Nan τιμές με 0, ενώ αφαιρέθηκαν και οι διπλότυπες. Στο αρχείο **clean_data.csv** είναι αποθηκευμένα τα αποτελέσματα. Ταυτόχρονα τυπώθηκαν οι στήλες dataset για καλύτερη κατανόηση της πληροφορίας.

Όπως είναι γνωστό, το dataset περιέχει πληροφορίες για την ροή διάφορων δικτύων δηλαδή την επικοινωνία μεταξύ πηγής και προορισμού (αντίστοιχες στήλες Src IP/Dst IP).

Ταυτόχρονα περιέχει και κάποιες άλλες αναγνωριστικές πληροφορίες για τα παραπάνω όπως το χρονικό της καταγραφής της ροής τους, το αναγνωριστικό τους, το πρωτόκολλο κλπ.

Πέρα από αυτά περιέχει στήλες που δίνουν πληροφορίες για την ένταση και το μέγεθος της ροής (Flow Duration/Total Fwd Packet/ Total Length of Bwd Packet).

Αντίστοιχα, περιέχονται πληροφορίες για το μέγεθος των πακέτων που μεταδίδονται μεταξύ πηγής και προορισμού (Average Packet Size, Fwd Packet Length Max/Min/Mean/Std) και αντίστοιχα πληροφορίες για τον ρυθμό μετάδοσης τους (Flow Bytes/s).

Επιπρόσθετα, έχει στήλες που αναλύουν την σύνδεση στο πέραςμα του χρόνου (Active Mean/Std/Max/Min) και πληροφορίες για υπορροές.

Τέλος, περιέχονται στήλες που δίνουν πληροφορίες για το είδος της επίθεσης (Malicious/Benign), το είδος και τις υποκατηγορίες τους.

➡ Οι παραπάνω πληροφορίες για την κάθε στήλη εξήχθησαν μετά από αναζήτηση στο διαδίκτυο, υπάρχουσες γνώσεις πάνω στην επικοινωνία των δικτύων καθώς και από την ιστοσελίδα από όπου έγινε η λήψη του dataset, όπου και αναλύονται εκτενέστερα οι κατηγορίες και οι υποκατηγορίες του κάθε Label ροής.

Εφόσον είναι εμφανές ότι οι ορισμένες στήλες έχουν κατηγορηματικά γνωρίσματα, έγινε ένας διαχωρισμός σε δύο ξεχωριστά αρχεία (**categorical_columns.csv/**
numeric_columns.csv) με σκοπό να υπολογιστούν τα αντίστοιχα στατιστικά μεγέθη για το κάθε είδος.

Όσον αφορά τις αριθμητικές στήλες υπολογίστηκαν

Μέσος Όρος: ο αριθμός γύρω από τον οποίο τείνουν να συγκεντρώνονται οι τιμές

Διακύμανση: δείχνει την διασπορά των τιμών γύρω από τον μέσο όρο

Λοξότητα: δείχνει την συμμετρία της κατανομής, δηλαδή αν οι τιμές τίνουν προς κάποια πλευρά (δεξιά ή αριστερά)

Κύρτωση: δείχνει την ύπαρξη ακαραίων τιμών ανάλογα με το πόσο πλατιά είναι η κορυφή

Αντίθετα για τις κατηγορηματικές στήλες, για τις οποίες δεν θα είχαν νόημα οι παραπάνω υπολογισμοί υπολογίστηκαν

- **Συχνότερη τιμή**
- **Πλήθος εμφάνισης της παραπάνω τιμής**

Τα αποτελέσματα βρίσκονται στα αντίστοιχα αρχεία `categorical_stats.csv` και `numeric_stats.csv`

Αφού υπολογίστηκαν τα παραπάνω στατιστικά μεγέθη για κάθε στήλη, σχεδιάστηκαν ορισμένες γραφικές παραστάσεις για κάποιες από αυτές ή συνδυαστικά με άλλες

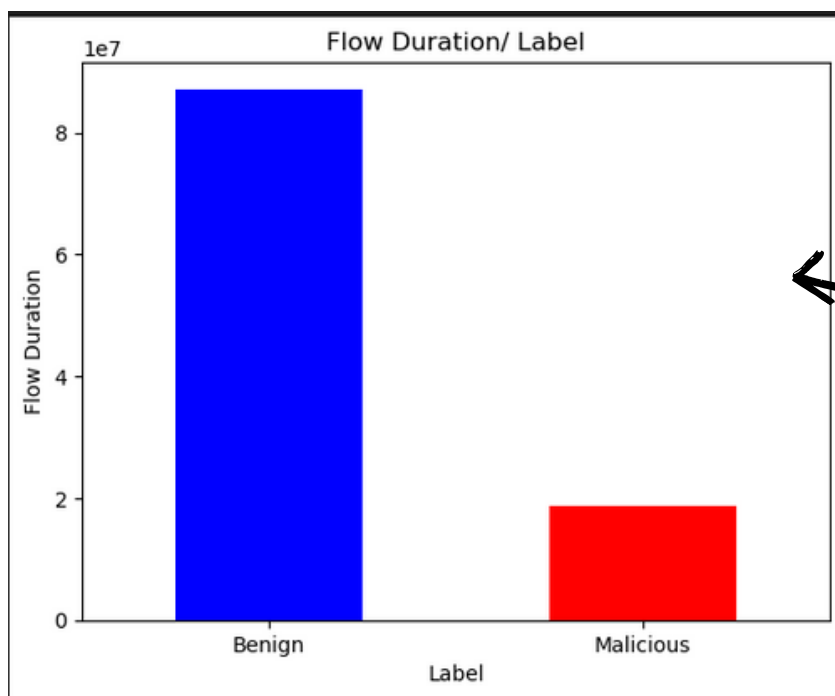
Πιο συγκεκριμένα:

Συνδυάστηκε η στήλη `Label` που αποτελεί και την βασική διάκριση του dataset ανάμεσα στις ροές (`Benign/Malicious`) με τις εξής πέντε στήλες:

- **Flow Duration:** αναφέρεται στην συνολική διάρκεια της ροής σε ms, ποικίλλει ανάλογα με το είδος της επίθεσης
- **Flow Bytes/s:** αποτελεί τον ρυθμό μεταφοράς δεδομένων
- **Total Fwd Packet:** είναι το σύνολο των πακέτων που στάλθηκαν μεταξύ client-server
- **Flow Packets/s:** είναι ο ρυθμός πακέτων ανά sec, δηλαδή πόσο γρήγορα στέλνονται τα πακέτα στη εκάστοτε ροή

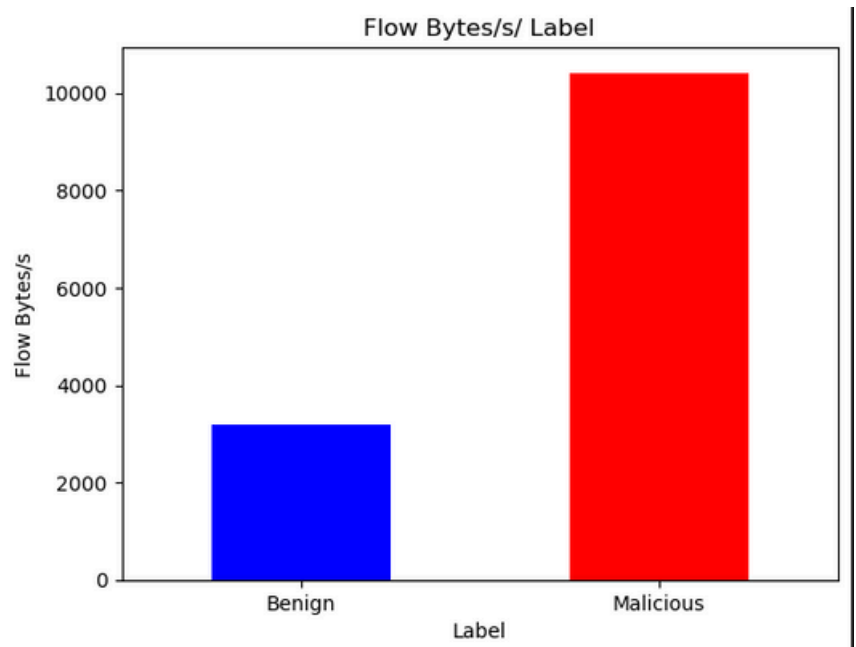
➡ Τα παραπάνω αποτελούν κάποια ενδεικτικά μεγέθη/στήλες του dataset και χρησιμοποιήθηκαν για να δείξουν πως μεταβάλλονται μεταξύ των 2 τιμών του label, ώστε να φανεί η όποια συσχέτιση υπάρχει μεταξύ τους

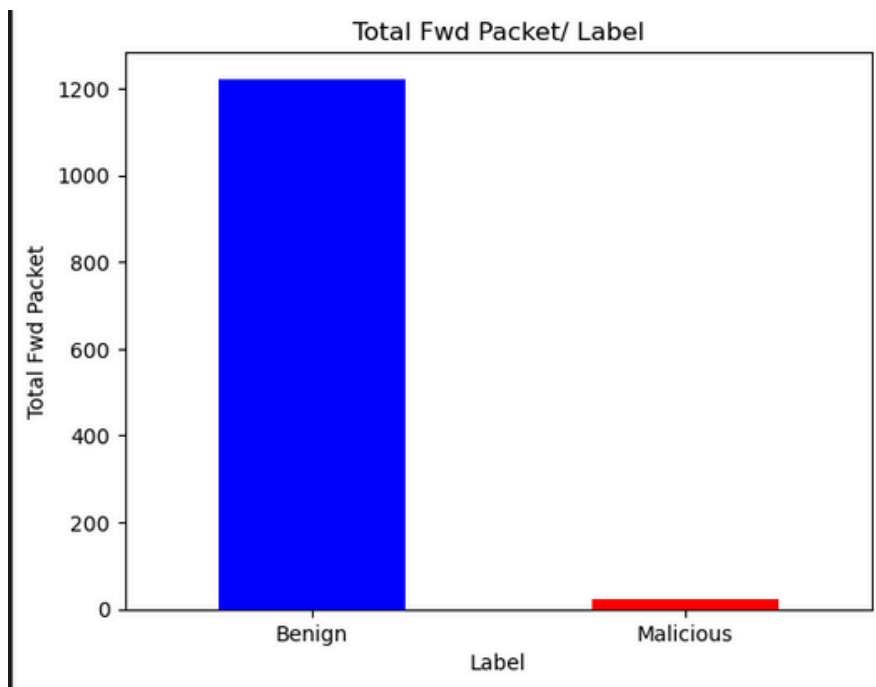
➡ Να σημειωθεί ότι το dataset έχει μεγάλη ανισορροπία στον αριθμό των Malicious ροών έναντι των Benign, και κατ' επέκταση επηρεάζονται και τα αποτελέσματα



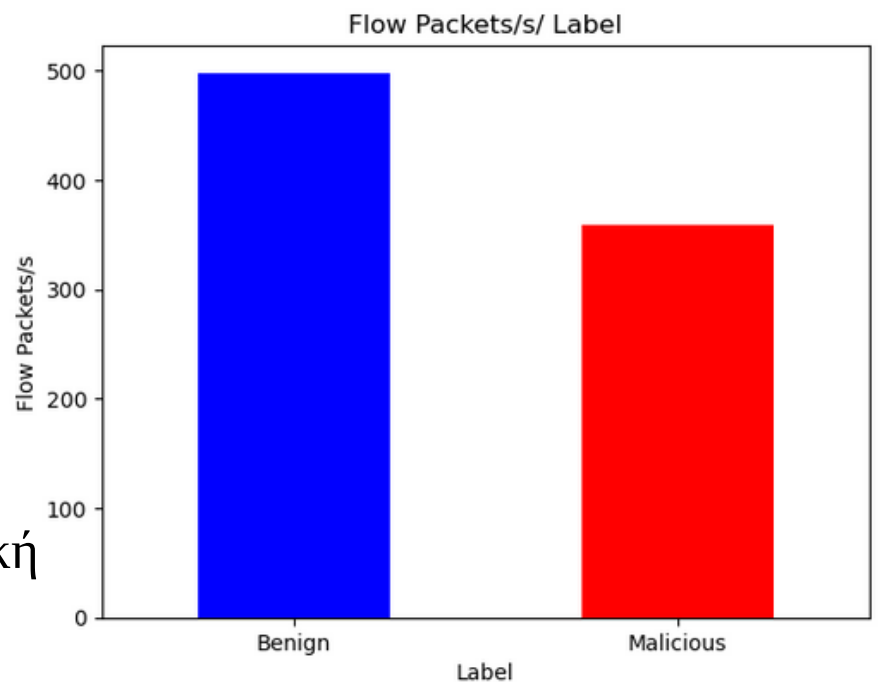
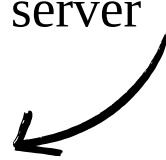
Benign έχουν πολύ μεγαλύτερη διάρκεια ροής

Στα Malicious γίνεται μαζική μεταφορά δεδομένων





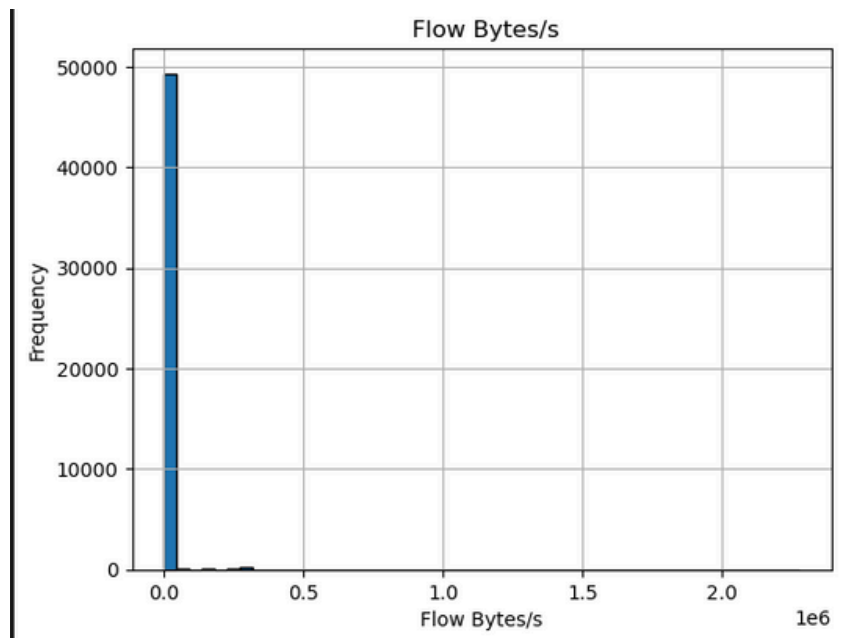
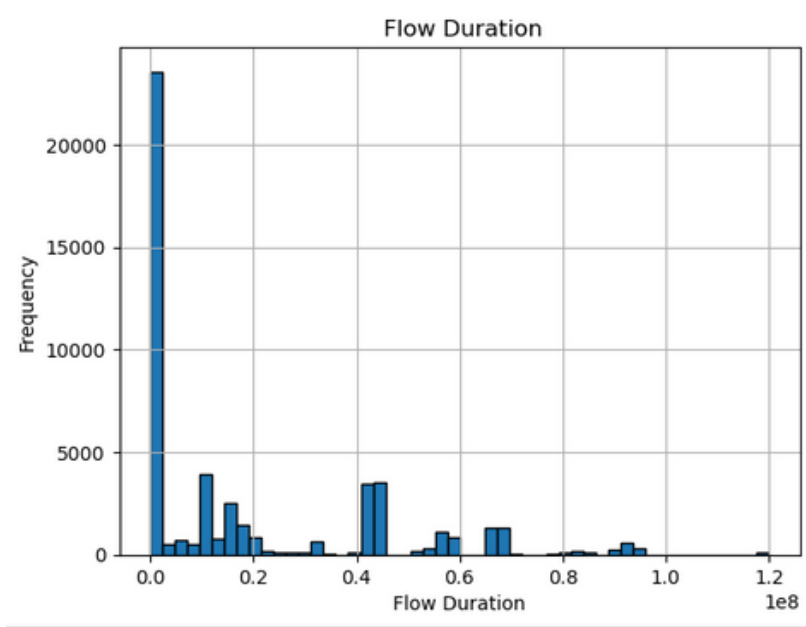
Στα Malicious
αποστέλλονται
συνολικά ελάχιστα
πακέτα προς τον
server

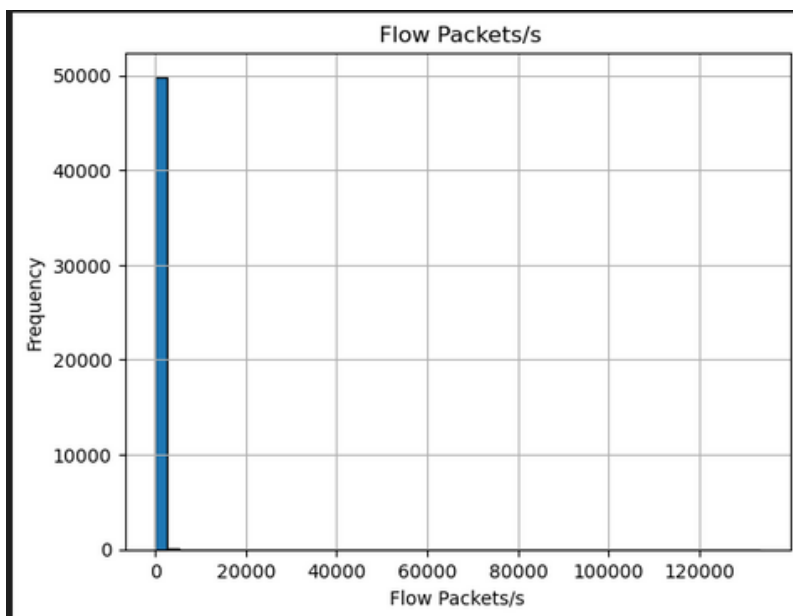
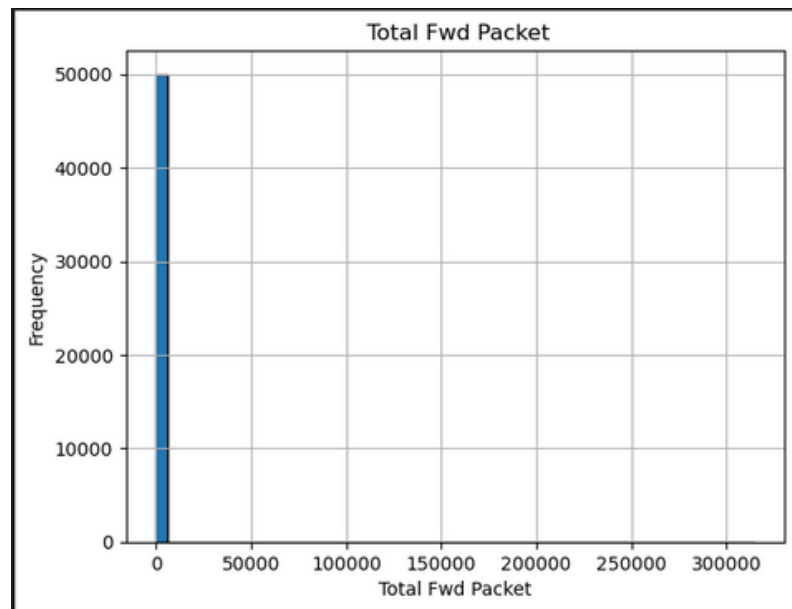


Δεν υπάρχει χαρακτηριστική
διαφορά στην μεταφορά
πακέτων ανά sec στις δύο
κατηγορίες, βέβαια στις
Benign είναι περισσότερα

Ταυτόχρονα δημιουργήθηκαν και τα αντίστοιχα ιστογράφηματα για τις παρακάτω στήλες, ώστε να φανεί συνολικά η κατανομή των τιμών για το καθένα.

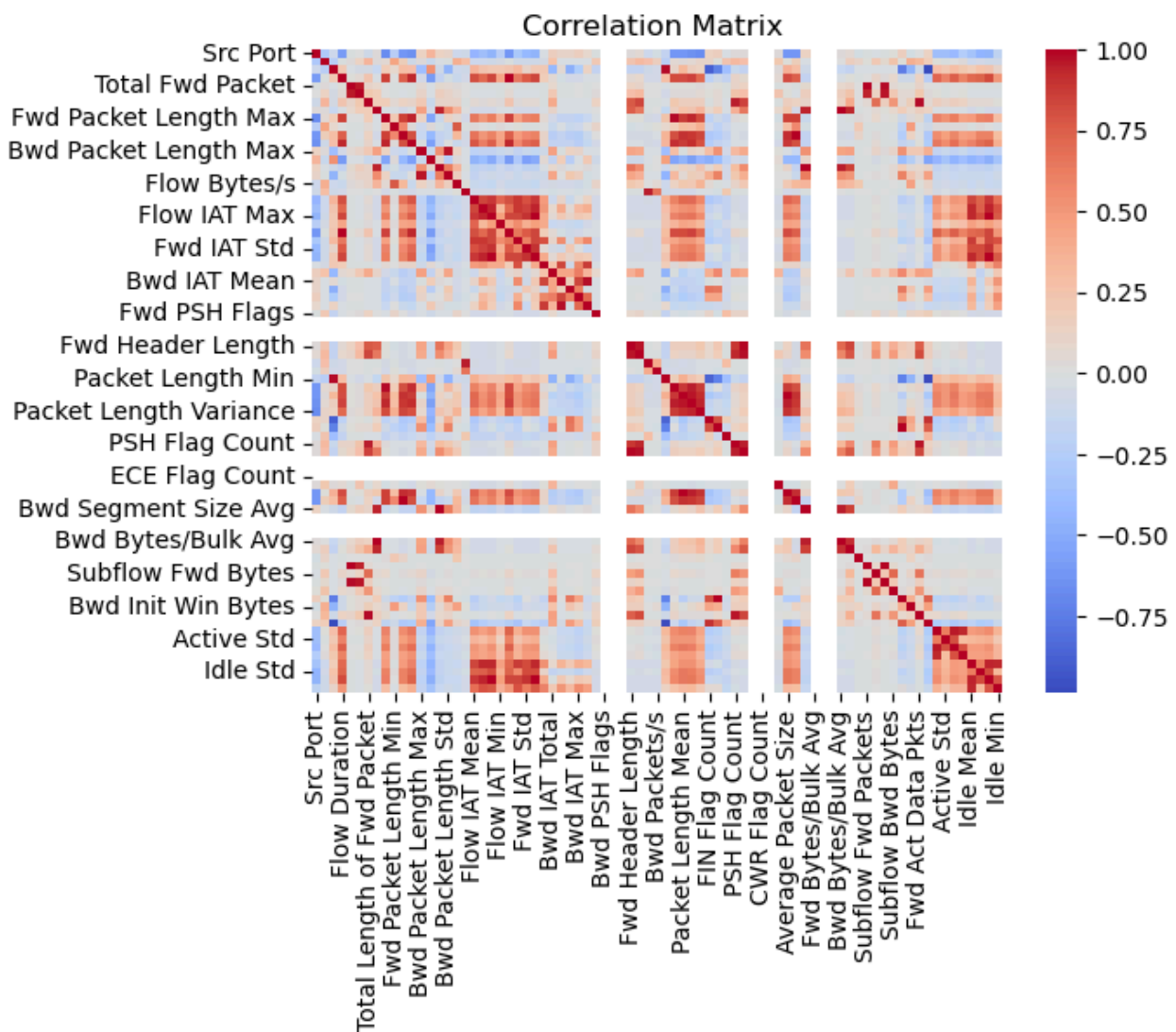
Τα γραφήματα φαίνονται παρακάτω:





Γενικά και στα 4 χαρακτηριστικά βλέπουμε το ίδιο μοτίβο: Οι κατανομές είναι κανονικές και χαμηλές, υπάρχουν μερικές ακραίες τιμές. Το συμπέρασμα είναι πως οι περισσότερες ροές είναι σύντομες, με λίγα δεδομένα και χαμηλό ρυθμό.

Στη συνέχεια προκειμένου να φανεί η συσχέτιση μεταξύ των στηλών του dataset, δημιουργήθηκε ένας correlation heatmap, ώστε να φανούν οι όποιες ισχυρές ή όχι σχέσεις μεταξύ των στηλών, ανάλογα με το πόσο θερμό είναι το χρώμα στο γράφημα.



Γενικά φαίνονται κάποιες περιοχές με ισχυρή συσχέτιση (έντονο κόκκινο) αλλά επειδή οι στήλες είναι πολλές, έπειτα τυπώθηκαν ζευγάρια στηλών από τον correlation matrix, με συσχέτιση μεγαλύτερη από 0,9 ώστε να βρεθούν στήλες που έχουν κοινή σχεδόν πληροφορία.

Παρακάτω φαίνονται τα αποτελέσματα:

```
- Flow IAT Std ↔ Flow IAT Max (corr = 0.94)
- Flow IAT Std ↔ Fwd IAT Max (corr = 0.94)
- Flow IAT Std ↔ Idle Mean (corr = 0.92)
- Flow IAT Std ↔ Idle Max (corr = 0.94)
- Flow IAT Max ↔ Fwd IAT Std (corr = 0.90)
- Flow IAT Max ↔ Fwd IAT Max (corr = 1.00)
- Flow IAT Max ↔ Idle Mean (corr = 0.93)
- Flow IAT Max ↔ Idle Max (corr = 0.99)
- Fwd IAT Std ↔ Fwd IAT Max (corr = 0.90)
- Fwd IAT Max ↔ Idle Mean (corr = 0.93)
- Fwd IAT Max ↔ Idle Max (corr = 0.99)
- Bwd IAT Mean ↔ Bwd IAT Min (corr = 0.98)
- Fwd Header Length ↔ Bwd Header Length (corr = 0.96)
- Fwd Header Length ↔ ACK Flag Count (corr = 0.99)
- Bwd Header Length ↔ PSH Flag Count (corr = 0.93)
- Bwd Header Length ↔ ACK Flag Count (corr = 0.98)
- Packet Length Min ↔ Fwd Seg Size Min (corr = -0.97)
- Packet Length Max ↔ Packet Length Std (corr = 0.95)
- Packet Length Mean ↔ Packet Length Std (corr = 0.96)
- Packet Length Mean ↔ Packet Length Variance (corr = 0.92)
- Packet Length Mean ↔ Average Packet Size (corr = 1.00)
- Packet Length Mean ↔ Fwd Segment Size Avg (corr = 0.91)
- Packet Length Std ↔ Packet Length Variance (corr = 0.95)
- Packet Length Std ↔ Average Packet Size (corr = 0.95)
- Packet Length Variance ↔ Average Packet Size (corr = 0.90)
- SYN Flag Count ↔ Fwd Init Win Bytes (corr = 1.00)
- PSH Flag Count ↔ ACK Flag Count (corr = 0.95)
- PSH Flag Count ↔ Fwd Act Data Pkts (corr = 0.96)
- Average Packet Size ↔ Fwd Segment Size Avg (corr = 0.91)
- Bwd Segment Size Avg ↔ Bwd Bytes/Bulk Avg (corr = 0.95)
- Bwd Bytes/Bulk Avg ↔ Bwd Packet/Bulk Avg (corr = 0.91)
```

```
- Flow IAT Std ↔ Flow IAT Max (corr = 0.94)
- Flow IAT Std ↔ Fwd IAT Max (corr = 0.94)
- Flow IAT Std ↔ Idle Mean (corr = 0.92)
- Flow IAT Std ↔ Idle Max (corr = 0.94)
- Flow IAT Max ↔ Fwd IAT Std (corr = 0.90)
- Flow IAT Max ↔ Fwd IAT Max (corr = 1.00)
- Flow IAT Max ↔ Idle Mean (corr = 0.93)
- Flow IAT Max ↔ Idle Max (corr = 0.99)
- Fwd IAT Std ↔ Fwd IAT Max (corr = 0.90)
- Fwd IAT Max ↔ Idle Mean (corr = 0.93)
- Fwd IAT Max ↔ Idle Max (corr = 0.99)
- Bwd IAT Mean ↔ Bwd IAT Min (corr = 0.98)
- Fwd Header Length ↔ Bwd Header Length (corr = 0.96)
- Fwd Header Length ↔ ACK Flag Count (corr = 0.99)
- Bwd Header Length ↔ PSH Flag Count (corr = 0.93)
- Bwd Header Length ↔ ACK Flag Count (corr = 0.98)
- Packet Length Min ↔ Fwd Seg Size Min (corr = -0.97)
- Packet Length Max ↔ Packet Length Std (corr = 0.95)
- Packet Length Mean ↔ Packet Length Std (corr = 0.96)
- Packet Length Mean ↔ Packet Length Variance (corr = 0.92)
- Packet Length Mean ↔ Average Packet Size (corr = 1.00)
- Packet Length Mean ↔ Fwd Segment Size Avg (corr = 0.91)
- Packet Length Std ↔ Packet Length Variance (corr = 0.95)
- Packet Length Std ↔ Average Packet Size (corr = 0.95)
- Packet Length Variance ↔ Average Packet Size (corr = 0.90)
- SYN Flag Count ↔ Fwd Init Win Bytes (corr = 1.00)
- PSH Flag Count ↔ ACK Flag Count (corr = 0.95)
- PSH Flag Count ↔ Fwd Act Data Pkts (corr = 0.96)
- Average Packet Size ↔ Fwd Segment Size Avg (corr = 0.91)
- Bwd Segment Size Avg ↔ Bwd Bytes/Bulk Avg (corr = 0.95)
- Bwd Bytes/Bulk Avg ↔ Bwd Packet/Bulk Avg (corr = 0.91)
- Subflow Fwd Packets ↔ Subflow Bwd Packets (corr = 0.99)
- Active Mean ↔ Active Max (corr = 0.96)
- Active Mean ↔ Active Min (corr = 0.93)
- Idle Mean ↔ Idle Max (corr = 0.94)
- Idle Mean ↔ Idle Min (corr = 0.94)
```

Ερώτημα 2

Προκειμένου να μειωθούν οι στήλες, θα χρησιμοποιηθούν τα αποτελέσματα που προέκυψαν από το προηγούμενο ερώτημα. Συγκεκριμένα από τα ισχυρώς συσχετιζόμενα ζευγάρια που προέκυψαν θα γίνει διαλογή και θα κρατηθεί ένα από το καθένα. Όσον αφορά τις κατηγορικές στήλες αφαιρέθηκαν οι:

- Flow ID
- Src IP
- Dst IP
- Timestamp

Άρα οι τελικές στήλες είναι οι παρακάτω:

- | | |
|------------------------------|-------------------|
| • Flow Duration | • Active Mean |
| • Total Fwd Packet | • Active Std |
| • Total Bwd packets | • Idle Mean |
| • Total Length of Fwd Packet | • Idle Std |
| • Total Length of Bwd Packet | • Label |
| • Flow Bytes/s | • Traffic Type |
| • Flow Packets/s | • Traffic Subtype |
| • Packet Length Mean | |
| • Subflow Fwd Bytes | |
| • Subflow Bwd Bytes | |

Όσον αφορά την μείωση των γραμμών χρησιμοποιήθηκαν 2 τρόποι όπως ζητήθηκε.

1. Με δειγματοληψία

Επειδή υπήρχε μεγάλη ανισορροπία μεταξύ των Malecious και Benign όπως αναφέρθηκε παραπάνω, η δειγματοληψία που έγινε αφαίρεσε κάποιες από τις αρχικές γραμμές των Malecious.

Αρχικά: **Malicious:** 8.654.324 και **Benign:** 1301

Μετά την δειγματοληψία: **Malicious:** 2500 και **Benign:** 1301

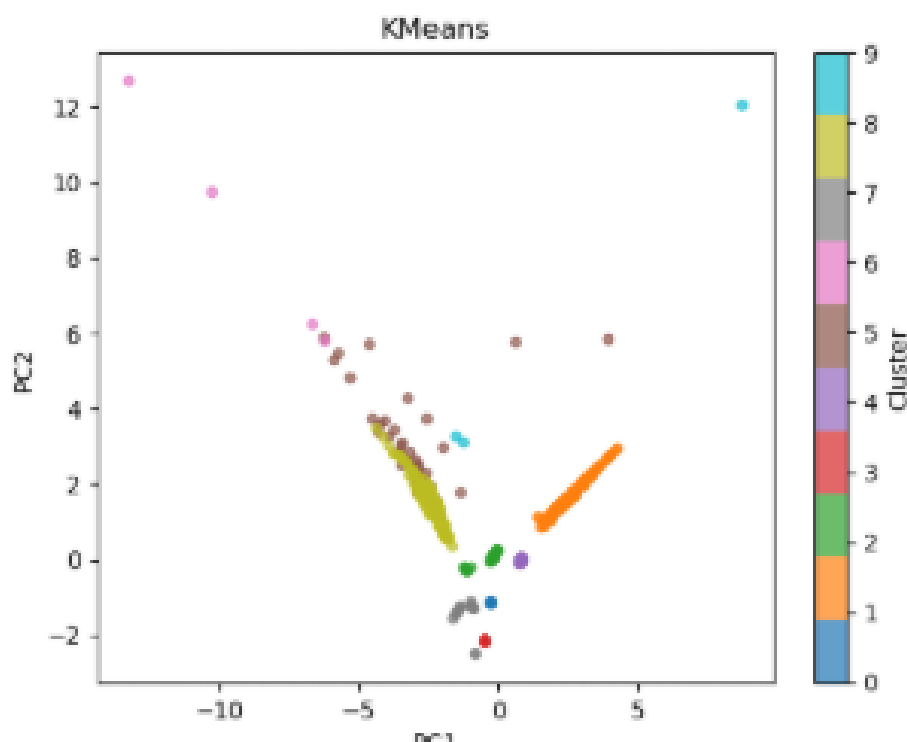
2. Με Clustering

- **K-Means**

Χρησιμοποιήθηκε ο συγκεκριμένος αλγόριθμος clustering καθώς είναι αποδοτικός και σε τόσο μεγάλα datasets. Ο αριθμός των clusters ορίστηκε χωρίς την χρήση κάποια άλλης μεθόδου στις 10, καθώς έχουμε πάνω από 10 στήλες και υπάρχουν πάνω από 2 υποκατηγορίες σε κάποιες από τις στήλες.

Οπότε για να μην φτιαχτούν τόσες ομάδες όσες και οι στήλες/υποκατηγορίες και μετά από δοκιμές επιλέχθηκε το $k=10$. Στη συνέχεια επιλέγουμε τις 200 πιο κοντινές τιμές στο κάθε cluster ώστε να υπάρχει αρκετό δείγμα για τα επόμενα ερωτήματα, αλλιώς οι γραμμές θα είναι τόσες όσες τα clusters. Στη συνέχεια εφαρμόζουμε την μέθοδο PCA για καλύτερη οπτικοποίηση των στηλών και κατ' επέκταση του γραφήματος.

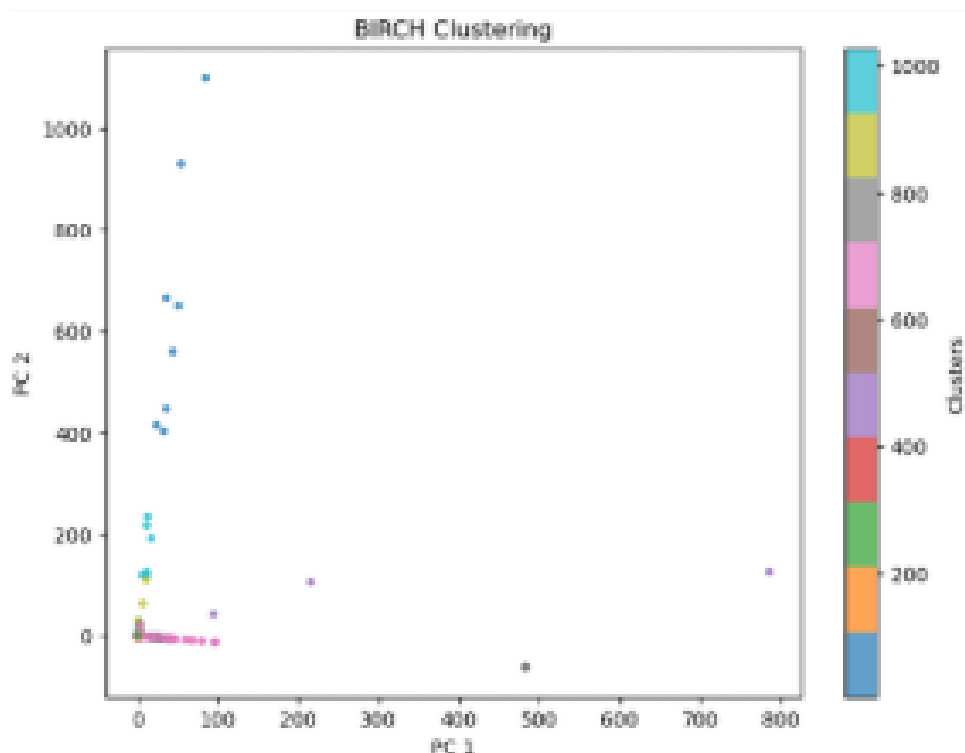
Παρακάτω φαίνεται το γράφημα:



- **Birch**

Σε αυτή τη μέθοδο χρειάστηκε να διαβαστεί το dataset σε chunks μεγέθους 2000 καθώς αλλιώς ήταν αδύνατο να διαβαστεί. Χρησιμοποιήθηκε γιατί λειτουργεί καλά σε τόσο μεγάλα datasets. Ανάλογα με το threshold δημιουργεί νέα subclusters να χρειαστεί αλλιώς μπαίνει σε κάποιο από όσα υπάρχουν ήδη. Το 1,5 που ορίστηκε είναι μια ενδιάμεση τιμή ώστε να μην δημιουργηθούν πολλά αλλά ούτε να δημιουργηθούν λίγα που θα χωρέσουν όλα τα data εκεί.

Παρακάτω φαίνεται το γράφημα:



➡ Γενικά και στις 2 μεθόδους clustering απομονώθηκαν οι αριθμητικές στήλες του dataset καθώς μεταξύ αυτών των τιμών. Όμως στο επόμενο ερώτημα τα αρχεία που προέκυψαν δεν μπορούσαν να χρησιμοποιηθούν για πρόβλεψη καθώς έλειπαν οι ζητούμενες στήλες. Έτσι λήφθηκαν υπόψιν στις μεθόδους clustering ώστε να χαρακτηριστούν τα clusters, ανάλογα με το ποιο εμφανίζεται περισσότερο

➡ Τα 3 αρχεία που προέκυψαν βρίσκονται αντίστοιχα στο παρών φάκελο και είναι αντίστοιχα τα

- sampled_data
- kmeans_reduced_data
- birch_reduced_data

Ερώτημα 3

Για αυτό το ερώτημα χρησιμοποιήθηκαν τα παραπάνω 3 νέα μειωμένα dataset και καθένα από αυτά εκπαιδεύτηκε με

- **SVM:** Στον αλγόριθμο χρησιμοποιήθηκε RBF kernel αφού ο διαχωρισμός του dataset δεν είναι γραμμικός
- **Neural Networks:** χρησιμοποιούνται 2 επίπεδα λόγω του ότι το dataset δεν είναι γραμμικό

Τα αποτελέσματα για την κάθε στήλη φαίνονται παρακάτω:

► Label

- **Sampled Data**

	precision	recall	f1-score	support
Benign	0.99	0.68	0.81	260
Malicious	0.86	1.00	0.92	501
accuracy			0.89	761
macro avg	0.92	0.84	0.87	761
weighted avg	0.90	0.89	0.88	761

Neural Network-Label

	precision	recall	f1-score	support
Benign	0.96	0.92	0.94	260
Malicious	0.96	0.98	0.97	501
accuracy			0.96	761
macro avg	0.96	0.95	0.95	761
weighted avg	0.96	0.96	0.96	761

Οι προβλέψεις Benign ήταν γενικά σωστές (precision=0.99)
Όμως χάθηκαν πολλά Benign (Recall = 0.68), ενώ η ισορροπία μεταξύ των δύο τιμών είναι καλή αλλά (F1-score = 0.81)
Όσον αφορά τις Malicious τιμές υπήρχε ένα ποσοστό 14% λάθος πρόβλεψης, αλλά δεν χάθηκε καμία τέτοια τιμή ενώ και το F1-score είναι πολύ καλό
Γενικά το μοντέλο είναι αρκετά καλό, αλλά χάνει αρκετά Benign.
Όσον αφορά το MLP νευρωνικό δίκτυο η πρόβλεψη είναι πολύ καλή και για τις δυο τιμές.

• **K-Means Data**

SVM-Label

	precision	recall	f1-score	support
Benign	0.97	1.00	0.99	35
Malicious	1.00	1.00	1.00	203
accuracy			1.00	238
macro avg	0.99	1.00	0.99	238
weighted avg	1.00	1.00	1.00	238

Neural Network-Label

	precision	recall	f1-score	support
Benign	1.00	0.97	0.99	35
Malicious	1.00	1.00	1.00	203
accuracy			1.00	238
macro avg	1.00	0.99	0.99	238
weighted avg	1.00	1.00	1.00	238

Σε αυτό το σύνολο δεδομένων και τα δυο μοντέλα αποδίδουν πάρα πολύ καλά. Όμως πρέπει να λάβουμε υπόψιν ότι ο αριθμός των Benign ήταν αρκετά μικρότερος, εάν άλλαζε αυτό το δείγμα τότε πιθανόν να είχαμε διαφορετικές μετρικές.

- **Birch Data**

SVM-Label

	precision	recall	f1-score	support
Benign	0.78	0.76	0.77	37
Malicious	0.93	0.94	0.93	128
accuracy			0.90	165
macro avg	0.85	0.85	0.85	165
weighted avg	0.90	0.90	0.90	165

Neural Network-Label

	precision	recall	f1-score	support
Benign	0.88	0.76	0.81	37
Malicious	0.93	0.97	0.95	128
accuracy			0.92	165
macro avg	0.90	0.86	0.88	165
weighted avg	0.92	0.92	0.92	165

Σε αυτό το σύνολο δεδομένων το MLP μοντέλο αποδίδει πολύ καλύτερα, αλλά γενικά και τα δυο μοντέλα έχουν πολύ καλή ακρίβεια στις Malicious τιμές ενώ οι Benign ξανά βρίσκονται σε δυσαναλογία.

➤ Traffic Type

- **Sampled Data**

	precision	recall	f1-score	support
Audio	0.95	0.53	0.68	38
Background	0.00	0.00	0.00	6
Bruteforce	0.75	0.99	0.85	180
DoS	1.00	0.99	1.00	321
Text	0.79	0.55	0.65	42
Video	0.86	0.75	0.80	174
accuracy			0.88	761
macro avg	0.73	0.64	0.66	761
weighted avg	0.89	0.88	0.88	761

Η ακρίβεια είναι αρκετά υψηλή, συγκεκριμένα στις τιμές DoS, Bruteforce έχουν υψηλό recall αλλά παρατηρείται χαμηλότερο precision. Γενικά αποδίδει καλά στα παραπάνω αλλά αποτυγχάνει εντελώς στο Background και δυσκολεύεται με Audio/Text.

- K-Means Data**

	precision	recall	f1-score	support
Bruteforce	0.07	0.75	0.13	4
DoS	0.93	0.39	0.55	164
Information Gathering	0.02	1.00	0.05	2
Mirai	1.00	0.15	0.26	33
Video	1.00	1.00	1.00	35
accuracy		0.46		238
macro avg	0.60	0.66	0.40	238
weighted avg	0.93	0.46	0.56	238

Neural Network -traffic type

	precision	recall	f1-score	support
Bruteforce	0.00	0.00	0.00	4
DoS	0.92	0.90	0.91	164
Information Gathering	0.00	0.00	0.00	2
Mirai	0.56	0.73	0.63	33
Video	1.00	0.97	0.99	35
accuracy		0.87		238
macro avg	0.50	0.52	0.51	238
weighted avg	0.86	0.87	0.86	238

Σε αυτό το dataset το μοντέλο έχει συνολικά χαμηλή απόδοση, καθώς οι περισσότερες κλάσεις αναγνωρίζονται με λάθη. Το νευρωνικό δίκτυο βελτίωσε σημαντικά την ακρίβεια, αλλά κάποιες κλάσεις εξακολουθούν να μην αναγνωρίζονται καθόλου.

- **Birch Data**

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

Audio	0.60	1.00	0.75	6	
Bruteforce	0.00	0.00	0.00	1	
DoS	0.95	0.50	0.66	103	
Information Gathering		0.09	0.83	0.17	6
Mirai	0.78	0.39	0.52	18	
Text	0.40	0.75	0.52	8	
Video	0.93	0.61	0.74	23	

accuracy			0.55	165
macro avg	0.54	0.58	0.48	165
weighted avg	0.85	0.55	0.63	165

Neural Network-Traffic Type

	precision	recall	f1-score	support	
Audio	0.60	1.00	0.75	6	
Bruteforce	0.00	0.00	0.00	1	
DoS	0.85	0.96	0.90	103	
Information Gathering		0.00	0.00	0.00	6
Mirai	0.78	0.39	0.52	18	
Text	0.62	0.62	0.62	8	
Video	0.95	0.91	0.93	23	
...					
accuracy			0.84	165	
macro avg	0.54	0.56	0.53	165	
weighted avg	0.80	0.84	0.81	165	

Το τελευταίο dataset αποδίδει ανάμεσα στα παραπάνω, έχει συνολικά μέτρια ακρίβεια. Το νευρωνικό δίκτυο ξανά βελτίωσε σημαντικά την ακρίβεια και κάποιες κατηγορίες, αλλά εξακολουθεί να έχει αδυναμία σε συγκεκριμένους τύπους επιθέσεων.

Γενικά οι διαφορές στην κάθε κατηγορία οφείλεται πιθανόν λόγω ανισορροπίας δεδομένων.