

ΜΗΧΑΝΙΚΗ ΜΑΘΗΣΗ

Απαλλακτική Εργασία

1.ΣΚΟΠΟΣ

Η εργασία αυτή θα ασχοληθεί με την ανάπτυξη μοντέλων μηχανικής μάθησης για το classification πακέτων δικτύου με βάση την εφαρμογή από την οποία προέρχονται (π.χ. HTTP, DHCP, DNS, TLS, Google κλπ). Σκοπός είναι η κατηγοριοποίηση της κυκλοφορίας δικτύου για τη βελτίωση της ασφάλειας, την αποτελεσματική κατανομή πόρων ή οποιαδήποτε άλλη ανάγκη απαιτείται. Ως dataset θα χρησιμοποιηθεί μια καταγραφή δικτύου από ένα πανεπιστήμιο της Κολομβίας¹ με συνολικά περίπου 3.5 εκατ. εγγραφές. Το dataset περιλαμβάνει 87 features με αρκετές λεπτομέρειες όπως μέγιστο/ελάχιστο/μέσο μέγεθος πακέτων, συνολικός αριθμός πακέτων προς κάθε κατεύθυνση κλπ.

Με βάση τα παραπάνω features το μοντέλο θα πρέπει να είναι ικανό να κατηγοριοποιεί παρόμοιες καταγραφές έτσι ώστε να υπάρχει ένα insight σχετικά με την κίνηση ενός δικτύου.

2.ΣΧΕΤΙΚΗ ΒΙΒΛΙΟΓΡΑΦΙΑ

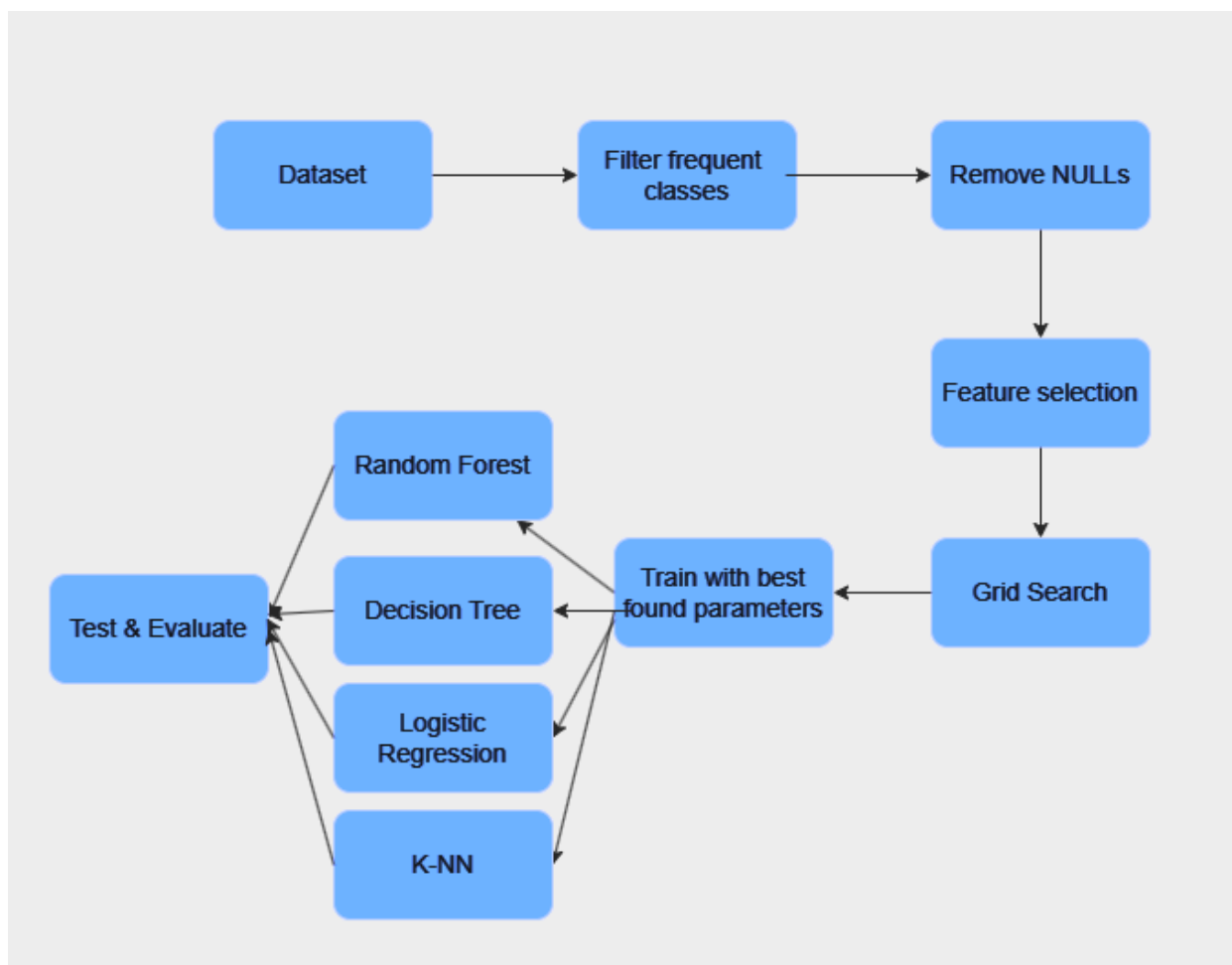
Η εργασία είναι εμπνευσμένη από την έρευνα *Network Traffic Classification Using Machine Learning for Software Defined Networks*²(Kuranage et al., 2019). Μια βασική διαφορά της εργασίας αυτής σε σχέση με την έρευνα είναι πως η εργασία βασίζεται σε supervised learning με έτοιμες κλάσεις. Επιπλέον, αξιοσημείωτο είναι ότι η έρευνα χρησιμοποίησε ακριβώς το ίδιο dataset όμως για διαφορετικούς σκοπούς και με διαφορετικό τρόπο. Συγκεκριμένα χρησιμοποιήθηκε ο αλγόριθμος k-means για το clustering και την δημιουργία των κλάσεων και έπειτα το classification με βάση αυτές. Αντίθετα, σε αυτή την εργασία θα αξιοποιήσουμε τις έτοιμες κλάσεις του dataset για το classification χωρίς την ανάγκη για clustering.

3.ΜΕΘΟΔΟΛΟΓΙΑ

Για την υλοποίηση των μοντέλων ακολουθήθηκε η διαδικασία της *Εικόνας 1*.

¹ <https://www.kaggle.com/datasets/jsrojas/ip-network-traffic-flows-labeled-with-87-apps>

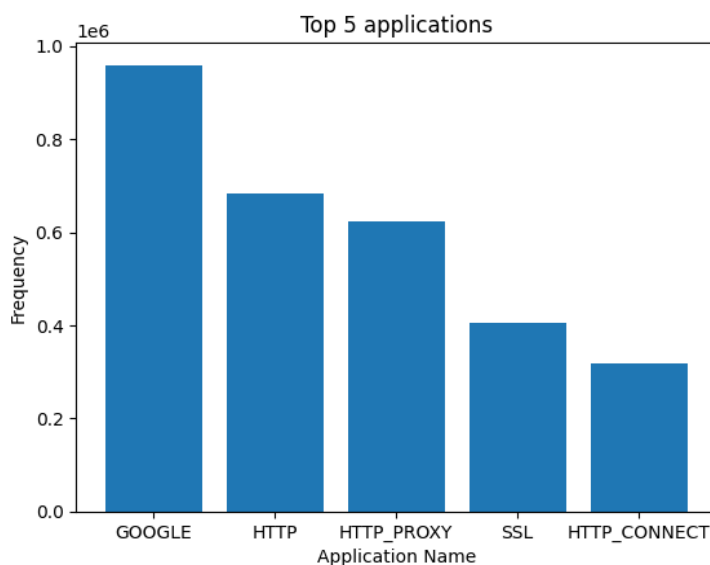
² Kuranage, M. P. J., Piamrat, K., & Hamma, S. *Network Traffic Classification Using Machine Learning for Software Defined Networks*. 2nd International Conference on Machine Learning for Networking (MLN), Dec 2019, Paris, France.



Εικόνα 1: Flowchart

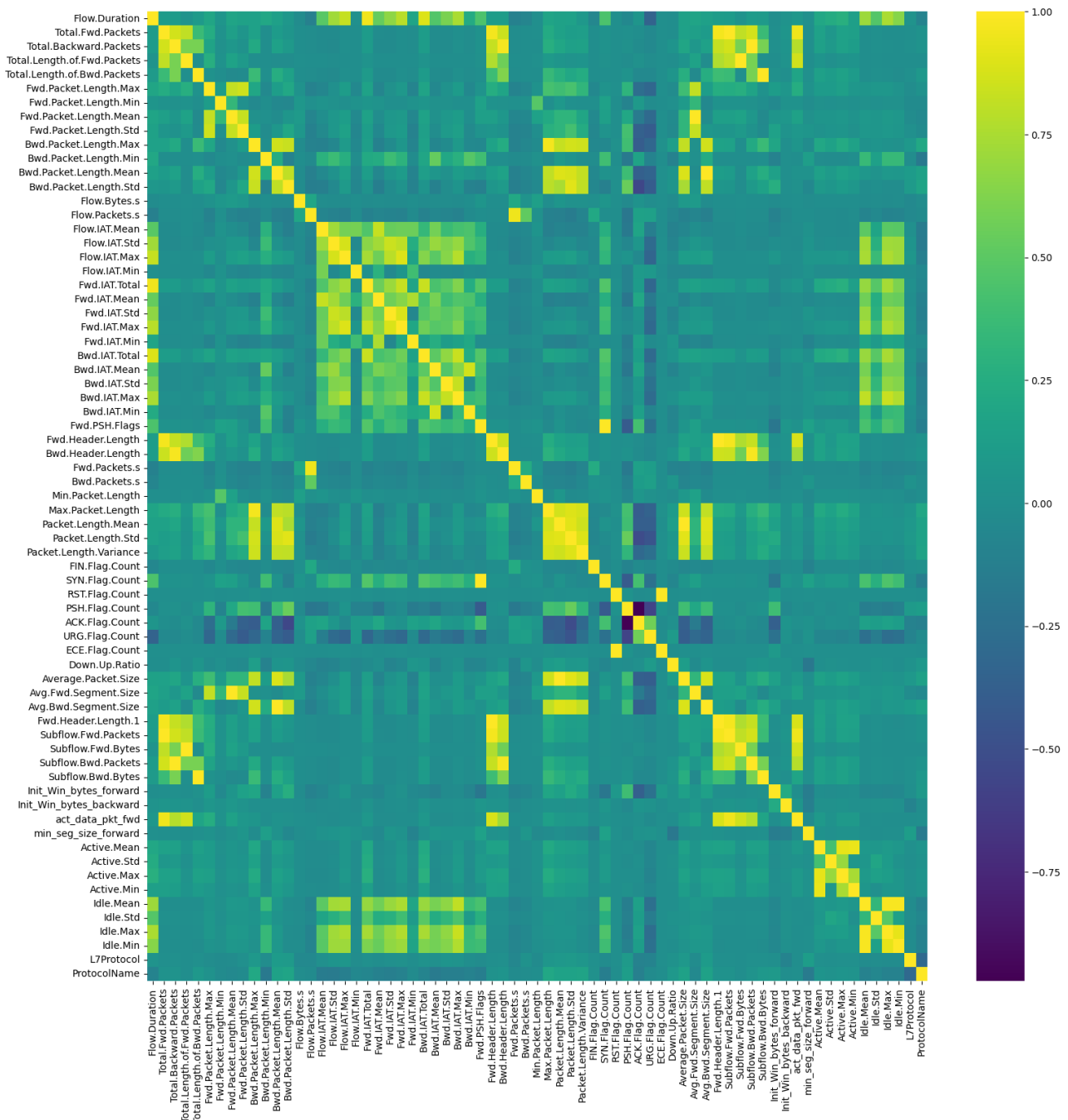
Προεπεξεργασία

Καθώς το μέγεθος του dataset ήταν αρκετά μεγάλο, επιλέχθηκαν μόνο κλάσεις με παραπάνω από 10.000 εγγραφές. Μετά από την αφαίρεση και των τιμών null το τελικό dataset περιείχε 170.000 εγγραφές με 87 features.



Εικόνα 2: Ιστόγραμμα κλάσεων

Για την επιλογή των features χρησιμοποιήθηκε η pearson correlation και παρέμειναν τα features με απόλυτη συσχέτιση μεγαλύτερη από 0,01. Εάν υπήρχαν παραπάνω από ένα features με την ίδια συσχέτιση επιλέχθηκε τυχαία ένα από αυτά. Επιπλέον, αφαιρέθηκαν και οι στήλες που δεν προσφέρουν χρήσιμες πληροφορίες για την συσχέτιση όπως Flow.ID, Source.IP, Timestamp, Destination.IP κλπ. Έτσι, το τελικό σχήμα του dataset ήταν (170.000, 29)



Εικόνα 3: Correlation Matrix

Training and Tuning

Τα μοντέλα μάθησης και οι παράμετροι που επιλέχθηκαν για την εκπαίδευση τους είναι τα εξής:

- **Random Forest**

- $n_estimators$: 100, 150, 200

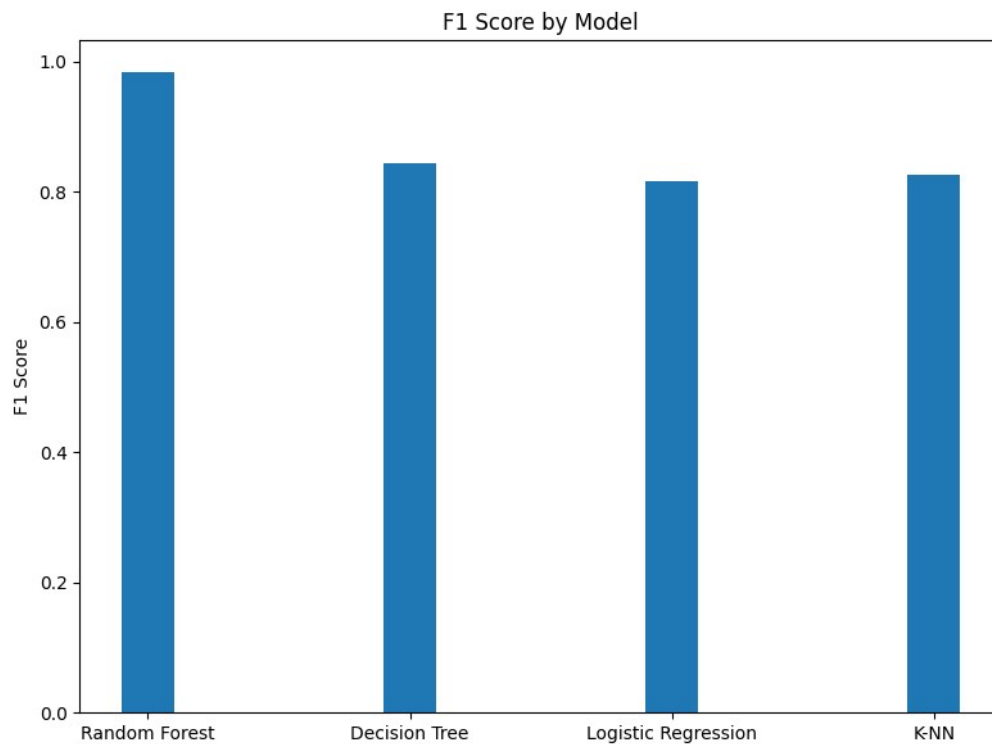
- **Decision Tree**
 - max_depth: 3, 4, 5, 6
 - min_samples_split: 2, 5, 10
 - min_samples_leaf: 1, 2, 4
- **Logistic Regression**
 - solver: lbfgs, sag
 - max_iter: 100, 200, 300
- **K-NN**
 - n_neighbors: 3, 5, 7, 9

Για την επιλογή των καλύτερων παραμέτρων χρησιμοποιήθηκε Grid Search με 5 folds στο training set. Έπειτα έγινε εκπαίδευση του μοντέλου με τις καλύτερες παραμέτρους και τέλος οι εκτιμήσεις των αποτελεσμάτων με βάση το test set.

4.ΑΠΟΤΕΛΕΣΜΑΤΑ

Παρακάτω φαίνονται συγκεντρωτικά οι αποδόσεις των τεσσάρων μοντέλων στο συγκεκριμένο πρόβλημα.

Model	F1 Score
Random Forest	98.3%
Decision Tree	84.2%
Logistic Regression	81.6%
K-NN	82.6%



Εικόνα 4: Σύγκριση F1 Score

Random Forest

Best parameters

- `n_estimators: 200`

		Predicted															
		GOOGLE	HTTP	HTTP_PROXY	SSL	HTTP_CONNECT	YOUTUBE	AMAZON	MICROSOFT	WINDOWS_UPDATE	GMAIL	FACEBOOK	SKYPE	DROPBOX	YAHOO	CLOUDFLARE	MSN
Actual	MSN	1	0	1	2	14	6	0	0	2	0	0	25	1	4	0	1921
	CLOUDFLARE	0	0	0	0	1	0	0	1	0	0	5	0	8	1	0	1969
	TWITTER	3	0	0	0	0	6	0	2	12	1	0	3	0	0	1974	0
	YAHOO	0	0	16	1	14	0	0	0	3	0	0	4	2	1957	2	0
	DROPBOX	0	0	0	2	1	0	0	0	1	0	0	0	2004	5	0	6
	FACEBOOK	1	0	2	0	1	9	0	0	3	0	0	1949	0	3	2	29
	SKYPE	0	0	0	1	0	0	0	0	0	0	1976	0	2	0	0	0
	WINDOWS_UPDATE	3	0	0	0	0	0	0	0	0	1989	0	0	0	0	1	0
	GMAIL	1	0	0	0	1	5	0	10	1965	0	0	1	1	0	3	3
	MICROSOFT	0	0	1	0	0	11	0	1974	12	0	0	7	0	0	2	3
	AMAZON	0	0	0	0	0	0	2017	0	0	0	2	0	3	0	0	0
	YOUTUBE	3	0	2	0	4	1954	0	5	4	0	0	30	0	0	2	9
	HTTP_CONNECT	0	0	12	2	1911	1	0	0	0	0	0	7	1	16	0	32
	SSL	0	0	3	1976	0	0	0	0	0	0	0	1	5	26	0	1
	HTTP_PROXY	0	0	1874	2	53	2	0	3	1	0	0	3	2	22	0	15
	HTTP	0	2033	0	0	0	0	0	0	1	0	0	0	0	0	0	0
	GOOGLE	1985	0	0	0	0	1	0	1	3	3	0	2	0	0	11	0

Confusion Matrix Random Forest

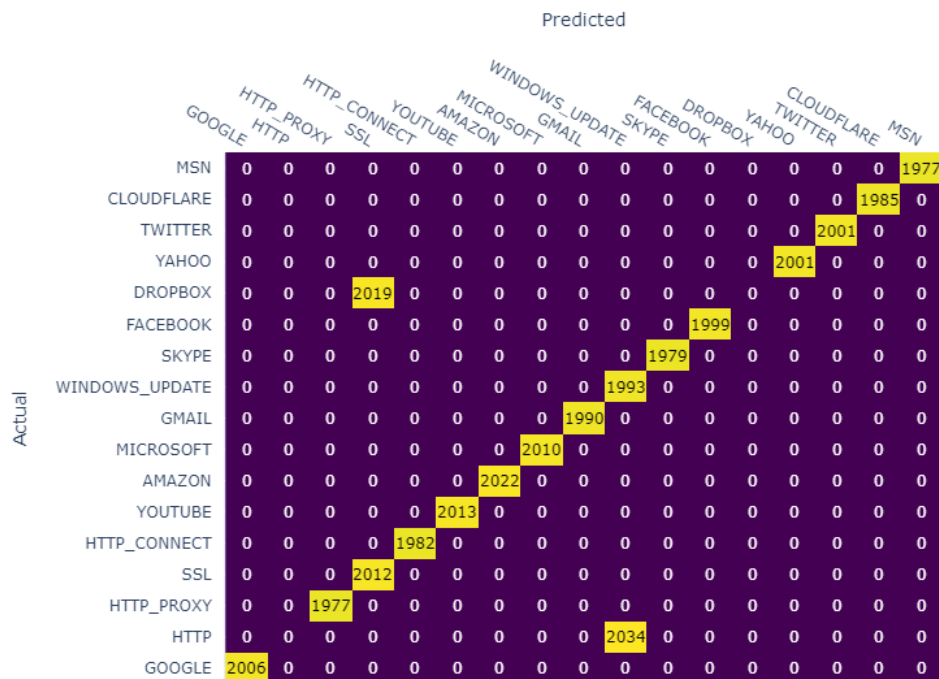
Εικόνα 5: Random Forest

Το μοντέλο Random Forest επιτυγχάνει υψηλή ακρίβεια και υπερτερεί των άλλων με ελάχιστο αριθμό misclassification. Άρα το συγκεκριμένο μοντέλο είναι καλύτερο στο να αναγνωρίζει διαφοροποιήσεις στα δεδομένα, οδηγώντας σε πιο ακριβείς προβλέψεις σε διαφορετικές κατηγορίες.

Decision Tree

Best parameters

- max_depth: 6
- min_samples_split: 2
- min_samples_leaf: 1



Confusion Matrix Decision Tree

Εικόνα 6: Decision Tree

Το μοντέλο Decision Tree δείχνει υψηλή ακρίβεια με μια πρώτη ματιά. Έχει ελάχιστες αλλά σοβαρές εσφαλμένες ταξινομήσεις. Κάθε κατηγορία προβλέπεται κυρίως σωστά, αλλά παρουσιάζει overfit εφόσον η ακρίβεια φτάνει το 100% στις περισσότερες κλάσεις. Επομένως, παρόλο που είναι το αμέσως καλύτερο μοντέλο μετά το Random Forest, με βάση το F1 score, δεν είναι το καταλληλότερο για αυτή την περίπτωση.

Logistic Regression

Best parameters:

- solver: lbfgs
- max_iter: 200

		Predicted															
		GOOGLE	HTTP_PROXY	HTTP_CONNECT	YOUTUBE	AMAZON	MICROSOFT	GMAIL	WINDOWS_UPDATE	FACEBOOK	SKYPE	DROPBOX	YAHOO	CLOUDFLARE	TWITTER	MSN	
Actual	MSN	0	0	17	7	166	215	0	0	0	0	0	287	0	4	1	1280
	CLOUDFLARE	0	0	0	0	0	0	0	0	0	260	0	0	0	0	1725	0
	TWITTER	0	0	0	0	0	0	2	17	0	0	0	0	0	1982	0	0
	YAHOO	0	0	66	714	223	0	0	0	0	0	0	0	987	0	0	11
	DROPBOX	0	0	1	2	0	0	0	0	0	1	0	2015	0	0	0	0
	FACEBOOK	0	0	25	10	11	403	0	0	2	0	0	1406	0	0	0	142
	SKYPE	0	0	0	0	0	0	0	0	0	1694	0	0	0	0	285	0
	WINDOWS_UPDATE	0	12	0	0	0	0	0	0	1981	0	0	0	0	0	0	0
	GMAIL	0	0	0	0	0	1	0	788	1174	0	0	4	0	0	20	3
	MICROSOFT	0	0	0	0	0	7	0	1650	337	0	0	7	0	0	1	8
	AMAZON	0	0	0	0	0	0	2022	0	0	0	0	0	0	0	0	0
	YOUTUBE	0	0	2	0	12	1488	0	30	15	0	0	385	0	0	0	81
	HTTP_CONNECT	0	0	86	383	1144	0	0	0	0	0	14	0	288	0	0	67
	SSL	0	0	18	1795	12	0	0	0	0	0	0	0	187	0	0	0
	HTTP_PROXY	0	0	1443	214	161	3	0	0	0	0	5	0	145	0	0	6
	HTTP	0	2032	0	0	0	0	0	0	2	0	0	0	0	0	0	0
	GOOGLE	1994	0	0	0	0	0	0	0	1	0	0	0	0	11	0	0
Confusion Matrix Logistic Regression																	

Confusion Matrix Logistic Regression

Εικόνα 7: Logistic Regression

Η Logistic Regression εμφανίζει ένα μέτριο επίπεδο ακρίβειας με αξιοσημείωτα misclassifications. Ενώ το μοντέλο ταξινομεί σωστά τις περισσότερες κατηγορίες, δυσκολεύεται με άλλες, εμφανίζοντας μεγάλες τιμές εκτός της διαγώνιου. Αυτό δείχνει ότι ενώ μπορεί να καταγράψει γραμμικές σχέσεις, δεν μπορεί να διαχειριστεί πιο περίπλοκα μοτίβα στα δεδομένα.

K-NN

Best parameters:

- n_neighbors: 3

		Predicted																	
		GOOGLE	HTTP_PROXY	HTTP_PROXY	SSL	HTTP_CONNECT	YOUTUBE	AMAZON	MICROSOFT	WINDOWS_UPDATE	GMAIL	FACEBOOK	SKYPE	DROPBOX	YAHOO	CLOUDFLARE	MSN		
Actual	MSN	4	0	37	19	212	190	0	57	63	0	0	80	13	42	12	0	1248	
	CLOUDFLARE	1	0	12	0	19	7	0	3	2	0	179	4	34	5	0	1713	6	
	TWITTER	14	0	14	2	18	9	0	34	16	3	0	7	2	3	1865	0	14	
	YAHOO	2	0	73	48	173	100	2	50	33	2	0	96	6	1337	5	0	74	
	DROPBOX	2	0	43	7	47	16	3	8	5	0	28	20	1775	17	1	32	15	
	FACEBOOK	3	0	102	14	97	158	0	80	27	1	2	1394	5	48	11	0	57	
	SKYPE	1	0	10	6	8	6	4	3	1	0	1746	7	19	1	0	166	1	
	WINDOWS_UPDATE	33	39	1	1	1	1	0	5	2	1902	0	1	0	1	4	0	2	
	GMAIL	13	1	66	27	150	152	1	142	1295	0	0	28	11	21	28	0	55	
	MICROSOFT	4	0	34	12	60	81	0	1573	114	0	2	53	2	24	6	0	45	
	AMAZON	0	0	1	1	2	4	1981	1	2	0	15	1	3	0	0	7	4	
	YOUTUBE	7	0	47	19	286	1256	0	47	108	0	0	70	6	33	2	0	132	
	HTTP_CONNECT	5	0	73	19	1474	125	0	28	43	0	2	48	3	62	12	0	88	
	SSL	2	1	12	1924	8	10	0	6	8	0	1	10	1	23	0	0	6	
	HTTP_PROXY	4	1	1773	8	50	28	0	11	9	1	1	46	2	20	7	1	15	
	HTTP	19	1949	1	0	1	0	0	0	2	57	0	0	0	0	3	0	2	
	GOOGLE	1891	19	6	1	3	9	0	6	19	21	0	6	0	0	20	0	5	

Confusion Matrix K-NN

Εικόνα 8: K-NN

Το μοντέλο K-NN δείχνει ένα πιο κατανομημένο σύνολο προβλέψεων, με αρκετές τιμές εκτός διαγώνιου. Ενώ η διαγώνιος εξακολουθεί να έχει τις υψηλότερες τιμές, υπάρχουν διαφορές μεταξύ ορισμένων κλάσεων. Αυτό υποδηλώνει ότι ο K-NN δυσκολεύεται με κλάσεις που είναι παρόμοιες ή έχουν παρόμοια χαρακτηριστικά, οδηγώντας σε ένα λιγότερο ακριβές μοντέλο σε σύγκριση με το Random Forest.

Συμπεράσματα

Εν ολίγοις, το μοντέλο Random Forest είναι το καλύτερο για το συγκεκριμένο πρόβλημα και υπερτερεί σημαντικά έναντι των υπολοίπων. Αντίθετα το μοντέλο Decision Tree δεν είναι το καταλληλότερο ενώ τα υπόλοιπα δύο μοντέλα έχουν περιθώρια βελτίωσης με καλύτερες παραμέτρους ή καλύτερης ποιότητας δεδομένα.