



Εθνικό Μετσόβιο Πολυτεχνείο
Σχολή Ηλεκτρολόγων Μηχανικών και Μηχανικών Υπολογιστών
Τομέας Επικοινωνιών, Ηλεκτρονικής και Συστημάτων Πληροφορικής

Μπεϋζιανά Παίγνια σε Ασύρματα Δίκτυα Ομόσπονδης
Μάθησης με κακόβουλους χρήστες

Διπλωματική Εργασία

της

ΣΟΦΙΑΣ ΜΠΑΡΚΑΤΣΑ

Επιβλέπων: Συμεών Παπαβασιλείου
Καθηγητής Ε.Μ.Π.

Αθήνα, Ιούνιος 2024



Εθνικό Μετσόβιο Πολυτεχνείο
Σχολή Ηλεκτρολόγων Μηχανικών και Μηχανικών Υπολογιστών
Τομέας Επικοινωνιών, Ηλεκτρονικής και Συστημάτων Πληροφορικής

Μπεϋζιανά Παίγνια σε Ασύρματα Δίκτυα Ομόσπονδης Μάθησης με κακόβουλους χρήστες

Διπλωματική Εργασία

της

ΣΟΦΙΑΣ ΜΠΑΡΚΑΤΣΑ

Επιβλέπων: Συμεών Παπαβασιλείου
Καθηγητής Ε.Μ.Π.

Εγκρίθηκε από την τριμελή εξεταστική επιτροπή την 27η Ιουνίου 2024.

(Υπογραφή)

(Υπογραφή)

(Υπογραφή)

.....

Συμεών Παπαβασιλείου
Καθηγητής Ε.Μ.Π.

.....

Ελένη Στάη
Επίκουρη Καθηγήτρια Ε.Μ.Π.

.....

Γεώργιος Ματσόπουλος
Καθηγητής Ε.Μ.Π.

Αθήνα, Ιούνιος 2024



Εθνικό Μετσόβιο Πολυτεχνείο
Σχολή Ηλεκτρολόγων Μηχανικών και Μηχανικών Υπολογιστών
Τομέας Επικοινωνιών, Ηλεκτρονικής και Συστημάτων Πληροφορικής

(Υπογραφή)

.....

Σοφία Μπαρκάτσα

Διπλωματούχος Ηλεκτρολόγος Μηχανικός και Μηχανικός Υπολογιστών Ε.Μ.Π.

Copyright © Σοφία Μπαρκάτσα, 2024.

Με επιφύλαξη παντός δικαιώματος. All rights reserved.

Απαγορεύεται η αντιγραφή, αποθήκευση και διανομή της παρούσας εργασίας, εξ ολοκλήρου ή τμήματος αυτής, για εμπορικό σκοπό. Επιτρέπεται η ανατύπωση, αποθήκευση και διανομή για σκοπό μη κερδοσκοπικό, εκπαιδευτικής ή ερευνητικής φύσης, υπό την προϋπόθεση να αναφέρεται η πηγή προέλευσης και να διατηρείται το παρόν μήνυμα. Ερωτήματα που αφορούν τη χρήση της εργασίας για κερδοσκοπικό σκοπό πρέπει να απευθύνονται προς τον συγγραφέα.

Οι απόψεις και τα συμπεράσματα που περιέχονται σε αυτό το έγγραφο εκφράζουν τον συγγραφέα και δεν πρέπει να ερμηνευθεί ότι αντιπροσωπεύουν τις επίσημες θέσεις του Εθνικού Μετσόβιου Πολυτεχνείου.

Η Ομόσπονδη Μάθηση (Federated Learning - FL) είναι μια αποκεντρωμένη προσέγγιση στη μηχανική μάθηση που επιτρέπει στις συσκευές να εκπαιδεύουν από κοινού μοντέλα μηχανικής μάθησης, χωρίς να αποκαλύπτουν τα τοπικά τους δεδομένα και να διακινδυνεύουν την ιδιωτικότητά τους. Αντί να συλλέγει τα προσωπικά δεδομένα των χρηστών σε έναν κεντρικό διακομιστή, η Ομόσπονδη Μάθηση δημιουργεί ένα συγκεντρωτικό μοντέλο αθροίζοντας τα τοπικά μοντέλα μηχανικής μάθησης του εκπαιδεύουν οι συμμετέχοντες. Παρά τα πλεονεκτήματά της, όμως, η Ομόσπονδη Μάθηση καλείται να αντιμετωπίσει πολλαπλές προκλήσεις. Μία σημαντική απειλή της Ομόσπονδης μάθησης είναι οι επιθέσεις δηλητηριασμού, όπου κακόβουλοι συμμετέχοντες προσπαθούν να διαφθείρουν τη διαδικασία μάθησης χρησιμοποιώντας δηλητηριασμένα δεδομένα για την εκπαίδευση των τοπικών μοντέλων τους. Επιπλέον, όταν η Ομόσπονδη Μάθηση υλοποιείται μέσω ασύρματων δικτύων, είναι επιπλέον ευάλωτη σε διάφορες δικτυακές επιθέσεις, όπως παρεμβολές (jamming), που υποβαθμίζουν την ποιότητα της επικοινωνίας των συμμετεχόντων με τον διακομιστή. Η παρούσα διπλωματική εργασία προτείνει νέες προσεγγίσεις για την ενίσχυση της ασφάλειας της Ομόσπονδης Μάθησης έναντι πολλαπλών απειλών. Η μελέτη μας αφορά στην διεξαγωγή ενός Μπεϋζιανού Παιγνίου (Bayesian game) μεταξύ των χρηστών, με σκοπό την επιλογή της βέλτιστης ισχύος εκπομπής, για την εξουδετέρωση των επιθέσεων παρεμβολών. Ταυτόχρονα, υλοποιείται στον διακομιστή ένα συνεργατικό παίγνιο με την χρήση μιας τροποποιημένης εκδοχής της μετρικής Sharpley, για την αξιολόγηση της ποιότητας των τοπικών μοντέλων και την ανίχνευση επιθέσεων δηλητηριασμού. Τέλος, προτείνουμε την χρήση ενός νέου αλγόριθμου συνάθροισης (ContrAvg), ο οποίος αθροίζει τα βάρη των τοπικών μοντέλων με βάση την ποιότητά τους, μειώνοντας έτσι τις επιπτώσεις των επιθέσεων δηλητηριασμού. Η μελέτη μας είναι ανάμεσα στις λίγες εργασίες που αντιμετωπίζουν δύο διαφορετικούς τύπους επιθέσεων στην Ομόσπονδη Μάθηση, υπό ένα ενιαίο πλαίσιο, προστατεύοντας τόσο το δίκτυο ασύρματης επικοινωνίας όσο και την διαδικασία της Ομόσπονδης Μάθησης, από την απειλή κακόβουλων συμμετεχόντων.

Λέξεις Κλειδιά

Ομόσπονδη Μάθηση, Μπεϋζιανά Παιγνία, τιμή Sharpley, επιθέσεις Παρεμβολών, επιθέσεις δηλητηριασμού

Abstract

Federated Learning (FL) is a decentralized approach to machine learning that allows devices to collaboratively train models without exposing their local data, thus preserving privacy and security. Instead of collecting private data, from various users on a centralized server, FL constructs a globally shared model by only aggregating clients' locally computed machine learning models. Despite its advantages, Federated Learning faces numerous security challenges. One significant threat comes from insider attacks, where malicious clients may attempt to corrupt the learning process by using poisoned data to manipulate their model updates. Additionally, when FL is implemented over wireless networks, it is further susceptible to various network attacks, notably jamming, which can significantly disrupt the communication and model update process. Our work introduces novel approaches to enhance the security of Federated Learning against diverse threats. To the best of our knowledge, we are the first to formulate a Bayesian game in FL, where users optimize their transmission power to neutralize jamming attacks. Secondly, we utilize a cooperative game framework with a modified version of the Shapley value, to assess the contribution index of local models and detect poisoning attacks. Lastly, we introduce a new aggregation algorithm, ContrAvg, which aggregates local models based on their contribution indices, thus mitigating the effects of poisoning attacks. Our method is amongst the few that address simultaneously two different types of attacks in FL, which threaten both the communication network and the produced machine learning model, thus making significant strides in safeguarding FL systems from different perspectives.

Keywords

Federated Learning, Bayesian Games, Shapley Value, Jamming Attack, Poison Attack

Ευχαριστίες

Καθώς ολοκληρώνεται η φοίτησή μου στην σχολή Ηλεκτρολόγων Μηχανικών και Μηχανικών Υπολογιστών του Εθνικού Μετσόβιου Πολυτεχνείου, θα ήθελα να ευχαριστήσω θερμά τον καθηγητή κ. Συμεών Παπαβασιλείου, όχι μόνο για την επίβλεψη της παρούσας διπλωματικής εργασίας, αλλά και για την εμπιστοσύνη που μου έδειξε, την καθοδήγησή του και τις γνώσεις που μου μετέδωσε. Επιπλέον θα ήθελα να ευχαριστήσω την Διδάκτορα Μαρία Διαμαντή και τον υποψήφιο Διδάκτορα Παναγιώτη Χαρατσάρη, για την σύλληψη του θέματος της πτυχιακής μου εργασίας και την βοήθεια που μου προσέφεραν σε όλα τα στάδια της εκπόνησής της.

Το ακαδημαϊκό μου αυτό ταξίδι δεν θα ήταν επιτυχές χωρίς την συμπαράσταση των οικείων μου προσώπων και συμφοιτητών μου. Δεν μπορώ παρά να εκφράσω την μεγαλύτερη ευγνωμοσύνη στους γονείς μου, που με υποστήριζαν πάντοτε σε κάθε μου βήμα και είναι ο λόγος που έχω φτάσει έως εδώ. Στην αδερφή μου Ειρήνη, οφείλω ένα μεγάλο ευχαριστώ, για την υπομονή της, την αγάπη της και την συναισθηματική συμπαράσταση που μου προσέφερε από την πρώτη μέρα των σπουδών μου. Τέλος, στον Πέτρο, την Σμαράγδα και τους συμφοιτητές που έγιναν φίλοι καρδιακοί με την πάροδο των χρόνων, θέλω να εκφράσω την ευγνωμοσύνη μου για όλες τις στιγμές που μοιραστήκαμε, τις εργασίες που ολοκληρώσαμε και τις δυσκολίες που ξεπεράσαμε.

Σοφία Μπαρκάτσα

Ιούλιος 2024

Περιεχόμενα

Περίληψη	1
Abstract	3
Ευχαριστίες	5
Περιεχόμενα	8
Κατάλογος Σχημάτων	10
Κατάλογος πινάκων	11
1 Εισαγωγή	13
1.1 Κίνητρο	13
1.2 Συνεισφορά	13
1.3 Περίγραμμα της Διπλωματικής Εργασίας	14
2 Θεωρητικό Υπόβαθρο	17
2.1 Ομόσπονδη Μάθηση	17
2.1.1 Τα είδη δικτύων Ομόσπονδης Μάθησης	18
2.1.2 Τα είδη επιθέσεων που απειλούν την Ομόσπονδη Μάθηση	20
2.2 Τηλεπικοινωνιακό Μοντέλο NOMA	23
2.3 Θεωρία Παιγνίων	25
2.3.1 Παίγνια σε Κανονική Μορφή Αναπαράστασης	26
2.3.2 Ισορροπία Nash	26
2.3.3 Μπεϋζιανά Παίγνια	27
2.3.4 Τιμή Sharpley	27
3 Σχετική Βιβλιογραφία	29
4 Μοντελοποίηση Συστήματος	33
4.1 Μοντελοποίηση Ομόσπονδης Μάθησης	33

4.1.1	Εκπαίδευση τοπικών Μοντέλων	33
4.1.2	Υπολογισμός Συνεισφοράς	34
4.1.3	Contribution Averaging	37
4.1.4	Βελτίωση των αλγορίθμων FedAvg και ContrAvg	37
4.2	Τηλεπικοινωνιακό Μοντέλο	38
4.2.1	Υπολογισμός Θεμελιωδών Μεγεθών	38
4.2.2	Κατώφλι Παρεμβολών	38
4.3	Διεξαγωγή Μπεϋζιανού Παιγνίου	39
4.3.1	Η συνάρτηση χρησιμότητας	39
4.3.2	Μπεϋζιανό Παίγνιο	40
4.3.3	Υπολογισμός Πιθανότητας Ιδιοτελούς χρήστη	41
5	Διεξαγωγή Πειραμάτων	43
5.1	Επιλογή τιμών σταθερών μεταβλητών	43
5.2	Μοντέλα Μηχανικής Μάθησης και Παράμετροι εκπαίδευσης	43
5.2.1	MNIST Dataset	44
5.2.2	Συνελκτικό Νευρωνικό Δίκτυο	44
5.2.3	Παράμετροι Εκπαίδευσης	45
5.3	Τα σενάρια που μελετήθηκαν	45
5.4	Τοπολογία και αρχικοποίηση χρηστών στον χώρο	46
5.5	Αρχικοποίηση Μεταβλητών	48
6	Αποτελέσματα	49
6.1	Εντοπισμός Κακόβουλου Χρήστη	49
6.1.1	Η μετρική a_i σε IID δεδομένα	50
6.1.2	Η μετρική a_i σε NON IID δεδομένα	51
6.2	Σύγκλιση Παιγνίου	52
6.2.1	Αρχικοποιήσεις χωρίς κακόβουλους χρήστες	52
6.2.2	Αρχικοποιήσεις με κακόβουλο χρήστη σε IID δεδομένα	53
6.2.3	Αρχικοποιήσεις με κακόβουλο χρήστη σε NON IID δεδομένα	57
6.3	Η απόδοση του Μπεϋζιανού Παιγνίου	59
6.3.1	Η μέση ισχύς μετάδοσης	59
6.3.2	Ο μέσος χρόνος που απαιτείται	61
6.3.3	Η μέση ενέργεια που απαιτείται	62
6.4	Το ποσοστό επιτυχίας του Συγκεντρωτικού Μοντέλου	64
6.4.1	Μελέτη σε IID δεδομένα	64
6.4.2	Μελέτη σε NON IID δεδομένα	67
7	Επίλογος και Μελλοντικές Επεκτάσεις	71
	Βιβλιογραφία	73
	Γλωσσάριο	77
	Παράρτημα	79

Κατάλογος σχημάτων

2.1	Η παραδοσιακή αρχιτεκτονική των μοντέλων μηχανικής μάθησης, όπου τα δεδομένα των συσκευών αποστέλλονται στον διακομιστή. Το σύστημα αυτό αποκαλείται κεντροποιημένο (centralized).	17
2.2	Ένα δίκτυο Ομόσπονδης Μάθησης. Οι χρήστες μεταδίδουν στον διακομιστή το τοπικό μοντέλο που εκπαιδύσαν, αντί για τα δεδομένα τους	18
2.3	Τα τρία είδη Ομόσπονδης Μάθησης, με κριτήριο τα δεδομένα που έχουν οι χρήστες A και B	20
2.4	Τα διαφορετικά είδη επιθέσεων στην Ομόσπονδη Μάθηση	22
4.1	Το μοντέλο που προτείνει η παρούσα εργασία. Ο κακόβουλος χρήστης αποτυπώνεται με κόκκινο χρώμα και πραγματοποιεί τόσο επιθέσεις παρεμβολών (2) όσο και δηλητηριασμού δεδομένων (1). Αφού ο διακομιστής λάβει τα τοπικά μοντέλα, υπολογίζει τον βαθμό συνεισφοράς (a_i) τους (3) και στην συνέχεια παράγεται το συγκεντρωτικό μοντέλο του γύρου (4) με σεβασμό στον βαθμό a_i .	34
4.2	Γραφική αναπαράσταση των πιθανοτήτων με $\mu = \pm 0.03, a = 15$	42
5.1	Μερικές από τις εικόνες του συνόλου MNIST [51]	44
5.2	Διαφορετικές αρχικοποιήσεις των χρηστών στον χώρο	47
6.1	Βαθμός συνεισφοράς a_i σε IID δεδομένα	50
6.2	Βαθμός συνεισφοράς a_i σε NON IID δεδομένα	51
6.3	Η σύγκλιση της ισχύος μετάδοσης p_i σε σενάρια με 5 ιδιοτελείς χρήστες	52
6.4	Η σύγκλιση της ισχύος μετάδοσης p_i σε σενάρια όπου ο χρήστης 1 είναι κακόβουλος, πραγματοποιεί επιθέσεις παρεμβολών, και το σύνολο δεδομένων των χρηστών έχει IID μορφή	54
6.5	Η ισχύς μετάδοσης των χρηστών στο βήμα 5 κατά τις πρώτες 10 εποχές, όπου ο χρήστης 0 είναι κακόβουλος και τα δεδομένα των χρηστών έχουν IID μορφή	55
6.6	Η σύγκλιση της ισχύος μετάδοσης p_i σε σενάρια όπου ο χρήστης 1 είναι κακόβουλος, πραγματοποιεί επιθέσεις παρεμβολών, και το σύνολο δεδομένων των χρηστών έχει NON IID μορφή	57

6.7	Η μέση ισχύς του κάθε χρήστη, ανάλογα αν τα δεδομένα είναι IID ή NON IID και το κατά πόσο υπάρχει κακόβουλος συμμετέχοντες που πραγματοποιεί επιθέσεις παρεμβολών και δηλητηριασμού	60
6.8	Ο μέσος χρόνος (σε seconds) που απαιτείται για την μετάδοση του κάθε χρήστη, ανάλογα αν τα δεδομένα είναι IID ή NON IID και το κατά πόσο υπάρχει κακόβουλος συμμετέχοντες που πραγματοποιεί επιθέσεις παρεμβολών και δηλητηριασμού	62
6.9	Η μέση ενέργεια που απαιτείται για την μετάδοση του κάθε χρήστη, ανάλογα αν τα δεδομένα είναι IID ή NON IID και το κατά πόσο υπάρχει κακόβουλος συμμετέχοντες που πραγματοποιεί επιθέσεις παρεμβολών και δηλητηριασμού	63
6.10	Ποσοστό επιτυχίας (Accuracy) 4 διαφορετικών σεναρίων, σε IID δεδομένα	65
6.11	Οι απώλειες (Loss) 4 διαφορετικών σεναρίων, σε IID δεδομένα	66
6.12	Ποσοστό επιτυχίας (Accuracy) 4 διαφορετικών σεναρίων, σε NON IID δεδομένα	67
6.13	Οι απώλειες (Loss) 4 διαφορετικών σεναρίων, σε IID δεδομένα	69

Κατάλογος πινάκων

4.1	Επεξήγηση συμβολισμών	40
4.2	Πίνακας τιμών εκτίμησης πιθανοτήτων με $\mu = \pm 0.03, a = 15$	42
5.1	Αρχικοποίηση Σταθερών	43
5.2	Πίνακας Αρχικοποίησης Μεταβλητών	48
6.1	Τιμές του μέσου βαθμού a_i	51
6.2	Αντιστοίχιση των ταξινομημένων κερδών καναλιού $G_0 \leq G_1 \leq G_2 \leq G_3 \leq G_4$ στους συμμετέχοντες του παιγνίου	59
6.3	Τα συγκεντρωτικά αποτελέσματα τεσσάρων διαφορετικών σεναρίων σε IID δεδομένα, στο τέλος της εποχής 50	66
6.4	Τα συγκεντρωτικά αποτελέσματα τεσσάρων διαφορετικών σεναρίων σε NON IID δεδομένα, στο τέλος της εποχής 50	68

1.1 Κίνητρο

Στην σύγχρονη εποχή τα μοντέλα μηχανικής μάθησης χρησιμοποιούνται από ένα διαρκώς αυξανόμενο πλήθος εφαρμογών, με αποτέλεσμα να απαιτείται μία υπέρογκη ποσότητα δεδομένων, όπως φωτογραφίες ή κείμενο, για την εκπαίδευσή τους. Ωστόσο, οι χρήστες της εκάστοτε εφαρμογής, είτε αυτοί είναι ιδιώτες είτε είναι επιχειρήσεις, δεν είναι διατεθειμένοι να θυσιάσουν τα ευαίσθητα δεδομένα τους στον βωμό της τεχνητής νοημοσύνης. Το πρόβλημα αυτό καλείται να αντιμετωπίσει η Ομόσπονδη Μάθηση (Federated Learning), στην οποία οι χρήστες συνδράμουν για την εκπαίδευση νευρωνικών δικτύων, χωρίς να κοινοποιήσουν τα προσωπικά τους δεδομένα.

Στην ομόσπονδη μάθηση οι συμμετέχοντες αναλαμβάνουν σε κάθε εποχή να εκπαιδεύσουν τοπικά ένα μοντέλο μηχανικής μάθησης, χρησιμοποιώντας τα προσωπικά τους δεδομένα. Έπειτα αποστέλλουν σε έναν κεντρικό εξυπηρετητή / διακομιστή (server) τα βάρη των νευρωνικών δικτύων που εκπαίδευσαν, χωρίς να μοιράζονται καμία ευαίσθητη πληροφορία. Τα βάρη των διαφόρων χρηστών αθροίζονται και το μοντέλο που προκύπτει αποστέλλεται πίσω στους συμμετέχοντες για περαιτέρω εκπαίδευση σε επόμενη εποχή.

Τα δίκτυα Ομόσπονδης Μάθησης είναι δυστυχώς ευάλωτα σε ένα πλήθος επιθέσεων από κακόβουλους χρήστες, ειδικά όταν η διαδικασία εκπαίδευσης πραγματοποιείται μέσω ενός ασύρματου καναλιού επικοινωνίας. Είναι πιθανό οι συμμετέχοντες στην Ομόσπονδη Μάθηση να επιθυμούν να δηλητηριάσουν νευρωνικό δίκτυο, στέλνοντας λανθασμένα βάρη στον κεντρικό εξυπηρετητή. Από την άλλη, η ασύρματη επικοινωνία των χρηστών με τον διακομιστή μπορεί να επιδεινωθεί ή να διακοπεί από σκόπιμες παρεμβολές κακόβουλων χρηστών.

Η διεθνής βιβλιογραφία αν και κάνει λόγο για διάφορα είδη επιθέσεων, σπάνια μελετάει την ταυτόχρονη αντιμετώπιση δύο διαφορετικών επιθέσεων, με την χρήση ενός ενιαίου μηχανισμού. Ένα τέτοιο σενάριο είναι πλήρως ρεαλιστικό, καθώς αν ένας από τους συμμετέχοντες είναι κακόβουλος, δεν έχει λόγο να περιοριστεί -για παράδειγμα- στην εκπομπή δηλητηριασμένων μοντέλων, όταν μπορεί συγχρόνως, με επιθέσεις παρεμβολών, να αποκλείσει χρήστες από την διαδικασία της Ομόσπονδης Μάθησης.

1.2 Συνεισφορά

Με αφορμή την παραπάνω έλλειψη στην διεθνή βιβλιογραφία, η παρούσα διπλωματική αναλαμβάνει την αντιμετώπιση δύο ταυτόχρονων επιθέσεων, με την χρήση ενός ενιαίου αμυντικού μηχανισμού. Για τον

σκοπό αυτό μοντελοποιούμε ένα ασύρματο δίκτυο Ομόσπονδης Μάθησης, στο οποίο ένας κακόβουλος συμμετέχοντας πραγματοποιεί επιθέσεις παρεμβολών και δηλητηριασμού, ώστε να μειώσει την αποτελεσματικότητα των συγκεντρωτικών μοντέλων Μηχανικής Μάθησης που παράγονται στον εξυπηρετητή.

Αρχικά ο κακόβουλος χρήστης εκπαιδεύει ένα τοπικό μοντέλο χρησιμοποιώντας ψευδή δεδομένα, με αποτέλεσμα το μοντέλο αυτό να επιτυγχάνει -σχεδόν- μηδενικά ποσοστά επιτυχίας. Πιο συγκεκριμένα, τα ψευδή αυτά δεδομένα είναι ασπρόμαυρες εικόνες (MNIST dataset), με χειρόγραφα ψηφία από το 0 έως το 9, στις οποίες ο επιτιθέμενος έχει αλλάξει την ετικέτα. Έτσι οι εικόνες με τον αριθμό 1 φαίνεται να έχουν τον αριθμό 2 και ούτω καθεξής. Αφού λοιπόν αλλάξουν οι ετικέτες του συνόλου δεδομένων, ο κακόβουλος χρήστης εκπαιδεύει το τοπικό του μοντέλο πάνω στο τροποποιημένο σύνολο δεδομένων. Η επίθεση αυτή ονομάζεται επίθεση δηλητηριασμού, αφού το μοντέλο που εκπέμπει ο κακόβουλος χρήστης στον εξυπηρετητή, θα επηρεάσει αρνητικά τον σχηματισμό του νέου συγκεντρωτικού μοντέλου.

Συγχρόνως με την επίθεση δηλητηριασμού, ο κακόβουλος χρήστης φροντίζει να μεταδώσει το μοντέλο του με περίσσεια ισχύ, αποσκοπώντας να προκαλέσει παρεμβολές στους κοντινούς χρήστες που μεταδίδουν ταυτόχρονα τα δικά τους τοπικά μοντέλα. Αν οι παρεμβολές αυτές είναι αρκετά ισχυρές, τότε τα σήματα των υπολοίπων συμμετεχόντων φτάνουν στον εξυπηρετητή με μειωμένο σηματοθορυβικό λόγο, με αποτέλεσμα αυτός να μην μπορέσει να τα αποκωδικοποιήσει ορθά. Αν αυτό συμβεί, τα τοπικά μοντέλα των χρηστών που δέχθηκαν παρεμβολές δεν προσμετρώνται στην διαδικασία της Ομόσπονδης Μάθησης. Έτσι το νέο συγκεντρωτικό μοντέλο δεν επηρεάζεται από χρήσιμα τοπικά μοντέλα και αυξάνεται ο βαθμός επιρροής του δηλητηριασμένου μοντέλου του κακόβουλου χρήστη.

Για την αντιμετώπιση των δύο αυτών προβλημάτων, διεξάγεται ένα Μπεϋζιανό Παιγνίο μεταξύ των συμμετεχόντων, στο οποίο ρυθμίζεται η ισχύς μετάδοσης των χρηστών. Ακόμα και αν ένας κακόβουλος συμμετέχοντας επιχειρήσει επίθεση παρεμβολών, οι υπόλοιποι χρήστες μπορούν τότε να ρυθμίσουν την ισχύ εκπομπής τους στον επόμενο γύρο, ώστε να αποφύγουν τις επιθέσεις αυτές. Επίσης, όταν ο εξυπηρετητής λάβει τα τοπικά μοντέλα των χρηστών, τα αξιολογεί με την χρήση της τιμής Sharpley, σε ένα συνεργατικό παίγνιο, με αποτέλεσμα να μπορεί να εντοπίσει τις επιθέσεις δηλητηριασμού και σε δεύτερο βήμα να τις περιορίσει με τον αλγόριθμο ασφαλούς συνάθροισης δεδομένων (ContrAvg) που προτείνουμε.

Μέσω της διεξαγωγής προσομοιώσεων αποδεικνύεται η αποδοτικότητα του αμυντικού μηχανισμού που μελετάμε, τόσο στην αποτροπή επιθέσεων παρεμβολών, όσο και στον περιορισμό του αντίκτυπου των επιθέσεων δηλητηριασμού.

1.3 Περίγραμμα της Διπλωματικής Εργασίας

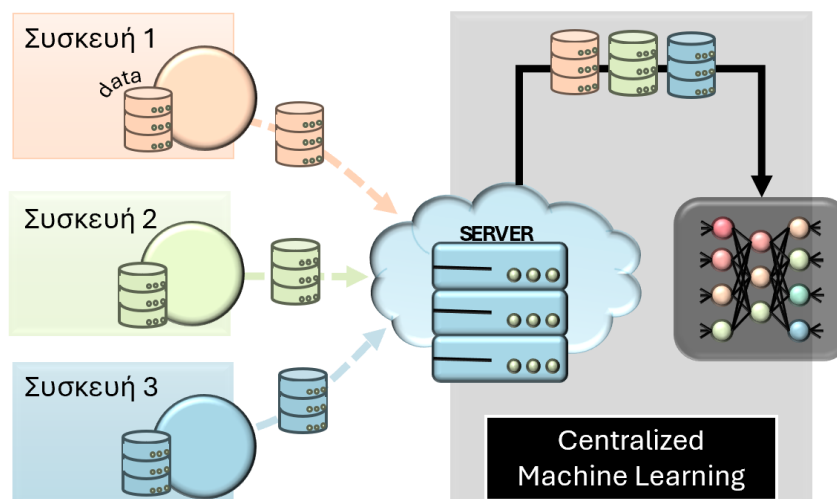
Η παρούσα διπλωματική εργασία παρουσιάζει την προτεινόμενη μέθοδο ανίχνευσης και αντιμετώπισης επιθέσεων παρεμβολών και δηλητηριασμού στην Ομόσπονδη Μάθηση στις ακόλουθες 6 ενότητες:

1. Στην ενότητα 2 παρουσιάζεται το θεωρητικό υπόβαθρο στο οποίο βασίζεται η διπλωματική εργασία, με έμφαση στην κατανόηση του Μηχανισμού της Ομόσπονδης Μάθησης και τις επιθέσεις που την απειλούν. Ακολουθεί επίσης μία πρώτη αναφορά στο τηλεπικοινωνιακό μοντέλο NOMA που χρησιμοποιείται για την ασύρματη επικοινωνία των χρηστών με τον διακομιστή, καθώς και σε βασικές έννοιες από την θεωρία παιγνίων, που αποτελούν την βάση του μηχανισμού αντιμετώπισης των επιθέσεων.
2. Στην ενότητα 3 παρουσιάζεται η διεθνής βιβλιογραφία σχετικά με την Ομόσπονδη Μάθηση, με έμφαση στην συμβολή της θεωρίας παιγνίων στην αντιμετώπιση κάποιων σημαντικών προβλημάτων.

3. Στην ενότητα 4 να παρουσιάζεται αναλυτικά η μοντελοποίηση του δικτύου Ομόσπονδης Μάθησης, οι παραδοχές που πραγματοποιήθηκαν καθώς και τα στάδια που απαιτούνται για τον εντοπισμό και τον περιορισμό του αντίκτυπου των επιθέσεων παρεμβολών και δηλητηριασμού. Μοντελοποιείται το Μπεϋζιανό Παιγνίο μεταξύ των συμμετεχόντων, παρουσιάζεται μια τροποποιημένη εκδοχή της τιμής Sharpley και περιγράφονται οι αλγόριθμοι FedAvg και ContrAvg που είναι υπεύθυνοι για τον σχηματισμό του συγκεντρωτικού μοντέλου.
4. Στην ενότητα 5 αναλύονται οι τεχνικές λεπτομέρειες της διεξαγωγής των προσομοιώσεων, όπως οι τιμές των σταθερών που επιλέχθηκαν, οι αρχικοποιήσεις των χρηστών στον χώρο, το πλήθος των πειραματικών διατάξεων και τα χαρακτηριστικά του μοντέλων μηχανικής μάθησης.
5. Στην ενότητα 6 παρουσιάζονται τα αποτελέσματα των προσομοιώσεων, με ξεχωριστή έμφαση σε κάθε τμήμα της μοντελοποίησης. Επαληθεύεται η δυνατότητα του εντοπισμού ενός κακόβουλου χρήστη, παρουσιάζεται η σύγκλιση του Μπεϋζιανού Παιγνίου, πραγματοποιείται σύγκριση των αλγορίθμων FedAvg και ContrAvg ως προς το ποσοστό επιτυχίας και τις απώλειες και παρουσιάζονται διάφορες μετρικές των παιγνίων.
6. Τέλος, στην ενότητα 7, πραγματοποιείται η σύνοψη της διπλωματικής εργασίας και αναφέρονται πιθανές μελλοντικές επεκτάσεις του συστήματος που παρουσιάστηκε.

2.1 Ομόσπονδη Μάθηση

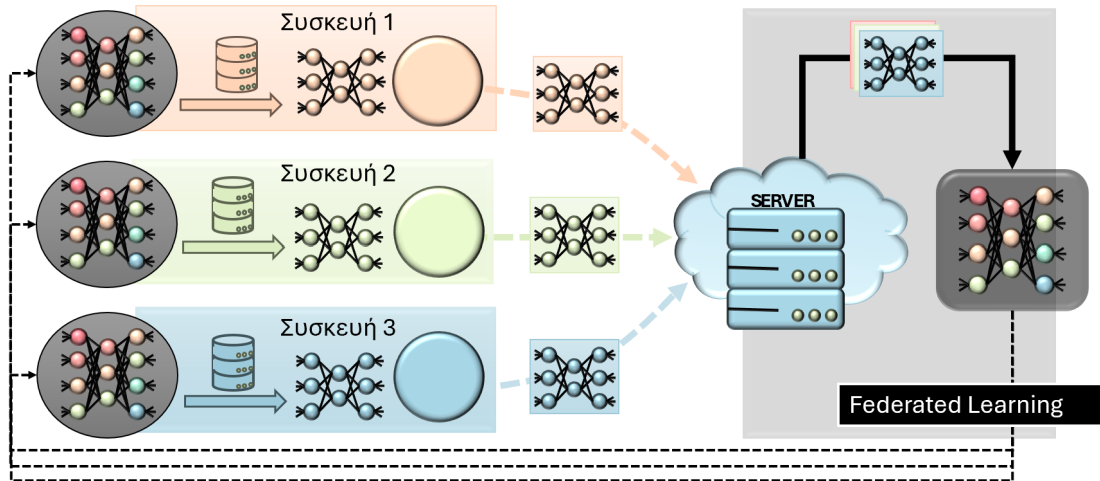
Στην σύγχρονη εποχή, οι ραγδαίες τεχνολογικές εξελίξεις έχουν αυξήσει σημαντικά την υπολογιστική ισχύ και τις δυνατότητες των φορητών συσκευών. Πλέον οι συσκευές άκρων (Edge Devices), όπως το κινητό τηλέφωνο και τα έξυπνα ρολόγια συλλέγουν και διαχειρίζονται μια πληθώρα διαφορετικών δεδομένων, που βρίσκονται σε μορφή εικόνων, κειμένου ή και ήχου. Ταυτόχρονα, η ανάγκη των καταναλωτών για ολοένα πιο έξυπνες και εξατομικευμένες εφαρμογές, έστρεψε την επιστημονική κοινότητα στην χρήση μοντέλων μηχανικής μάθησης, τα οποία απαιτούν ένα υπέρογκο πλήθος δεδομένων για την εκπαίδευσή τους.



Σχήμα 2.1: Η παραδοσιακή αρχιτεκτονική των μοντέλων μηχανικής μάθησης, όπου τα δεδομένα των συσκευών αποστέλλονται στον διακομιστή. Το σύστημα αυτό αποκαλείται κεντροποιημένο (centralized).

Τα παραδοσιακά μοντέλα μηχανικής μάθησης (centralized machine learning) απαιτούν την συλλογή, αποθήκευση και επεξεργασία ενός πλήθους πληροφοριών, οι οποίες προέρχονται συχνά από τους χρήστες του διαδικτύου. Η ανάγκη των χρηστών αυτών να προστατέψουν τα προσωπικά τους δεδομένα, οδήγησε τους ερευνητές της Google το 2016 στην δημιουργία ενός δικτύου Ομόσπονδης Μάθησης (Federated Learning), το οποίο επιτρέπει την αποκεντροποιημένη εκπαίδευση νευρωνικών δικτύων [1]. Χάρη στην Ομόσπονδη Μάθηση τα δεδομένα χρηστών, αντί να συλλέγονται σε έναν κεντρικό διακομιστή (server),

παραμένουν στις συσκευές των χρηστών και δεν κοινοποιούνται περαιτέρω. Οι συσκευές άκρων αναλαμβάνουν σε κάθε εποχή να εκπαιδεύσουν τοπικά ένα μοντέλο μηχανικής μάθησης, χρησιμοποιώντας ως δεδομένα, αυτά που έχουν συλλέξει από τον χρήστη τους. Έπειτα αποστέλλουν στον κεντρικό εξυπηρετητή (server) αποκλειστικά τα βάρη των νευρωνικών δικτύων που εκπαιδεύσαν. Ο διακομιστής, αφού λάβει τα μοντέλα όλων των συμμετεχόντων στην Ομόσπονδη Μάθηση, αθροίζει τις παραμέτρους και επιστρέφει στους χρήστες το **συγκεντρωτικό μοντέλο** (aggregated model) για περαιτέρω εκπαίδευση στην επόμενη εποχή.



Σχήμα 2.2: Ένα δίκτυο Ομόσπονδης Μάθησης. Οι χρήστες μεταδίδουν στον διακομιστή το τοπικό μοντέλο που εκπαιδεύσαν, αντί για τα δεδομένα τους

Η Ομόσπονδη Μάθηση έχει πλέον αξιοποιηθεί από πολλές διαφορετικές εφαρμογές, με χαρακτηριστικό παράδειγμα το GBoard, το εικονικό πληκτρολόγιο της Google που είναι προ-εγκατεστημένο στα περισσότερα Android κινητά τηλέφωνα. Το GBoard προκειμένου να βελτιώσει τις προτάσεις του, εκπαιδεύει το γλωσσικό του μοντέλο (language model) με βάση το κείμενο που πληκτρολογούν οι χρήστες του [2]. Η χρήση Ομόσπονδης Μάθησης καθιστά την εκπαίδευση αυτή εφικτή, χωρίς να χρειαστεί η κοινοποίηση ευαίσθητων δεδομένων που μπορεί να περιλαμβάνουν κωδικούς, διευθύνσεις ή ακόμη και αριθμούς πιστωτικών καρτών.

2.1.1 Τα είδη δικτύων Ομόσπονδης Μάθησης

Ένας διαχωρισμός των δικτύων Ομόσπονδης Μάθησης μπορεί να πραγματοποιηθεί ανάλογα με τα δεδομένα που έχουν οι χρήστες. Οι τρεις βασικές κατηγορίες Ομόσπονδης Μάθησης είναι η Οριζόντια (Horizontal Federated Learning), η Κάθετη (Vertical Federated Learning) και η Ομόσπονδη Μεταφορά Μάθησης (Federated Transfer Learning) [3]. Έστω ότι ο πίνακας D_i περιλαμβάνει τα δεδομένα που έχει ο χρήστης i . Κάθε γραμμή του πίνακα αντιπροσωπεύει ένα δείγμα, για παράδειγμα μια εικόνα, και κάθε στήλη περιγράφει κάποιο χαρακτηριστικό του δείγματος αυτού. Σε κάποια σύνολα δεδομένων, όπως στο MNIST που χρησιμοποιεί η παρούσα εργασία, περιλαμβάνονται και ετικέτες ή επισημειώσεις (labels) που δείχνουν σε ποιά κλάση ανήκει το εκάστοτε δείγμα. Ορίζεται ως X το σύνολο των χαρακτηριστικών (features), ως Y ο χώρος των ετικετών και ως I ο χώρος στον οποίο ανήκει το αναγνωριστικό των δειγμάτων (sample ID).

Τα χαρακτηριστικά X , οι ετικέτες Y και τα αναγνωριστικά των δειγμάτων I συνθέτουν ένα ολοκληρωμένο dataset (I, X, Y) . Οι χώροι των χαρακτηριστικών και των δειγμάτων των δεδομένων που έχουν οι

συμμετέχοντες στην Ομόσπονδη Μάθηση μπορεί όμως να μην είναι ταυτόσημοι. Οι τρεις βασικές κατηγορίες Οριζόντιας, Κάθετης Ομόσπονδης Μάθησης και Ομόσπονδης Μεταφοράς Μάθησης προκύπτουν με γνώμονα τον τρόπο κατανομής των δεδομένων στους χρήστες.

1. Οριζόντια Ομόσπονδη Μάθηση (Horizondal Federated Learning)

Στην Οριζόντια Ομόσπονδη Μάθηση οι συμμετέχοντες έχουν κοινό χώρο ετικετών, αλλά διαφορετικά δείγματα. Αυτό είναι και το είδος Ομόσπονδης Μάθησης που χρησιμοποιείται στην παρούσα εργασία. Οι χρήστες έχουν ο καθένας τους διαφορετικές εικόνες (δηλαδή διαφορετικά δείγματα), αλλά όλες οι εικόνες μοιράζονται το ίδιο σύνολο χαρακτηριστικών και αντιστοιχούν στο ίδιο σύνολο ετικετών. Συμβολικά, η Οριζόντια Ομόσπονδη Μάθηση αναπαριστάται ως εξής:

$$X_i = X_j, Y_i = Y_j, I_i \neq I_j, \forall D_i, D_j, i \neq j \quad (2.1)$$

2. Κάθετη Ομόσπονδη Μάθηση (Vertical Federated Learning)

Στην Κάθετη Ομόσπονδη Μάθηση οι συμμετέχοντες έχουν τα ίδια δείγματα, αλλά ο καθένας γνωρίζει διαφορετικά χαρακτηριστικά των δειγμάτων αυτών. Παράδειγμα της συγκεκριμένης κατηγορίας Ομόσπονδης Μάθησης θα μπορούσε να αποτελέσει η συνεργασία τραπεζών με εταιρίες τιμολόγησης, ώστε να εκπαιδεύσουν μοντέλα χρηματοοικονομικού κινδύνου για τους εταιρικούς τους πελάτες [4]. Η μεν τράπεζα έχει πρόσβαση στις καταθέσεις, τις αποταμιεύσεις και τα δάνεια ενός πελάτη, ενώ οι εταιρίες τιμολόγησης γνωρίζουν διαφορετικές πληροφορίες για τους πελάτες αυτούς, όπως τις συναλλαγές τους. Συμβολικά, η Κάθετη Ομόσπονδη Μάθηση αναπαριστάται ως εξής:

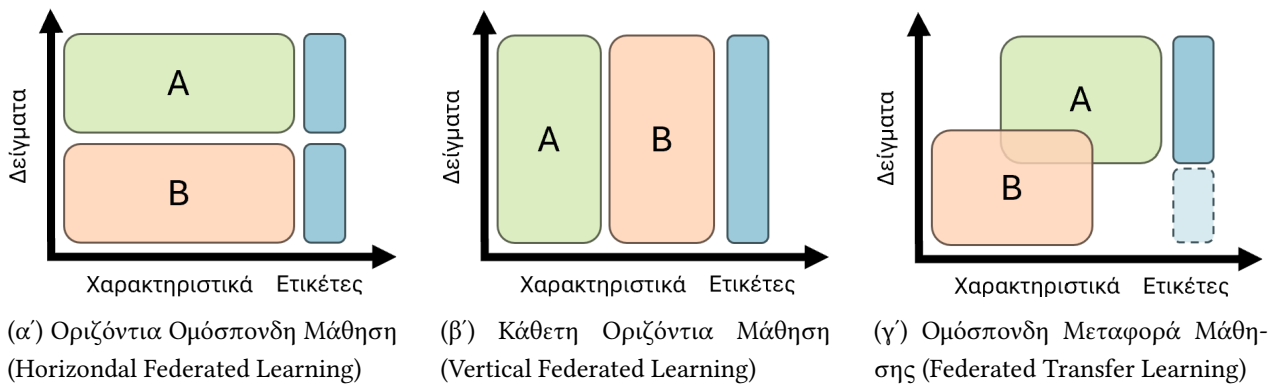
$$X_i \neq X_j, Y_i \neq Y_j, I_i = I_j, \forall D_i, D_j, i \neq j \quad (2.2)$$

3. Ομόσπονδη Μεταφορά Μάθησης (Federated Transfer Learning)

Στην Ομόσπονδη Μεταφορά Μάθησης οι συμμετέχοντες έχουν τελείως διαφορετικά σύνολα δεδομένων, όχι μόνο ως προς τα δείγματα αλλά και ως προς τα χαρακτηριστικά των δειγμάτων τους. Μια εφαρμογή της συγκεκριμένης κατηγορίας είναι η εκπαίδευση ενός μοντέλου σε ηλεκτροεγκεφαλογραφικά σήματα (ECG signals), τα οποία επειδή συνήθως προέρχονται από διαφορετικά μηχανήματα, έχουν διαφορετικά χαρακτηριστικά. Αυτό σημαίνει πως από δείγμα σε δείγμα παρατηρείται -για παράδειγμα- διαφορετικός αριθμός ηλεκτροδίων και διαφορετικός ρυθμός λήψης των σημάτων, ανάλογα με το μέρος και τον εξοπλισμό που χρησιμοποιήθηκε [5]. Συμβολικά, η Οριζόντια Ομόσπονδη Μάθηση αναπαριστάται ως εξής:

$$X_i \neq X_j, Y_i \neq Y_j, I_i \neq I_j, \forall D_i, D_j, i \neq j \quad (2.3)$$

Εκτός από τον διαχωρισμό των δικτύων Ομόσπονδης Μάθησης ανάλογα με το πλήθος των δεδομένων τους, μπορεί να υπάρξει και διαχωρισμός με βάση το πλήθος (και το είδος) των συμμετεχόντων σε αυτό [6]. Ενώ η Google όταν δημιούργησε την Ομόσπονδη Μάθηση, ο στόχος ήταν η συμμετοχή ενός μεγάλου αριθμού συμμετεχόντων. Για παράδειγμα, η εφαρμογή της Google GBoard έχει περίπου 1 δισεκατομμύριο λήψεις [7], με αποτέλεσμα το δίκτυο Ομόσπονδης μάθησης, από το οποίο εκπαιδεύονται οι προτάσεις του πληκτρολογίου, να μπορεί να επιστρατεύσει έναν υπέρογκο αριθμό συσκευών. Οι συσκευές αυτές -ωστόσο- έχουν ελάχιστα δεδομένα, μόνο όσα συνέλεξαν από τον χρήστη τους. Υπάρχουν όμως σενάρια, όπου η Ομόσπονδη Μάθηση χρησιμοποιείται από ένα μικρό πλήθος εταιριών, οι οποίες όμως έχουν στην



Σχήμα 2.3: Τα τρία είδη Ομόσπονδης Μάθησης, με κριτήριο τα δεδομένα που έχουν οι χρήστες A και B

διάθεσή τους μια πληθώρα δεδομένων, για την εκπαίδευση των τοπικών μοντέλων. Έτσι προκύπτουν τα δύο ακόλουθα είδη Ομόσπονδης Μάθησης [7], [8]

1. Cross-Device Ομόσπονδη Μάθηση

είναι η κατηγορία όπου το πλήθος των χρηστών είναι μεγάλο (από μερικές χιλιάδες έως δισεκατομμύρια), αλλά μόνο ένα μικρό μέρος των χρηστών αυτών επιλέγεται (τυχαία) να συμμετέχει σε έναν γύρο της Ομόσπονδης Μάθησης.

2. Cross-Silo Ομόσπονδη Μάθηση

ή αλλιώς επιχειρησιακή Ομόσπονδη Μάθηση [6], είναι η κατηγορία Ομόσπονδης Μάθησης όπου το πλήθος συμμετεχόντων είναι μικρό, από 2 έως 100 εταιρίες, ωστόσο η κάθε εταιρία έχει ένα μεγάλο πλήθος από δεδομένα.

Η παρούσα διπλωματική εργασία θα μελετήσει δίκτυα Οριζόντιας Cross-Silo Ομόσπονδης Μάθησης. Ο λόγος μιας τέτοιας επιλογής έγκειται αφενός στον περιορισμό που μας υποβάλει το τηλεπικοινωνιακό πλαίσιο NOMA, που δεν επιτρέπει την ταυτόχρονη μετάδοση μεγάλου πλήθους χρηστών. Αφετέρου, λόγω περιορισμένων υπολογιστικών πόρων, δεν είναι εφικτή η υποστήριξη μεγάλου αριθμού συμμετεχόντων, καθώς με την αύξηση των χρηστών της Ομόσπονδης Μάθησης, αυξάνεται ο αριθμός τοπικών μοντέλων που πρέπει να εκπαιδευτούν και να αξιολογηθούν, όπως αναλύεται στις επόμενες ενότητες.

2.1.2 Τα είδη επιθέσεων που απειλούν την Ομόσπονδη Μάθηση

Όπως προαναφέρθηκε, η εκπαίδευση του συγκεντρωτικού μοντέλου στην Ομόσπονδη Μάθηση εξαρτάται από τα τοπικά μοντέλα που μεταδίδουν οι χρήστες στον διακομιστή. Το γεγονός αυτό καθιστά ένα δίκτυο Ομόσπονδης Μάθησης επιρρεπές σε ένα ευρύ πλήθος επιθέσεων. Όταν όμως το δίκτυο Ομόσπονδης Μάθησης είναι ασύρματο, τότε είναι επιπλέον ευάλωτο σε ασύρματες δικτυακές επιθέσεις [9]. Κάποιες από αυτές είναι:

1. Επίθεση Λαθρακρόασης (Eavesdropping Attack)

Η επίθεση αυτή πραγματοποιείται όταν ένας κακόβουλος χρήστης κρυφακούει το ασύρματο κανάλι επικοινωνίας και καταγράφει τις μεταδόσεις χρηστών. Έτσι είναι πιθανό να αποκτήσει πρόσβαση σε ευαίσθητες πληροφορίες, τις οποίες θα αξιοποιήσει στην συνέχεια για άλλες σκοπιμότητες. Η επίθεση αυτή, αν και είναι πολύ εύκολο να πραγματοποιηθεί, μπορεί να αντιμετωπιστεί εύκολα με την κρυπτογράφηση των δεδομένων των χρηστών. Σε περίπτωση που χρησιμοποιούνται αλγόριθμοι κρυπτογράφησης κατά την επικοινωνία ενός συμμετέχοντα με τον Διακομιστή, τότε ακόμα

και αν ένας κακόβουλος χρήστης κρυφακούσει το μήνυμα που μεταδίδεται, δεν θα καταφέρει να αποκρυπτογραφήσει το περιεχόμενό του και -συνεπώς- δεν θα λάβει οποιαδήποτε ευαίσθητη πληροφορία.

2. Επίθεση Επανάληψης (Replay Attack)

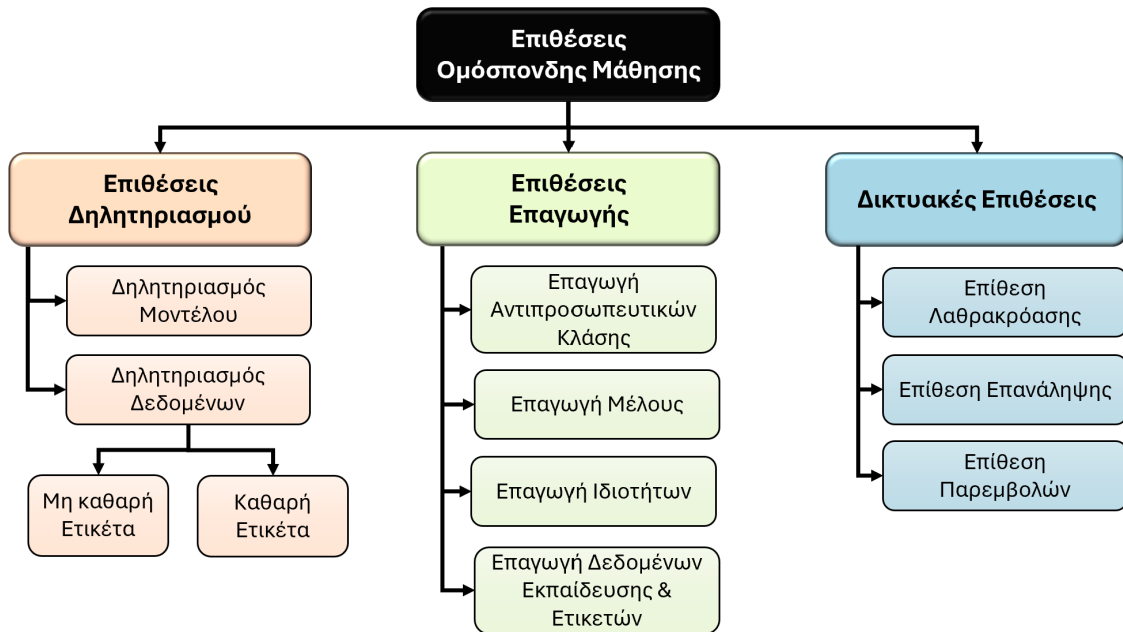
Στην επίθεση αυτή, ένας κακόβουλος χρήστης καταγράφει την κίνηση του Δικτύου Ομόσπονδης Μάθησης, με σκοπό να αναπαράγει τα μηνύματα μεταξύ συμμετεχόντων και διακομιστή κάποια μελλοντική στιγμή. Για παράδειγμα, ένας κακόβουλος χρήστης θα μπορούσε να κρυφακούσει και να αποθηκεύσει το πρώτο συγκεντρωτικό μοντέλο που έστειλε ο εξυπηρετητής στους συμμετέχοντες της Ομόσπονδης Μάθησης και να αναπαράγει το μήνυμα αυτό -προσποιούμενος τον εξυπηρετητή- μετά το πέρας ενός αριθμού γύρων. Αν οι υπόλοιποι συμμετέχοντες δεχθούν το μήνυμα του κακόβουλου χρήστη ως το νέο συγκεντρωτικό μοντέλο που έστειλε ο εξυπηρετητής, τότε ουσιαστικά θα είναι σαν να ξαναρχίζει η Ομόσπονδη Μάθηση από την αρχή. Μια τέτοια επίθεση είναι εύκολο να αποτραπεί, αρκεί να προγραμματιστούν όλα τα εμπλεκόμενα μέρη της Ομόσπονδης Μάθησης ώστε να μην αποδέχονται το ίδιο μήνυμα δύο φορές, είτε με την χρήστη ετικετών χρόνου (time stamps), τα οποία θα δείχνουν πότε στάλθηκε ένα μήνυμα.

3. Επίθεση Παρεμβολών (Jamming Attack)

Όταν κάποιος συμμετέχοντας της Ομόσπονδης Μάθησης επιχειρεί την μετάδοση των παραμέτρων του στον εξυπηρετητή, τότε υπάρχει η πιθανότητα ένας κακόβουλος χρήστης να προσπαθήσει να κάνει παρεμβολές στην μετάδοση αυτή. Για να το επιτύχει αυτό πραγματοποιεί μία ευρεία εκπομπή (broadcasting) στην ίδια συχνότητα με τον στόχο του, χρησιμοποιώντας πολύ μεγάλη ισχύ. Η επίθεση αυτή οδηγεί σε παρεμβολές, με αποτέλεσμα όταν το σήμα του συμμετέχοντα φτάσει στον διακομιστή, αυτός να μην μπορεί να το αποκωδικοποιήσει σωστά και συνεπώς να το απορρίψει. Έτσι είναι εφικτό ένα πλήθος χρηστών να αποκλειστεί από έναν γύρο της Ομόσπονδης Μάθησης, αφού το σήμα τους δεν έφτασε αρκετά "καθαρά" στον εξυπηρετητή. Αυτά τα φαινόμενα μπορούν να επιλυθούν με την χρήση κατευθυνόμενων κεραιών (directional antennas), τεχνικών beamforming και MIMO (Multiple Input, Multiple Output), ωστόσο απαιτούν την χρήση επιπλέον εξοπλισμού, κάτι που αυξάνει τόσο το κόστος όσο και την πολυπλοκότητα της εγκατάστασης.

Πέρα όμως από τις επιθέσεις που λαμβάνουν χώρα σε επίπεδο δικτύου, η Ομόσπονδη Μάθηση πλήττεται και από άλλα είδη επιθέσεων, που έχουν σκοπό είτε να μειώσουν το ποσοστό επιτυχίας του παραγόμενου συγκεντρωτικού μοντέλου, είτε να αποσπάσουν ευαίσθητες πληροφορίες χρησιμοποιώντας τα τοπικά μοντέλα των χρηστών.

Οι **επιθέσεις επαγωγής** (inference attacks) αποτελούν σημαντική απειλή για την προστασία των προσωπικών δεδομένων των χρηστών στην Ομόσπονδη Μάθηση. Σκοπός των επιθέσεων αυτών είναι να αποσπάσουν πληροφορίες για τα σύνολα δεδομένων που χρησιμοποιήθηκαν κατά την εκπαίδευση των τοπικών μοντέλων των συμμετεχόντων [10]. Παρόλο που τα σύνολα δεδομένων αυτά είναι αποθηκευμένα τοπικά, στις συσκευές των χρηστών, οι επιτιθέμενοι μπορούν να αξιοποιήσουν τα μοντέλα μηχανικής μάθησης που μεταδίδονται προς τον διακομιστή, ώστε να λάβουν πληροφορίες για αυτά [11]. Για παράδειγμα, οι εργασίες [12], [13] απέδειξαν ότι είναι εφικτή η ανακατασκευή μίας εικόνας που χρησιμοποιήθηκε για την εκπαίδευση ενός νευρωνικού δικτύου, με την αξιοποίηση των μεταβολών (gradient) των μοντέλων. Το είδος αυτής της επίθεσης επαγωγής ονομάζεται **Επαγωγή Δεδομένων Εκπαίδευσης και Ετικετών** (Training Inputs and Labels Inference Attack). Η δημοσίευση [14] πραγματοποιεί δύο διαφορετικά είδη επιθέσεων επαγωγής. Από την μία, η **Επίθεση Επαγωγής Μέλους** (membership inference attack) επιτρέπει



Σχήμα 2.4: Τα διαφορετικά είδη επιθέσεων στην Ομόσπονδη Μάθηση

σε κάποιον που απλώς έχει πρόσβαση στα τοπικά μοντέλα, να επαληθεύσει αν κάποιο συγκεκριμένο δεδομένο χρησιμοποιήθηκε κατά την εκπαίδευση του μοντέλου. Για παράδειγμα είναι εφικτό να αναγνωρίσει κάποιος αν ένα συγκεκριμένο πρόσωπο χρησιμοποιήθηκε για την εκπαίδευση ενός μοντέλου δυαδικής ταξινόμησης φύλων (binary gender classifier). Οι **Επιθέσεις Επαγωγής Ιδιοτήτων** (Property Inference Attacks), από την άλλη, σχετίζονται με την δυνατότητα εύρεσης συγκεκριμένων χαρακτηριστικών που είναι ανεξάρτητα από αυτά που χρησιμοποιούνται για την εκπαίδευση. Στο παράδειγμα του ταξινομητή φύλων, μία τέτοια ιδιότητα θα μπορούσε να είναι το αν ένα πρόσωπο φοράει γυαλιά μυωπίας. Η εργασία [15] παρουσίασε πρώτη την χρήση GANs (Generative Adversarial Networks) έτσι ώστε να επανακατασκευάσει αντιπροσωπευτικά δείγματα της κάθε κλάσης που χρησιμοποιήθηκαν κατά την εκπαίδευση μοντέλων μηχανικής μάθησης, μια επίθεση που αποκαλείται ως **Επίθεση Επαγωγής Αντιπροσωπευτικών Κλάσης** (Class Representatives Inference Attack) [10] ή βασιζόμενη σε GAN Επίθεση Επαγωγής (GANs-based inference attack).

Για την αντιμετώπιση επιθέσεων επαγωγής έχουν αναπτυχθεί διάφορες τεχνικές, όπως για παράδειγμα το Differential Privacy, του οποίου σκοπός είναι η εισαγωγή θορύβου σε ευαίσθητα χαρακτηριστικά των χρηστών, πριν αυτοί κοινοποιήσουν τα μοντέλα τους στον διακομιστή [16], [17]. Μια άλλη μέθοδος που αντιμετωπίζει επιτυχώς τις επιθέσεις επαγωγής αναπτύχθηκε από τους ερευνητές της Google [18], οι οποίοι δημιούργησαν έναν αλγόριθμο ασφαλούς συνάθροισης δεδομένων, χρησιμοποιώντας την τεχνική Secure Multiparty Computation προκειμένου να υπολογίσουν τον σταθμισμένο μέσο όρο των απεσταλμένων τοπικών μοντέλων, χρησιμοποιώντας κρυπτογραφικές μεθόδους.

Οι **επιθέσεις δηλητηριασμού** περιλαμβάνουν έναν κακόβουλο χρήστη που σκοπίμως επιδιώκει να μειώσει το ποσοστό επιτυχίας του συγκεντρωτικού μοντέλου. Διεξάγει επιθέσεις είτε τροποποιώντας το τοπικό του μοντέλο, είτε αλλοιώνοντας τα δεδομένα που χρησιμοποιούνται για την εκπαίδευσή του. Οι επιθέσεις αυτές δεν είναι πάντα εύκολο να εντοπιστούν, λόγω της κατακεκομμένης φύσης της Ομόσπονδης Μάθησης, και μπορούν να είναι ιδιαίτερα επιβλαβείς, ειδικά σε σενάρια Cross-Silo Ομόσπονδης Μάθησης, όπου το πλήθος συμμετεχόντων δεν είναι ιδιαίτερα μεγάλο. Υπάρχουν δύο βασικά είδη επιθέσεων

δηλητηριασμού:

1. Επιθέσεις Δηλητηριασμού Μοντέλου

Ο κακόβουλος συμμετέχοντας επιδιώκει να τροποποιήσει τις παραμέτρους του μοντέλου που αποστέλλονται στον διακομιστή, κατά την διαδικασία εκπαίδευσης, προκειμένου το συγκεντρωτικό μοντέλο που θα παραχθεί, να έχει μειωμένο ποσοστό επιτυχίας ή να αλλάξει η συμπεριφορά του όταν δοκιμαστεί σε συγκεκριμένες εισόδους.

2. Επιθέσεις Δηλητηριασμού Δεδομένων

Οι επιθέσεις δηλητηριασμού πραγματοποιούνται με την προσθήκη λανθασμένων ή κακόβουλων δεδομένων στο σύνολο εκπαίδευσης του τοπικού μοντέλου των χρηστών. Το αποτέλεσμα της επίθεσης είναι τα τοπικά μοντέλα που εκπαιδεύονται είτε να έχουν χαμηλότερο ποσοστό επιτυχίας είτε να συμπεριφέρονται με τρόπο που ωφελεί τον κακόβουλο χρήστη. Χωρίζονται σε δύο βασικές κατηγορίες:

(α) Επιθέσεις Καθαρής Ετικέτας

Στις επιθέσεις καθαρής ετικέτας ένας κακόβουλος συμμετέχοντας εισάγει στο σύνολο εκπαίδευσης ειδικά δεδομένα, χωρίς να βάζει λανθασμένη ετικέτα (label) σε αυτά. Τα δεδομένα αυτά, αν και φαίνονται σωστά, είναι σχεδιασμένα έτσι ώστε εκπαιδεύσουν το τοπικό και κατ' επέκταση το συγκεντρωτικό μοντέλο, να εμφανίζει λανθασμένες απαντήσεις σε συγκεκριμένες εισόδους [19].

(β) Επιθέσεις μη Καθαρής Ετικέτας

Στην επίθεση αυτή, ο κακόβουλος συμμετέχοντας αλλάζει τις ετικέτες των δεδομένων του συνόλου εκπαίδευσης, με αποτέλεσμα το τοπικό μοντέλο που εκπαιδεύεται να συσχετίζει συγκεκριμένα χαρακτηριστικά (features) με λανθασμένες ετικέτες [20]. Ένα παράδειγμα τέτοιας επίθεσης είναι η **αλλαγή ετικέτας** (label flipping), κατά την οποία αλλάζουν οι ετικέτες μιας κλάσης δεδομένων σε μία άλλη. Για παράδειγμα, στο MNIST dataset που περιλαμβάνει χειρόγραφα ψηφία, είναι πιθανό ένας κακόβουλος συμμετέχοντας να αλλάξει την ετικέτα των εικόνων με τον αριθμό 5 στον αριθμό 3 και αντίθετα.[21]

Η συνεισφορά της παρούσας διπλωματικής εργασίας έγκειται στην αποτροπή των **επιθέσεων παρεμβολών** και ταυτόχρονα στον περιορισμό των επιθέσεων δηλητηριασμού. Πιο συγκεκριμένα, μελετήθηκαν **επιθέσεις δηλητηριασμού δεδομένων μη καθαρής ετικέτας**, ωστόσο η μέθοδος που προτείνουμε, επεκτείνεται εύκολα και στα υπόλοιπα είδη επιθέσεων δηλητηριασμού.

2.2 Τηλεπικοινωνιακό Μοντέλο NOMA

Η Μη Ορθογώνια Πολλαπλή Πρόσβαση (Non Orthogonal Multiple Access - NOMA) είναι μια καινοτόμος τεχνική στις ασύρματες επικοινωνίες, που εφαρμόζεται σε κυψελοειδή δίκτυα 5^{ης} γενιάς (5G) και σε μελλοντικά ασύρματα επικοινωνιακά συστήματα (beyond 5G). Η βασική ιδέα της NOMA είναι η ταυτόχρονη εξυπηρέτηση πολλαπλών χρηστών στις ίδιες συχνότητες φάσματος, με τις ελάχιστες δυνατές παρεμβολές μεταξύ των διαφορετικών σημάτων. Σε αντίθεση με τα παραδοσιακά σχήματα Ορθογώνιας Πολλαπλής Πρόσβασης (Orthogonal Multiple Access - OMA), όπου οι χρήστες κατανέμονται σε ξεχωριστές συχνότητες ο καθένας, η NOMA επιτρέπει σε πολλούς χρήστες να μοιράζονται την ίδια συχνότητα αξιοποιώντας τα διαφορετικά κέρδη καναλιού τους [22].

Στην κατερχόμενη ζεύξη (downlink NOMA), όταν δηλαδή ένας ασύρματος σταθμός προσπαθεί να εκπέμψει διαφορετικά σήματα σε πολλούς χρήστες, είναι δυνατή η επιλογή συγκεκριμένων τιμών ισχύος μετάδοσης για κάθε χρήστη, ώστε να μπορέσουν όλα τα σήματα να μεταδοθούν ταυτόχρονα και να αποκωδικοποιηθούν ορθά από τους παραλήπτες. Αναλυτικές προσομοιώσεις απέδειξαν ξεκάθαρα πλεονεκτήματα στην χρήση τεχνικών NOMA έναντι OMA [23].

Η εργασία [24] ήταν η πρώτη που μελέτησε την χρήση της τεχνικής NOMA στην ανερχόμενη ζεύξη (uplink NOMA). Σε αυτό το σενάριο, οι χρήστες μεταδίδουν ταυτόχρονα στον σταθμό βάσης το μήνυμά τους και ο αποδέκτης το αποκωδικοποιεί χρησιμοποιώντας το ελάχιστο μέσο τετραγωνικό σφάλμα στην διαδικασία διαδοχικής ακύρωσης παρεμβολών (Successive Interference Cancellation - SIC). Ωστόσο, αν οι χρήστες που εκπέμπουν τα δεδομένα τους δεν χρησιμοποιούν συγκεκριμένο ρυθμό μετάδοσης, είναι πολύ πιθανό -λόγω παρεμβολών- να απορριφθούν τα σήματά τους, κάτι που δεν συμβαίνει, ως επί το πλείστον, στις παραδοσιακές τεχνικές OMA [25].

Το γεγονός ότι τα σήματα κάποιων χρηστών είναι πιθανό να απορριφθούν στην τεχνική NOMA οφείλεται στον τρόπο που λειτουργεί διαδικασία διαδοχικής ακύρωσης παρεμβολών (SIC). Η SIC είναι μια τεχνική επεξεργασίας των σημάτων, που χρησιμοποιώντας τα κέρδη καναλιού των χρηστών καταφέρνει με επαναληπτικό τρόπο να αφαιρέσει το σήμα με το μεγαλύτερο κέρδος καναλιού, ώστε να παραμείνουν τα υπόλοιπα. Πιο αναλυτικά, η μέθοδος SIC λειτουργεί με τον ακόλουθο τρόπο [26], [27]:

1. Το υπέρθετο σήμα

Στην NOMA ανερχόμενης ζεύξης, οι χρήστες εκπέμπουν ταυτόχρονα τα σήματά τους, στην ίδια συχνότητα αλλά με διαφορετική ισχύ. Έτσι ο αποδέκτης σταθμός βάσης λαμβάνει ένα υπέρθετο σήμα (superimposed signal), που αποτελείται από τα επιμέρους σήματα των χρηστών

2. Η σειρά αποκωδικοποίησης

Ο παραλήπτης αποκωδικοποιεί τα σήματα σε φθίνουσα σειρά κέρδους καναλιού. Το κέρδος καναλιού στην επικοινωνία ενός πομπού με έναν δέκτη επηρεάζεται από παραμέτρους όπως την μεταξύ τους απόσταση, την εξασθένηση και τις απώλειες. Ως εκ τούτου, τα σήματα που φτάνουν με μεγάλο κέρδος καναλιού, είναι συνήθως τα πιο ισχυρά, αφού ο πομπός τους βρίσκεται πιο κοντά στον σταθμό βάσης.

3. Η αφαίρεση των σημάτων

Όταν το ισχυρότερο σήμα, αυτό με το μεγαλύτερο κέρδος καναλιού, αποκωδικοποιηθεί, τότε αφαιρείται από το υπέρθετο σήμα. Έτσι μειώνονται οι παρεμβολές που επηρεάζουν τα υπόλοιπα σήματα.

4. Επαναληπτική Αποκωδικοποίηση

Στην συνέχεια ο δέκτης μεταβαίνει στην αποκωδικοποίηση του αμέσως ισχυρότερου σήματος από το υπέρθετο σήμα, το αφαιρεί, και συνεχίζει επαναληπτικά έως ότου ανακτηθούν όλα τα σήματα των χρηστών.

Η παραπάνω διαδικασία ακύρωσης παρεμβολών έχει ως αποτέλεσμα ο χρήστης με το μεγαλύτερο κέρδος καναλιού να υπόκειται σε παρεμβολές από όλους τους άλλους χρήστες, αφού το σήμα του είναι το πρώτο που πρέπει να αποκωδικοποιηθεί από το υπέρθετο σήμα. Όταν όμως το σήμα του αφαιρεθεί από τα υπόλοιπα, τότε δεν προκαλεί παρεμβολές σε κανέναν. Έτσι κατά την αποκωδικοποίηση του δεύτερου ισχυρότερου σήματος, δεν υπάρχουν παρεμβολές από το ισχυρότερο, μονάχα από τα ασθενέστερα σήματα που δεν έχουν αφαιρεθεί ακόμα. Η διαδικασία αυτή ευνοεί τον χρήστη με το χαμηλότερο κέρδος καναλιού, καθώς την στιγμή που θα αποκωδικοποιηθεί το σήμα του, θα έχουν αφαιρεθεί τα σήματα όλων των υπολοίπων χρηστών, με αποτέλεσμα να μην αισθάνεται καθόλου παρεμβολές από αυτούς.

2.3 Θεωρία Παιγνίων

Η Θεωρία Παιγνίων αποτελεί τον κλάδο των εφαρμοσμένων μαθηματικών που παρέχει το κατάλληλο υπόβαθρο για την ανάλυση στρατηγικών αλληλεπιδράσεων μεταξύ λογικών οντοτήτων, που αποκαλούνται παίκτες [28]. Η θεωρία παιγνίων περιλαμβάνει εργαλεία για τη μελέτη καταστάσεων όπου άτομα ή ομάδες λαμβάνουν αλληλοεξαρτώμενες αποφάσεις, δηλαδή περιπτώσεις στις οποίες το αποτέλεσμα για κάθε παίκτη εξαρτάται από τις επιλογές των υπολοίπων. Μάλιστα κάθε παίκτης καλείται να εξετάσει τις πιθανές αποφάσεις των υπολοίπων, προκειμένου να διαμορφώσει την δική του στρατηγική, με γνώμονα την μεγιστοποίηση του κέρδους του. Όπως εξηγεί ο Peters στο έργο του [29], η θεωρία παιγνίων έχει πλέον ένα ευρύ φάσμα εφαρμογών, όπως στην πολιτική επιστήμη και τη βιολογία, παρόλο που πρωτοεφαρμόστηκε στον κλάδο των οικονομικών.

Όπως ανέφερε ο John Harsanyi, που τιμήθηκε με βραβείο Νόμπελ Οικονομικών για την συνεισφορά του στην θεμελίωση της θεωρίας παιγνίων [30]:

"Game theory is a theory of strategic interaction. That is to say, it is a theory of rational behavior in social situations in which each player has to choose his moves on the basis of what he thinks the other players' countermoves are likely to be."

Δηλαδή η θεωρία Παιγνίων είναι η θεωρία στρατηγικών αλληλεπιδράσεων. Αυτό σημαίνει ότι είναι η θεωρία των λογικών συμπεριφορών σε κοινωνικές καταστάσεις όπου κάθε παίχτης πρέπει να επιλέξει την δράση του, βασιζόμενος στην εκτίμησή του για το ποία είναι η πιο πιθανή συμπεριφορά των υπολοίπων παικτών.

Κάποιες από τις βασικότερες κατηγορίες διαχωρισμού των παιγνίων είναι οι ακόλουθες:

1. Συνεργατικά (Cooperative) και μη Συνεργατικά Παιγνία (Non-Cooperative)

Στα Συνεργατικά Παιγνία οι παίκτες συνεργάζονται προκειμένου να επιτύχουν αμοιβαία οφέλη. Η έμφαση δίνεται στην συλλογική στρατηγική και στην κατανομή των αποδόσεων μεταξύ των παικτών. Αντίθετα στα μη Συνεργατικά Παιγνία, οι παίκτες λαμβάνουν ατομικές αποφάσεις, με σκοπό την μεγιστοποίηση των δικών τους αποδόσεων. Έτσι οι στρατηγικές που υιοθετούνται αναπτύσσονται ανεξάρτητα και τα αποτελέσματα εξαρτώνται από τις κινήσεις των υπολοίπων [31], [32].

2. Συμμετρικά (Symmetric) και Μη Συμμετρικά (Asymmetric)

Τα συμμετρικά παίγνια είναι εκείνα στα οποία όλοι οι παίκτες έχουν το ίδιο σύνολο στρατηγικής και το όφελος της κάθε στρατηγικής εξαρτάται μόνο από την ίδια, όχι από το ποιος την επιλέγει. Όταν οι παίκτες, πραγματοποιώντας την ίδια επιλογή, έχουν το ίδιο ακριβώς όφελος, τότε το παίγνιο είναι συμμετρικό. Στα ασύμμετρα παίγνια, οι αποδόσεις διαφέρουν από παίκτη σε παίκτη, ανάλογα τον ρόλο και την στρατηγική που επιλέγει. [33], [34]

3. Τέλεια (Perfect) και Ατελούς πληροφορίας (Imperfect Information)

Τα παιχνίδια τέλει πληροφορίας είναι εκείνα στα οποία όλοι οι παίκτες γνωρίζουν ανά πάσα στιγμή, τα πάντα για το παίγνιο που διεξάγεται. Αντίθετα, στα παιχνίδια ατελούς πληροφόρησης περιλαμβάνονται καταστάσεις όπου οι συμμετέχοντες δεν έχουν πλήρη γνώση για τις προηγούμενες ενέργειες ή στρατηγικές που επιλέχθηκαν από άλλους παίκτες [28], [35].

2.3.1 Παίγνια σε Κανονική Μορφή Αναπαράστασης

Η συνάρτηση ωφέλειας (payoff function) είναι μια από τις σημαντικότερες έννοιες στην θεωρία παιγνίων. Η συνάρτηση αυτή αντιπροσωπεύει το κέρδος του εκάστοτε παίκτη, που εξαρτάται από την στρατηγική που θα επιλέξει σε αυτό. Για παράδειγμα, σε επιχειρησιακά παίγνια, η συνάρτηση ωφέλειας είναι πιθανόν να αντιστοιχεί στο οικονομικό όφελος μιας εταιρίας, όταν επιλέξει μια συγκεκριμένη στρατηγική.

Στην θεωρία Παιγνίων τα παίγνια μπορούν να περιγραφούν με πολλές μορφές. Μια από τις πιο συνηθισμένες είναι η κανονική μορφή, στην οποία τα παίγνια περιγράφονται με την χρήση ενός πίνακα, που περιλαμβάνει τις διαθέσιμες στρατηγικές και την ωφέλεια των χρηστών [36].

Ορισμός 2.1. Η κανονική μορφή αναπαράστασης ενός μη συνεργατικού παιγνίου, έχει την εξής μορφή:

1. Ένα σύνολο παικτών $N = \{1, 2, \dots, n\}$
2. Ένα σύνολο στρατηγικής S_i για κάθε παίκτη i
3. Μια συνάρτηση χρησιμότητας U_i που υπολογίζει την ωφέλεια $U_i(s)$ για κάθε στρατηγική $s = (s_1, s_2, \dots, s_n)$

Όσον αφορά στα μη συνεργατικά παίγνια τέλειας πληροφωρίας, εκτός από την υπόθεση ότι οι παίκτες είναι λογικοί και επιλέγουν ταυτόχρονα την στρατηγική τους, θεωρούμε επιπλέον ότι η δομή του παιχνιδιού είναι απόλυτα γνωστή. Δηλαδή κάθε παίκτης γνωρίζει για κάθε άλλον το σύνολο των δυνατών στρατηγικών του και την συνάρτηση ωφέλειάς του. Αυτό ωστόσο δεν ισχύει στα Μπεϋζιανά Παιγνια ατελούς πληροφωρίας που θα παρουσιαστούν στην συνέχεια.

2.3.2 Ισορροπία Nash

Η ισορροπία Nash (Nash Equilibrium) είναι μια θεμελιώδης έννοια των μη συνεργατικών παιγνίων. Σε ένα μη συνεργατικό παίγνιο, οι συμμετέχοντες ανταγωνίζονται μεταξύ τους σε κάθε γύρο και αναπροσαρμόζουν τις στρατηγικές τους, με στόχο να αυξήσουν την ωφέλειά τους. Όταν ένα παίγνιο φτάσει στην ισορροπία Nash, τότε όλοι οι παίκτες επιτυγχάνουν το μέγιστο δυνατό κέρδος τους, δεδομένων υπόλοιπων στρατηγικών. Πρόκειται δηλαδή για μία κατάσταση, στην οποία κανένας παίκτης δεν θα ωφεληθεί αν αλλάξει την στρατηγική του, δεδομένου ότι όλοι οι υπόλοιποι συμμετέχοντες πραγματοποιούν τις ίδιες επιλογές.

Έστω παίγνιο n παικτών και S_i το σύνολο στρατηγικών του παίκτη i . Αν $u_i(s_1, s_2, \dots, s_n)$ είναι η συνάρτηση χρησιμότητας του παίκτη i , όπου $s_i \in S_i$ αντιπροσωπεύει την στρατηγική που επιλέχθηκε από τον παίκτη i , τότε

Ορισμός 2.2. Το στρατηγικό προφίλ $\mathbf{s}^* = (s_1^*, s_2^*, \dots, s_n^*)$ είναι ισορροπία Nash αν για κάθε παίκτη i :

$$u_i(s_i^*, \mathbf{s}_{-i}^*) \geq u_i(s_i, \mathbf{s}_{-i}^*) \quad \forall s_i \in S_i$$

Όπου το διάνυσμα $\mathbf{s}_{-i}^* = (s_1^*, \dots, s_{i-1}^*, s_{i+1}^*, \dots, s_n^*)$ αναπαριστά τις στρατηγικές όλων των υπόλοιπων παικτών, πλην του i . Ο παραπάνω ορισμός εκφράζει ότι το στρατηγικό προφίλ \mathbf{s}^* είναι Ισορροπία Nash, αν δεν συμφέρει κανέναν παίκτη (έστω i) να διαλέξει οποιαδήποτε άλλη στρατηγική, πέρα από την s_i^* , όταν και οι υπόλοιποι συμμετέχοντες επιλέγουν την βέλτιστη για εκείνους στρατηγική \mathbf{s}_{-i}^* . Ισοδύναμα, αν ο παίκτης i επιλέξει διαφορετική στρατηγική από την s_i^* της ισορροπίας Nash, θα λάβει το ίδιο ή λιγότερο κέρδος [35], [37].

2.3.3 Μπεϋζιανά Παίγνια

Τα Μπεϋζιανά Παίγνια (Bayesian Games) είναι μια υποκατηγορία της θεωρίας Παιγνίων που ενσωματώνει την αβεβαιότητα κάποιων πτυχών του περιβάλλοντος, όπως τα είδη ή τις προτιμήσεις παικτών. Σε αντίθεση με τα κλασικά μη συνεργατικά παίγνια, στα οποία θεωρείται ότι οι συμμετέχοντες έχουν πλήρη γνώση των δρώμενων, τα Μπεϋζιανά Παίγνια είναι ατελούς πληροφορίας. Σε αυτά κάθε παίκτης έχει μια κρυφή πληροφορία ή έναν κρυφό τύπο, που επηρεάζει την συνάρτηση χρησιμότητάς του, συνεπώς και το κέρδος του [35].

Συνεπώς, τα Μπεϋζιανά Παίγνια είναι πιο ρεαλιστικά, αφού περιλαμβάνουν την αβεβαιότητα που συχνά χαρακτηρίζει τις αλληλεπιδράσεις ατόμων στην κοινωνία, για αυτό και επιλέγονται στην παρούσα εργασία για την προσομοίωση των συμμετεχόντων σε ένα δίκτυο Ομόσπονδης Μάθησης. Όπως θα αναλυθεί και στην ενότητα 4, θεωρούμε ότι ο κάθε παίκτης έχει διαφορετικό τύπο (type), ο οποίος εκφράζει τις προτιμήσεις και τους στόχους του. Για παράδειγμα, ο τύπος ενός παίκτη θα μπορούσε να είναι κακόβουλος, αν επιδιώκει να βλάψει τους υπόλοιπους, ή καλόβουλος, αν επιθυμεί να τους βοηθήσει. Ο τύπος ενός παίκτη διαμορφώνει το κέρδος του και διατηρείται κρυφός από τους υπόλοιπους παίκτες.

Εξαιτίας λοιπόν της ατελούς αυτής πληροφορίας, στα Μπεϋζιανά παίγνια υπάρχουν ως κυρίαρχο χαρακτηριστικό οι πεποιθήσεις των συμμετεχόντων, σχετικά με τους τύπους των άλλων χρηστών. Κάθε χρήστης, δηλαδή, ξεκινάει με μια αρχική εικασία, για το ποιος μπορεί να είναι ο τύπος των υπολοίπων συμμετεχόντων και σταδιακά ανανεώνει την εκτίμησή του αυτή, όσο εξελίσσεται το παίγνιο. Η βέλτιστη λύση αυτών των παιγνίων ονομάζεται Bayesian Nash equilibrium και (κατ' αντιστοιχία με το Nash Equilibrium) εκφράζει την κατάσταση, στην οποία οι παίκτες έχουν επιλέξει στρατηγικές που μεγιστοποιούν το κέρδος τους και δεν επιθυμούν να παρεκκλίνουν από αυτές [36].

Bayesian Nash Equilibrium

Ορισμός 2.3. Το στρατηγικό προφίλ $\mathbf{s}^* = (s_1^*, s_2^*, \dots, s_n^*)$ είναι Bayesian Nash Equilibrium αν για κάθε παίκτη i με τύπο $t_i \in T_i$:

$$\mathbb{E}_{t_{-i} \sim \mu_i(t_{-i})} [u_i(s_i^*(t_i), s_{-i}^*(t_{-i}), t_i, t_{-i})] \geq \mathbb{E}_{t_{-i} \sim \mu_i(t_{-i})} [u_i(s_i(t_i), s_{-i}^*(t_{-i}), t_i, t_{-i})] \quad \forall s_i \in S_i$$

όπου $t_{-i} \sim \mu_i(t_{-i})$ εκφράζει την πεποίθηση του παίκτη i για τους τύπους των υπολοίπων χρηστών, όπως αυτή υπολογίζεται με μια πιθανοτική κατανομή $\mu_i : T_{-i} \rightarrow [-1, 1]$.

2.3.4 Τιμή Shapley

Η τιμή Shapley (Shapley Value) έλαβε το όνομά της από τον Lloyd Shapley, που την πρωτοεισηγάγε το 1953 σαν λύση σε συνεργατικά παίγνια [38]. Είναι μία μέθοδος που κατανέμει δίκαια το συνολικό κέρδος (ή κόστος) ανάμεσα στους παίκτες, ανάλογα με την ατομική τους συνεισφορά στο παίγνιο. Επειδή η τιμή Shapley έχει κάποιες σημαντικές ιδιότητες, όπως συμμετρία και γραμμικότητα, αποτελεί έναν μοναδικό και δίκαιο τρόπο κατανομής πόρων.

Αναλυτικότερα, σε ένα συνεργατικό παίγνιο, οι παίκτες διαμορφώνουν συμμαχίες και η τιμή κάθε συμμαχίας υπολογίζεται από μία χαρακτηριστική συνάρτηση. Για παράδειγμα, στην ομόσπονδη Μάθηση, η συνάρτηση αυτή θα μπορούσε να είναι το ποσοστό επιτυχίας κάποιου μοντέλου Μηχανικής Μάθησης. Η τιμή Shapley αντιστοιχεί ένα κέρδος σε κάθε παίκτη, υπολογίζοντας την συνεισφορά ενός παίκτη σε κάθε πιθανή συμμαχία παικτών [36].

Ορισμός 2.4. Σε ένα συνεργατικό παίγνιο n συμμετεχόντων, η τιμή *Shapley* $\phi_i(v)$ του παίκτη i υπολογίζεται ως εξής:

$$\phi_i(v) = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(n - |S| - 1)!}{n!} (v(S \cup \{i\}) - v(S))$$

όπου:

- N το σύνολο των παικτών,
- f είναι η χαρακτηριστική συνάρτηση που αντιστοιχεί τιμή σε μία συμμαχία,
- S είναι ένα υποσύνολο των παικτών N , που δεν περιλαμβάνει τον χρήστη i ,
- $|S|$ το πλήθος των παικτών της συμμαχίας S .

Ουσιαστικά ο βαθμός συνεισφοράς του κάθε παίκτη i , υπολογίζεται μέσω της σύγκρισης των αποτελεσμάτων μιας συμμαχίας S , με την συμμαχία $S \cup \{i\}$, που περιλαμβάνει τον i . Αν η συμμαχία $S \cup \{i\}$ επιτυγχάνει καλύτερα αποτελέσματα από την συμμαχία S , τότε ο χρήστης i συνεισέφερε σημαντικά σε αυτό, με αποτέλεσμα να επιβραβεύεται αντίστοιχα.

Σχετική Βιβλιογραφία

Η Ομόσπονδη Μάθηση έχει αναδειχθεί ως ένα ισχυρό εργαλείο για την εκπαίδευση μοντέλων μηχανικής μάθησης με αποκεντρωμένο τρόπο, εξασφαλίζοντας την ιδιωτικότητα των δεδομένων των χρηστών και μειώνοντας το κόστος επικοινωνίας. Ωστόσο, η κατανομημένη φύση της Ομόσπονδης Μάθησης την καθιστά ευάλωτη σε επιθέσεις, αφού προϋποθέτει την συνεργασία και αλληλεπίδραση πολλών εμπλεκόμενων πλευρών, αυξάνοντας έτσι την πολυπλοκότητα των συστημάτων. Η παρούσα ενότητα παρουσιάζει μέρος της διεθνούς βιβλιογραφίας σε θέματα που άπτονται της Ομόσπονδης Μάθησης και της Θεωρίας Παιγνίων, επισημαίνοντας τις υπάρχουσες συνεισφορές και εντοπίζοντας τα κενά που αποτέλεσαν εφαλτήριο για την μελέτη και συγγραφή της παρούσας διπλωματικής.

Κίνητρα Συμμετοχής στην Ομόσπονδη Μάθηση

Η Θεωρία παιγνίων αποτελεί ένα πολύ ισχυρό εργαλείο που χρησιμοποιείται για την διευθέτηση ποικίλων διαφορετικών προβλημάτων στην Ομόσπονδη Μάθηση, όπως την ενίσχυση της συνεργασίας των εμπλεκόμενων μερών, την παροχή κινήτρων και αμοιβών στους συμμετέχοντες καθώς και την εξασφάλιση της ασφάλειας και ιδιωτικότητας των χρηστών [39]. Τα δίκτυα Ομόσπονδης Μάθησης απαρτίζονται συνήθως από ένα μεγάλο πλήθος συσκευών άκρων (edge devices), όπως είναι τα κινητά τηλέφωνα, που έχουν περιορισμένους πόρους να διαθέσουν για την εκπαίδευση μοντέλων μηχανικής μάθησης, καθώς τόσο η μπαταρία τους όσο και η υπολογιστική τους ισχύς είναι χαμηλή. Εργασίες όπως αυτή του Hu [40] αναγνωρίζουν τους περιορισμένους πόρους των κινητών συσκευών και προτείνουν ένα σύστημα βασισμένο σε μηχανισμούς κινήτρων, που επιβραβεύει τους χρήστες όταν συμμετέχουν σε έναν γύρο της Ομόσπονδης Μάθησης. Μάλιστα, μοντελοποιούν ένα Μπεϋζιανό παίγνιο ελλιπούς πληροφορίας, χάρη στο οποίο ο κάθε χρήστης μπορεί να καταφύγει στην βέλτιστη για εκείνον επιλογή, χωρίς να γνωρίζει πόσοι και ποιοι από τους υπόλοιπους χρήστες θα συμμετέχουν στην εκάστοτε εποχή της Ομόσπονδης Μάθησης. Σε παρόμοια λογική κινήθηκε και η εργασία [41], όπου μοντελοποιείται ένα εξελικτικό παίγνιο για να βοηθήσει τις κινητές συσκευές να αποφασίσουν πως θα καταναείμουν την υπολογιστική τους ισχύ και τα δεδομένα τους ανάμεσα σε ένα πλήθος από διαφορετικές εφαρμογές της Ομόσπονδης Μάθησης. Η δημοσίευση [42] κάνει λόγο για εφαρμογή της θεωρίας συμβολαίων στην ομόσπονδη μάθηση, έτσι ώστε να παρέχει κίνητρο σε χρήστες με χρήσιμα σύνολα δεδομένων να εκπαιδεύουν τοπικά μοντέλα και να συμμετέχουν ενεργά στην δημιουργία των συγκεντρωτικών μοντέλων.

Εφαρμογές της τιμής Shapley

Οι παραπάνω εργασίες μελετούν την παροχή κινήτρων και τον διαμοιρασμό των πόρων στην Ομόσπονδη Μάθηση. Παρόμοιους στόχους θέτει όμως και η δημοσίευση [6], στην οποία διεξάγεται ένα συνεργατι-

κό παίγνιο, στο οποίο υπολογίζεται η τιμή Sharpley (Sharpley Value). Η τιμή αυτή αντικατοπτρίζει το κατά πόσο ένα τοπικό μοντέλο είναι αξιόλογο, δηλαδή πόσο συνεισφέρει στην δημιουργία ενός συγκεντρωτικού μοντέλου με υψηλό ποσοστό επιτυχίας (accuracy). Αφού, λοιπόν, υπολογιστεί ο βαθμός συνεισφοράς για κάθε συμμετέχοντα, ο διακομιστής τον χρησιμοποιεί προκειμένου να επιβραβεύσει τους χρήστες ανάλογα με την συμβολή τους. Η εργασία [43] αξιολογεί, με βάση την τιμή Sharpley, τα τοπικά μοντέλα των χρηστών, έτσι ώστε να επιλέξει για την εκπαίδευση του συγκεντρωτικού μοντέλου το πλήθος χρηστών που παρουσιάζουν τα καλύτερα αποτελέσματα, προκειμένου να αυξήσει το τελικό ποσοστό επιτυχίας της Ομόσπονδης Μάθησης. Τέλος, η συνάρτηση συνάθροισης Federated Sharpley Value (FedSv) που προτείνει η εργασία των [44] αποδεδειγμένα βελτιώνει την επιτυχία του συγκεντρωτικού μοντέλου, ακόμα και σε Non-IID δεδομένα, ωστόσο στην μελέτη αυτή δεν συμπεριλαμβάνονται κακόβουλοι χρήστες, ούτε μελετάται η ανθεκτικότητα του FedSV σε επιθέσεις δηλητηριασμού.

Επιθέσεις Παρεμβολών στην Ομόσπονδη Μάθηση

Ένα είδος επιθέσεων που συχνά απειλούν τα ασύρματα δίκτυα, είναι οι επιθέσεις παρεμβολών. Οι εργασίες [45], [46] μελετάνε δίκτυα Ομόσπονδης Μάθησης, υπό την παρουσία συσκευών που προκαλούν παρεμβολές (παρεμβολείς). Ερευνώνται διάφορες στρατηγικές επίθεσης, τόσο στην ανερχόμενη όσο και στην κατερχόμενη ζεύξη, που απέδειξαν ότι οι επιθέσεις αυτές μπορούν να μειώσουν σημαντικά την απόδοση του συγκεντρωτικού μοντέλου της Ομόσπονδης Μάθησης. Η δημοσίευση [47] μελετάει την άμυνα σε τέτοιες επιθέσεις παρεμβολών στην Ομόσπονδη Μάθηση, με την διεξαγωγή ενός Stackelberg παιγνίου. Στο παίγνιο αυτό θεωρείται ως αρχηγός ένας παρεμβολέας, ενώ ως ακόλουθοι οι υπόλοιποι συμμετέχοντες. Ενώ η εργασία αυτή είναι η μόνη που μελετά την άμυνα σε επιθέσεις παρεμβολών σε ασύρματα δίκτυα Ομόσπονδης Μάθησης, η μοντελοποίησή τους δεν καλύπτει σενάρια όπου οι παρεμβολές πραγματοποιούνται συγχρόνως με τις μεταδόσεις των χρηστών ή του διακομιστή, αλλά προϋποθέτουν ότι ο παρεμβολέας ξεκινάει πρώτος και οι υπόλοιποι χρήστες ρυθμίζουν την ισχύ τους αντίστοιχα.

Παρά τις παραπάνω δημοσιεύσεις, υπάρχει ένα εμφανές κενό στην βιβλιογραφία, όσον αφορά στις επιθέσεις παρεμβολών. Ακόμα και δημοσιεύσεις που τις αναγνωρίζουν ως πρόβλημα ([9]), προτείνουν ως λύση είτε την αύξηση της ισχύος μετάδοσης των χρηστών, ώστε να αυξηθεί ο σηματοθορυβικός λόγος (SNR), είτε την χρήση κατευθυνόμενων κεραιών (directional antennas), ώστε να αγνοηθεί ο θόρυβος από διαφορετικές κατευθύνσεις. Οι λύσεις αυτές όμως οδηγούν σε σπατάλη ενέργειας, που είναι από τους πλέον περιορισμένους πόρους μίας κινητής συσκευής και προσθέτουν επιπλέον κόστος στην εγκατάσταση. Η χρήση της θεωρίας παιγνίων μπορεί να βοηθήσει να αντιμετωπιστούν επιθέσεις παρεμβολών, χωρίς να προϋποθέτει την χρήση επιπλέον υλικού ή να υπαγορεύει την αλόγιστη αύξηση της ισχύος μετάδοσης.

Επιθέσεις Δηλητηριασμού στην Ομόσπονδη Μάθηση

Ενώ λοιπόν η χρήση της θεωρίας παιγνίων σε δίκτυα Ομόσπονδης Μάθησης έχει πολλές εφαρμογές, όπως για κατανομή πόρων, παροχή κινήτρων και προστασία των δεδομένων [48], μια εφαρμογή που δεν έχει μελετηθεί αρκετά είναι η άμυνα απέναντι σε επιθέσεις δηλητηριασμού, με την χρήση μεθόδων θεωρίας παιγνίων. Η δημοσίευση [49] προτείνει έναν νέο αλγόριθμο, αντί του Federated Averaging (FedAvg), που βασίζεται σε παιγνιοθεωρητικές τεχνικές. Ο αλγόριθμος αυτός, χρησιμοποιώντας την ιδιότητα της Ισοροπίας Nash, επιτρέπει στον διακομιστή να εκτιμήσει την πιθανότητα ένας χρήστης να πραγματοποιεί επιθέσεις δηλητηριασμού, με αποτέλεσμα να περιορίσει τον αντίκτυπό του. Ωστόσο, πέρα από την παραπάνω προσέγγιση, δεν υπάρχει, από όσο γνωρίζουμε, περαιτέρω συμβολή στην αντιμετώπιση επιθέσεων δηλητηριασμού, με τεχνικές της θεωρίας παιγνίων.

Συνοψίζοντας, είναι εμφανές ότι υπάρχει κενό στη διεθνή βιβλιογραφία όσον αφορά τη χρήση της

θεωρίας παιγνίων για την αντιμετώπιση των σημαντικότερων επιθέσεων που απειλούν την Ομόσπονδη Μάθηση. Η παρούσα διπλωματική εργασία επιχειρεί να καλύψει την ανάγκη για ανθεκτικότερα δίκτυα Ομόσπονδης Μάθησης, παρουσιάζοντας έναν αποτελεσματικό τρόπο αντιμετώπισης τόσο των επιθέσεων παρεμβολών όσο και των επιθέσεων δηλητηριασμού, με ενιαίο τρόπο. Οι καινοτομίες όμως δεν περιορίζονται αποκλειστικά στο γεγονός ότι χρησιμοποιήθηκαν παιγνιοθεωρητικές τεχνικές, αφού η από κοινού αντιμετώπιση δύο ειδών επιθέσεων είναι από μόνη της αρκετά σπάνια και πρωτότυπη στον χώρο της Ομόσπονδης Μάθησης.

Μοντελοποίηση Συστήματος

Το σύστημα που μοντελοποιεί η παρούσα εργασία περιλαμβάνει ένα σύνολο ιδιοτελών χρηστών N καθώς και ένα σύνολο από κακόβουλους συμμετέχοντες M , όπου

$$M = \{1, \dots, m, \dots, |M|\}$$

$$N = \{1, \dots, n, \dots, |N|\}$$

Όλοι οι χρήστες που συμμετέχουν στο δίκτυο ομόσπονδης μάθησης, καλούνται σε κάθε γύρο να εκπαιδεύσουν τοπικά ένα μοντέλο μηχανικής μάθησης, χρησιμοποιώντας τα διαφορετικά σύνολα δεδομένων (dataset) που έχουν στην διάθεσή τους. Θέτουμε

$$D_i = \{ \mathbf{x}_j, y_j \}_{j=1}^{|D_i|}, \quad i \in N \cup M$$

το σύνολο δεδομένων του χρήστη i , όπου \mathbf{x}_j, y_j είναι τα δείγματα και οι αντίστοιχες κλάσεις, που απαρτίζουν το σύνολο αυτό. Η ένωση όλων των διαφορετικών δεδομένων του συστήματος ομόσπονδης μάθησης, συμβολίζεται με

$$D = \bigcup_{i=1}^{|N|+|M|} D_i = \{ \mathbf{x}_j, y_j \}_{j=1}^{|D|}$$

4.1 Μοντελοποίηση Ομόσπονδης Μάθησης

4.1.1 Εκπαίδευση τοπικών Μοντέλων

Σε κάθε γύρο, ή αλλιώς εποχή, της ομόσπονδης μάθησης οι συμμετέχοντες εκπαιδεύουν το μοντέλο που τους έστειλε ο διακομιστής, χρησιμοποιώντας το σύνολο δεδομένων που έχουν στην διάθεσή τους. Έτσι στον εκάστοτε γύρο t , ο χρήστης i ενημερώνει τα βάρη του μοντέλου του με βάση τον αναδρομικό τύπο:

$$\mathbf{W}_i^t = \mathbf{W}_i^{t-1} - \eta \nabla F_i(D_i, \mathbf{W}_i^{t-1}) \quad (4.1)$$

όπου $\eta \in [0, 1]$ είναι ο ρυθμός μάθησης (learning rate) και F_i είναι η εμπειρική συνάρτηση απωλειών (empirical loss function), που προκύπτει από το σύνολο δεδομένων του χρήστη i και ορίζεται ως εξής:

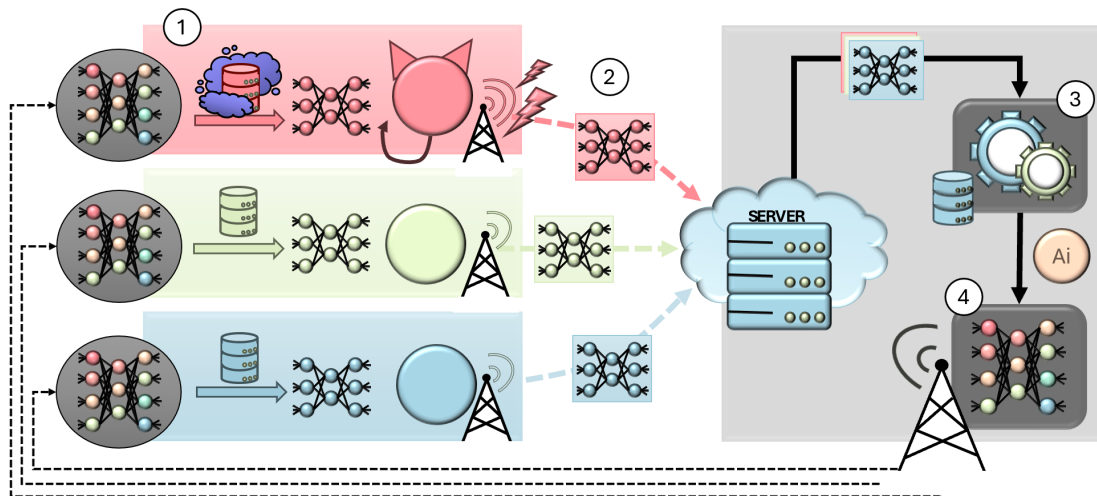
$$F_i(D_i, \mathbf{W}_i) = \frac{1}{|D_i|} \sum_{j=1}^{|D_i|} f(\mathbf{x}_j, y_j, \mathbf{W}_i) \quad (4.2)$$

Η συνάρτηση f στον παραπάνω τύπο αντιπροσωπεύει την εμπειρική συνάρτηση απωλειών του j -οστού δείγματος του σύνολου δεδομένων του χρήστη i και αντιπροσωπεύει το σφάλμα πρόβλεψης του μοντέλου \mathbf{W}_i , όταν στο μοντέλο αυτό δίνεται ως είσοδος το στοιχείο j . Κατά την εκπαίδευση των τοπικών μοντέλων χρησιμοποιούμε ως συνάρτηση f την συνάρτηση απωλειών διασταυρούμενης εντροπίας (Cross Entropy Loss Function)

Μόλις ολοκληρωθεί η εκπαίδευση των τοπικών μοντέλων ενός γύρου t , ανανεώνοντας τα βάρη με την χρήση των εξισώσεων 4.1 και 4.2, τα μοντέλα αποστέλλονται στον διακομιστή, όπου πραγματοποιείται πρώτα η αξιολόγηση και κατόπιν η συνάθροισή τους.

4.1.2 Υπολογισμός Συνεισφοράς

Αφού αποσταλούν οι παράμετροι των τοπικών μοντέλων στον κεντρικό εξυπηρετητή, αυτός αναλαμβάνει σε πρώτη φάση την αξιολόγηση των μοντέλων που έλαβε. Για να συμβεί κάτι τέτοιο, θεωρούμε ότι ο διακομιστής έχει και αυτός ένα δικό του σύνολο δεδομένων, έστω D_{server} , το οποίο χρησιμοποιεί για να εκτιμήσει την συνεισφορά του κάθε συμμετέχοντα, δηλαδή το κατά πόσο το μοντέλο που στέλνει ο κάθε χρήστης βοηθάει στην διαμόρφωση του συγκεντρωτικού μοντέλου. Αυτό το βήμα είναι ύψιστης σημασίας, καθώς χάρη σε αυτό είναι εφικτός ο εντοπισμός επιθέσεων δηλητηριασμού που διενεργούνται από κακόβουλους χρήστες, έτσι ώστε σε επόμενα βήματα να περιοριστεί η ζημιά που προκαλούν.



Σχήμα 4.1: Το μοντέλο που προτείνει η παρούσα εργασία. Ο κακόβουλος χρήστης αποτυπώνεται με κόκκινο χρώμα και πραγματοποιεί τόσο επιθέσεις παρεμβολών (2) όσο και δηλητηριασμού δεδομένων (1). Αφού ο διακομιστής λάβει τα τοπικά μοντέλα, υπολογίζει τον βαθμό συνεισφοράς (a_i) τους (3) και στην συνέχεια παράγεται το συγκεντρωτικό μοντέλο του γύρου (4) με σεβασμό στον βαθμό a_i .

Η παρούσα πτυχιακή εργασία, εμπνευσμένη από την δουλειά των [6], χρησιμοποιεί την τιμή Shapley (Shapley Value) στο πλαίσιο ενός συνεργατικού παιχνιδιού που διοργανώνεται από τον εξυπηρετητή, ώστε να αξιολογήσει την απόδοση των μοντέλων μηχανικής μάθησης που αποστέλλονται από τους χρήστες. Αν θεωρηθεί ότι $W_i^{(t)}$ είναι το μοντέλο που αποστέλλει ο συμμετέχοντα i στον γύρο t της ομόσπονδης μάθησης και $S \subseteq N \cup M, S \neq \emptyset$ ένα οποιοδήποτε μη κενό υποσύνολο όλων των συμμετεχόντων, τότε το αθροιστικό μοντέλο του υποσυνόλου S προκύπτει από τον τύπο:

$$\tilde{W}_s^{(t)} = \sum_{i \in S} \frac{|D_i|}{\sum_{j \in S} |D_j|} W_i^{(t)} \quad (4.3)$$

Το μοντέλο $\tilde{W}_s^{(t)}$ ενός συνόλου χρηστών S για τον γύρο t είναι ουσιαστικά ο σταθμισμένος μέσος όρος των μοντέλων που έστειλαν οι συμμετέχοντες με βάρη ανάλογα του πλήθους δεδομένων που διατέθηκαν για τοπική εκπαίδευση. Αν θεωρήσουμε πως $U(\tilde{W}_S^t)$ είναι η απόδοση του μοντέλου ενός πλήθους $|S|$ χρηστών, ως προς κάποια μετρική, όπως το ποσοστό επιτυχίας (Accuracy), τότε η συνεισφορά του εκάστοτε συμμετέχοντα ορίζεται ως:

$$\phi_i^t = \sum_{S \subseteq (N \cup M) \setminus \{i\}} \frac{U(\tilde{W}_{S \cup \{i\}}^t) - U(\tilde{W}_S^t)}{\binom{|N|+|M|-1}{|S|}} \quad (4.4)$$

Η εξίσωση 4.4 είναι ο ορισμός της τιμής Shapley, όπως αυτή περιγράφηκε στην ενότητα 2.3.4. Υπολογίζει ουσιαστικά την συνεισφορά του χρήστη i , ελέγχοντας αφενός την ακρίβεια του μοντέλου όλων των συνόλων χρηστών S , που δεν περιλαμβάνουν τον i , και αφετέρου την ακρίβεια του ίδιου συνόλου χρηστών με την προσθήκη του συμμετέχοντα i . Σε αυτό το βήμα είναι που εντοπίζονται πιθανές επιθέσεις δηλητηριασμού του συγκεντρωτικού μοντέλου, αφού το τοπικό μοντέλο που αποστέλλει στον διακομιστή ένας κακόβουλος χρήστης δεν έχει καλή απόδοση. Το αποτέλεσμα είναι ότι το δηλητηριασμένο μοντέλο αυτό, μειώνει την επιτυχία κάθε άλλου συγκεντρωτικού μοντέλου στο οποίο συμμετέχει, με αποτέλεσμα η διαφορά $U(\tilde{W}_{S \cup \{i\}}^t) - U(\tilde{W}_S^t)$ να είναι -ως επί το πλείστον- αρνητική.

Αλγόριθμος 1 Shapley Value Contribution Calculation

Για κάθε γύρο $t \leftarrow 1, 2, \dots, R-1$:

/* Υπολογισμός τοπικών μοντέλων */

$W_i^{(t)} \leftarrow \text{ClientUpdate}(i, W^{(t)})$ για κάθε χρήστη $i \in N \cup M$

Για $S \subseteq (N \cup M)$:

$\tilde{W}_S^{(t)} \leftarrow \sum_{i \in S} \frac{|D_i|}{\sum_{j \in S} |D_j|} W_i^{(t)}$

Τέλος Επανάληψης

/* Υπολογισμός Βαθμού συνεισφοράς */

Για $i \leftarrow 1$ έως $(|N| + |M|)$:

$\phi_i^t \leftarrow \sum_{S \subseteq (N \cup M) \setminus \{i\}} \frac{U(\tilde{W}_{S \cup \{i\}}^t) - U(\tilde{W}_S^t)}{\binom{|N|+|M|-1}{|S|}}$

Τέλος Επανάληψης

/* Εφαρμογή Συνάρτησης Softmax */

Για $i \leftarrow 1$ έως $(|N| + |M|)$:

$a_i \leftarrow \frac{e^{\phi_i^t}}{\sum_{j=1}^{|N|+|M|} e^{\phi_j^t}}$

Τέλος Επανάληψης

Τέλος Επανάληψης

Επειδή είναι επιθυμητό ο βαθμός συνεισφοράς των χρηστών να είναι θετικός και όλοι οι βαθμοί να αθροίζονται στο 1, τα ϕ_i^t τίθενται ως όρισμα της συνάρτησης Softmax. Έτσι υπολογίζουμε τον τελικό βαθμό συνεισφοράς a_i^t , στην μορφή που θα χρησιμοποιηθεί για τον υπολογισμό του νέου συγκεντρωτικού μοντέλου και που θα κοινοποιηθεί ύστερα στους χρήστες.

Η συνάρτηση Softmax σε ένα διάνυσμα $\mathbf{z} = [z_1, z_2, \dots, z_k]$ ορίζεται ως:

$$\text{softmax}(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

Οπότε τελικά ο βαθμός συνεισφοράς ενός χρήστη i προκύπτει ως:

$$a_i^t = \frac{e^{\phi_i^t}}{\sum_{j=1}^{|N|+|M|} e^{\phi_j^t}} \quad (4.5)$$

Από το παρόν σημείο και ως εξής, ο βαθμός συνεισφοράς a_i^t θα αναφέρεται κυρίως ως a_i , εκτός αν δεν είναι προφανής ο γύρος t στον οποίο αντιστοιχεί. Συνολικά τα βήματα που χρησιμοποιούνται για τον υπολογισμό της συνεισφοράς των χρηστών μπορούν να αναπαρασταθούν όπως φαίνεται στον Αλγόριθμο 1. Ο αλγόριθμος αυτός είναι ο πιο έγκυρος όσον αφορά στον υπολογισμό της συνεισφοράς των χρηστών, καθώς συγκρίνει το μοντέλο ενός χρήστη, με το άθροισμα οποιουδήποτε άλλου αριθμού χρηστών. Επειδή λοιπόν χρησιμοποιούνται όλα τα πιθανά υποσύνολα του συνόλου $|N|+|M|$ συμμετεχόντων, ο υπολογισμός της τιμής Shapley έχει πολυπλοκότητα $O(2^{(|N|+|M|)})$ σε κάθε γύρο της ομόσπονδης μάθησης, γεγονός που τον καθιστά αρκετά αργό, παρά το μικρό πλήθος συμμετεχόντων που μελετάται στην παρούσα διπλωματική εργασία.

Ως εκ τούτου αναζητήθηκε μια προσεγγιστική λύση πολυπλοκότητας $O(|N| + |M|)$, ώστε να εκτιμηθεί η συνεισφορά των χρηστών, χωρίς να απαιτηθούν $2^{(|N|+|M|)}$ διαφορετικοί υπολογισμοί. Αντί να χρησιμοποιηθούν όλα τα πιθανά υποσύνολα χρηστών, μελετάται μόνο η διαφορά στην ακρίβεια των μοντέλων που περιλαμβάνουν όλους τους συμμετέχοντες, με τα μοντέλα στα οποία δεν έχει αξιοποιηθεί ένας χρήστης. Ταυτόχρονα, για να ενισχυθεί η απόδοση του νέου προσεγγιστικού αλγορίθμου, χρησιμοποιείται μια επιπλέον σύγκριση, μεταξύ του αθροιστικού μοντέλου του προηγούμενου γύρου (W^{t-1}) και του μοντέλου που στάλθηκε από κάθε χρήστη στον παρόν γύρο. Συνεπώς σε έναν γύρο t έχουμε ως συνεισφορά του συμμετέχοντα i :

$$\tilde{\phi}_i^t = (U(\tilde{W}_{(N \cup M)}^t) - U(\tilde{W}_{(N \cup M) - \{i\}}^t)) + (U(W_i^t) - U(W^{t-1})) \quad (4.6)$$

Αλγόριθμος 2 Approximated Shapley Contribution Calculation

Για κάθε γύρο $t \leftarrow 1, 2, \dots, R-1$:

/ Υπολογισμός τοπικών μοντέλων */*

$W_i^{(t)} \leftarrow \text{ClientUpdate}(i, W^{(t)})$ για κάθε χρήστη $i \in N \cup M$

Για $S \subseteq (N \cup M)$ και $|S| = |N| + |M| - 1$:

$\tilde{W}_S^{(t)} \leftarrow \sum_{i \in S} \frac{|D_i|}{\sum_{j \in S} |D_j|} W_i^{(t)}$

Τέλος Επανάληψης

/ Υπολογισμός Βαθμού συνεισφοράς */*

Για $i \leftarrow 1$ έως $(|N| + |M|)$:

$\tilde{\phi}_i^t \leftarrow (U(\tilde{W}_{(N+M)}^t) - U(\tilde{W}_{(N+M) - \{i\}}^t)) + (U(W_i^t) - U(W^{t-1}))$

Τέλος Επανάληψης

/ Εφαρμογή Συνάρτησης Softmax */*

Για $i \leftarrow 1$ έως $(|N| + |M|)$:

$a_i \leftarrow \frac{e^{\tilde{\phi}_i^t}}{\sum_{j=1}^{|N|+|M|} e^{\tilde{\phi}_j^t}}$

Τέλος Επανάληψης

Τέλος Επανάληψης

4.1.3 Contribution Averaging

Αφού υπολογιστούν στον διακομιστή οι βαθμοί συνεισφοράς όλων των συμμετεχόντων στην εκάστοτε εποχή t , απομένει μόνο η εύρεση του νέου συγκεντρωτικού μοντέλου, ώστε να ολοκληρωθεί ο γύρος της ομόσπονδης μάθησης. Ο βασικός αλγόριθμος που χρησιμοποιείται για τον υπολογισμό αυτό αποκαλείται FedAvg (Federated Averaging) και ορίζεται ως [1]:

$$\mathbf{w}^{(t+1)} = \sum_{i \in \text{NUM}} \frac{|D_i|}{|D|} \mathbf{w}_i^t \quad (4.7)$$

Η παρούσα εργασία προτείνει την χρήση ενός εναλλακτικού αλγόριθμου υπολογισμού του συγκεντρωτικού μοντέλου, τον αλγόριθμο Contribution Averaging (ContrAvg), ο οποίος αξιοποιεί την γνώση που προκύπτει κατά το στάδιο υπολογισμού συνεισφοράς (εξισώσεις 4.3, 4.4 και 4.5). Με βάση αυτόν, το συγκεντρωτικό μοντέλο που προκύπτει υπολογίζεται ως:

$$\mathbf{w}^{(t+1)} = \sum_{i \in \text{NUM}} a_i \mathbf{w}_i^t \quad (4.8)$$

Το πλεονέκτημα του αλγορίθμου ContrAvg είναι ότι το μοντέλο που προκύπτει επηρεάζεται περισσότερο από τα καλύτερα τοπικά μοντέλα και λιγότερο από μοντέλα που έχουν χαμηλότερη συνεισφορά. Έτσι μειώνεται η επίδραση των κακόβουλων χρηστών στην εκπαίδευση των νευρωνικών δικτύων, με αποτέλεσμα να αυξάνεται η συνολική επίδοση της διαδικασίας ομόσπονδης μάθησης.

Αφού υπολογισθεί το συγκεντρωτικό μοντέλο του νέου γύρου $t + 1$, αποστέλλεται στους χρήστες, μαζί με τους βαθμούς συνεισφοράς που επέτυχαν, έτσι ώστε να ξεκινήσει η νέα εποχή ομόσπονδης μάθησης.

4.1.4 Βελτίωση των αλγορίθμων FedAvg και ContrAvg

Επειδή η παρούσα διπλωματική εργασία μελετάει και επιθέσεις παρεμβολών, υπάρχει ο κίνδυνος σε κάποιους γύρους της ομόσπονδης μάθησης να μην καταφέρει ένα πλήθος συμμετεχόντων να στείλει το μοντέλο του με αρκετή ισχύ, ώστε να παρακάμψει τις παρεμβολές. Λαμβάνοντας υπόψιν τον μικρό αριθμό χρηστών που μελετάται στην συγκεκριμένη μοντελοποίηση της ομόσπονδης μάθησης, η απώλεια ιδιοτελών χρηστών ευνοεί τον κακόβουλο χρήστη, αφού το δηλητηριασμένο μοντέλο του καταλήγει να επηρεάζει περισσότερο το αθροιστικό μοντέλο.

Για παράδειγμα, αν θεωρήσουμε 5 συμμετέχοντες, εκ των οποίων ο ένας είναι κακόβουλος, με βάση τον αλγόριθμο FedAvg το αθροιστικό μοντέλο θα προκύψει κατά 1/5 από το δηλητηριασμένο μοντέλο. Ωστόσο αν σε κάποια άλλη εποχή της Ομόσπονδης Μάθησης δύο χρήστες δεν καταφέρουν να στείλουν στον εξυπηρετητή λόγω παρεμβολών, τότε το δηλητηριασμένο μοντέλο του κακόβουλου χρήστη θα επηρεάσει κατά 1/3 το αθροιστικό.

Προκειμένου να αυξήσουμε την ανθεκτικότητα των αλγορίθμων FedAvg (4.7) και ContrAvg (4.8) ο διακομιστής αποθηκεύει τα μοντέλα που έστειλαν οι χρήστες στους προηγούμενους γύρους, με αποτέλεσμα να αξιοποιήσει το τελευταίο μοντέλο που έστειλαν επιτυχώς, αν σε κάποια εποχή οι χρήστες εκτεθούν σε παρεμβολές. Συνεπώς το μοντέλο που αποκαλούμε W_i^t στις προηγούμενες ενότητες, πρόκειται είτε για το τοπικό μοντέλο του χρήστη i στον γύρο t (αν δεν υπήρξαν παρεμβολές), είτε για το αμέσως προηγούμενο μοντέλο του συμμετέχοντα i που μεταδόθηκε επιτυχώς.

Φυσικά, μια τέτοια παραδοχή προϋποθέτει ότι δεν υπάρχει κίνδυνος για επιθέσεις free-rider, στις οποίες κάποιοι συμμετέχοντες επιλέγουν να μην συμμετέχουν σε μερικούς γύρους της Ομόσπονδης Μάθησης

[21]. Το κίνητρο μιας τέτοιας επίθεσης είναι συνήθως η εξοικονόμηση πόρων, αφού αποφεύγεται η διαδικασία εκπαίδευσης ενός τοπικού μοντέλου καθώς και η εκπομπή αυτού στον διακομιστή. Η παραδοχή ότι κανένας χρήστης δεν θα απέχει σκόπιμα από έναν γύρο της ομόσπονδης Μάθησης είναι ιδιαίτερα συνηθισμένη σε Cross-Silo δίκτυα, που έχουν μικρό πλήθος συμμετεχόντων.

4.2 Τηλεπικοινωνιακό Μοντέλο

Η παρούσα διπλωματική εργασία χρησιμοποιεί την τεχνική μη ορθογώνιας πολλαπλής πρόσβασης (Non Orthogonal Multiple Access - NOMA) για την ασύρματη επικοινωνία των χρηστών με τον εξυπηρετητή. Σύμφωνα με αυτή την τεχνική πολυπλεξίας, οι χρήστες μεταδίδουν τα τοπικά τους μοντέλα στην ίδια μπάντα συχνοτήτων, εύρους ζώνης B [Hz].

Αν θεωρήσουμε γνωστές και σταθερές τις θέσεις των συμμετεχόντων στο δίκτυο ομόσπονδης μάθησης, τότε μπορούμε να ταξινομήσουμε τα κέρδη καναλιού (channel gain) των χρηστών σε αύξουσα σειρά $G_1 \leq G_2 \leq \dots \leq G_{|N|+|M|}$. Έτσι με την βοήθεια της φόρμουλας του Shannon, ο ρυθμός μετάδοσης (data rate) που επιτυγχάνει ένας χρήστης i , υπολογίζεται ως εξής:

$$R_i = B \log_2 \left(1 + \frac{G_i p_i}{\sum_{j=1}^{i-1} (G_j p_j) + I_0 B} \right) \quad (4.9)$$

όπου με p_i [Watts] συμβολίζουμε την ισχύ με την οποία μεταδίδει το σήμα του ο χρήστης i και με I_0 [dBm/Hz] την πυκνότητα Ενεργειακού Φάσματος του προσθετικού λευκού γκαουσιανού θορύβου (Additive White Gaussian Noise - AWGN) με μηδενικό μέσο όρο.

4.2.1 Υπολογισμός Θεμελιωδών Μεγεθών

Με βάση το τηλεπικοινωνιακό μοντέλο NOMA, ο χρόνος και η ενέργεια που απαιτείται από τους χρήστες για την μετάδοση των $Z(\mathbf{W}_i)$ bits, που περιλαμβάνουν τα βάρη του τοπικού μοντέλου W_i^t , υπολογίζονται ως εξής:

$$\text{Ο χρόνος} \quad T_i = \frac{Z(\mathbf{W}_i)}{R_i} \quad (4.10)$$

$$\text{Η ενέργεια} \quad E_i = p_i * T_i = \frac{p_i * Z(\mathbf{W}_i)}{R_i} \quad (4.11)$$

4.2.2 Κατώφλι Παρεμβολών

Προτού μελετηθεί ο τρόπος διαχείρισης των επιθέσεων παρεμβολών είναι σημαντικό να διευκρινιστεί το πότε θεωρείται μια επίθεση παρεμβολής επιτυχής. Η παρούσα εργασία, εμπνευσμένη από την δουλειά των [22], χρησιμοποιεί ένα κατώφλι παρεμβολών, προκειμένου να αξιολογήσει κατά πόσο μια εκπομπή είναι επιτυχής, δηλαδή κατά πόσο το σήμα που φτάνει στον διακομιστή έχει αρκετή ισχύ ώστε να αποκοδικοποιηθεί ορθά. Θέτοντας λοιπόν ως P_{tol} την ελάχιστη απαιτούμενη διαφορά σε ισχύ που απαιτείται για την ορθή μετάδοση ενός σήματος, τότε η συνάρτηση που πρέπει να πληροί η ισχύς μετάδοσης του χρήστη i , προκύπτει από την εξίσωση:

$$p_i \gamma_i - \sum_{j=0}^{i-1} (p_j \gamma_j) \geq P_{tol} \quad (4.12)$$

Όπου γ_i είναι το κανονικοποιημένο κέρδος καναλιού $\gamma_i = \frac{G_i}{I_0 B}$. Φυσικά για τον τύπο 4.12 έχει θεωρηθεί, όπως και στο 4.9, ότι οι χρήστες είναι ταξινομημένοι σε αύξουσα σειρά κέρδους καναλιού $G_1 \leq G_2 \leq \dots \leq G_{|N|+|M|}$ ή ισοδύναμα $\gamma_1 \leq \gamma_2 \leq \dots \leq \gamma_{|N|+|M|}$.

Σε κάθε γύρο της Ομόσπονδης Μάθησης, ο εκάστοτε συμμετέχοντας i , μπορεί να υπολογίσει την χαμηλότερη τιμή ισχύος, με την οποία η μετάδοση θα είναι επιτυχής, επιλύοντας την εξίσωση 4.12 ως προς την μεταβλητή απόφασης p_i .

$$p_i \geq \frac{\sum_{j=0}^{i-1} (p_j^{t-1} \gamma_j) + P_{tol}}{\gamma_i} \triangleq P_{thres}(\mathbf{p}_{-i}^{t-1}) \quad (4.13)$$

όπου p_i^{t-1} είναι η ισχύς που επιλέγεται από τον i για μετάδοση, στον γύρο $t-1$ της ομόσπονδης μάθησης και πως $\mathbf{p}_{-i}^{t-1} = [p_0^{t-1}, \dots, p_{i-1}^{t-1}, p_{i+1}^{t-1}, \dots, p_{|N|+|M|}^{t-1}]$ οι τιμές ισχύος που επέλεξαν οι υπόλοιποι χρήστες στον προηγούμενο γύρο.

4.3 Διεξαγωγή Μπεϋζιανού Παιγνίου

4.3.1 Η συνάρτηση χρησιμότητας

Στις προηγούμενες ενότητες του κεφαλαίου 4, μελετήθηκε ο τρόπος να αντιμετωπιστεί η επίθεση δηλητηριασμού που πιθανώς να επιχειρήσει ένας κακόβουλος αντίπαλος. Μέσω ενός συνεργατικού παιγνίου που διεξάγεται στον εξυπηρετητή, επιτυγχάνεται η αξιολόγηση των μοντέλων που αποστέλλουν οι χρήστες σε κάθε γύρο, έτσι ώστε να μετριάσει ο αντίκτυπος που έχει ένα κακόβουλο μοντέλο στον σχηματισμό του συγκεντρωτικού. Ωστόσο η παραπάνω μελέτη δεν αρκεί για να αποτρέψει τις επιθέσεις παρεμβολών, που πραγματοποιούνται αν ένας χρήστης επιλέξει να στείλει το μήνυμά του με πολύ μεγαλύτερη ισχύ, από όση απαιτείται. Προκειμένου να λυθεί το πρόβλημα αυτό σχεδιάστηκε ένα μη συνεργατικό Μπεϋζιανό παίγνιο, όπου το όφελος του κάθε χρήστη προκύπτει μέσα από την συνάρτηση χρησιμότητάς του (utility function). Η συνάρτηση αυτή είναι καθολική, δεν διαφοροποιείται δηλαδή ανάλογα με το αν ένας χρήστης είναι κακόβουλος ή όχι, και αντικατοπτρίζει το κέρδος ενός συμμετέχοντα, ανάλογα με την ισχύ εκπομπής του:

$$U_i(p_i^t, \mathbf{p}_{-i}^{t-1}) = \begin{cases} 0 & \text{αν } P_{thres} > P_{max}, \\ c_1 a_i^{t-1} + c_2 \frac{T - T_i(p_i^t, \mathbf{p}_{-i}^{t-1})}{T} - c_3 p_i^t & \text{αλλιού.} \end{cases} \quad (4.14)$$

Πρώτου ερμηνευτεί η παραπάνω συνάρτηση χρησιμότητας, είναι σημαντικό επισημανθεί τι θεωρείται γνωστό σε κάθε γύρο, από τον εκάστοτε χρήστη. Όταν σε έναν γύρο t καλείται ο χρήστης i να επιλέξει την τιμή p_i^t που μεγιστοποιεί την συνάρτηση 4.14, δεν μπορεί να γνωρίζει με ποιες τιμές ισχύος θα μεταδώσουν οι υπόλοιποι χρήστες το μήνυμά τους ή τι αξιολόγηση a_i^t θα λάβει το μοντέλο που θα στείλει. Το μόνο που γνωρίζει είναι τα αποτελέσματα του προηγούμενου γύρου $t-1$, με βάση τα οποία μπορεί να εκτιμήσει τι ρυθμό μεταφοράς T_i θα έχει και ποιο είναι το κατώτατο όριο για επιτυχή εκπομπή P_{thres} .

Η παραπάνω συνάρτηση χρησιμότητας 4.14 επιλέχθηκε, καθώς πληροί κάποιες βασικές προϋποθέσεις

1. Η συνάρτηση χρησιμότητας 4.14 αυξάνεται όσο όταν ένας χρήστης στέλνει ποιοτικό μοντέλο στον εξυπηρετητή. Αυτό είναι προφανές σαν ζητούμενο για τον μέσο χρήστη που επιθυμεί να εκπαιδεύσει το νευρωνικό δίκτυο, ωστόσο εξυπηρετεί και το συμφέρον ενός κακόβουλου συμμετέχοντα, καθώς

Πίνακας 4.1: Επεξήγηση συμβολισμών

Συμβολισμός	Ερμηνεία
c1, c2, c3	σταθερές
a_i^t	η αξιολόγηση του μοντέλου του i από τον εξυπηρετητή, στον γύρο t (4.5, 4.6)
T_i	ο χρόνος μετάδοσης του χρήστη i , δεδομένων των ισχύων όλων των χρηστών (4.10)
T	σταθερά που κανονικοποιεί το $T_i()$

αν το μοντέλο του λάβει τιμή a_i πολύ κοντά στο 0, τότε δεν θα επηρεάσει καθόλου τον σχηματισμό του αθροιστικού μοντέλου. Ως εκ τούτου το αποτέλεσμα θα είναι πολύ παρόμοιο με το να μην είχε μεταδώσει τίποτα ο χρήστης και να μην είχε σπαταλήσει καθόλου πόρους.

2. Η ερμηνεία του όρου $\frac{T-T_i}{T}$ είναι ότι όσο μεγαλύτερος είναι ο χρόνος μετάδοσης T_i τόσο μικρότερη πρέπει να είναι η αμοιβή του χρήστη i .
3. Η συνάρτηση χρησιμότητας είναι κοίλη, με αποτέλεσμα να έχει μοναδική βέλτιστη λύση, για κάποια τιμή p_i . Η απόδειξη βρίσκεται στο τέλος της διπλωματικής, στο Παράρτημα.

4.3.2 Μπεϋζιανό Παίγνιο

Η παρούσα διπλωματική εργασία θεωρεί δύο διαφορετικά είδη συμμετεχόντων στην διαδικασία της Ομόσπονδης μάθησης. Από την μία υπάρχουν οι ιδιοτελείς χρήστες, που επιδιώκουν να μεγιστοποιήσουν το κέρδος τους, ανεξάρτητα από την επιτυχία των υπολοίπων παιχτών και από την άλλη βρίσκονται οι κακόβουλοι χρήστες, που όχι μόνο επιθυμούν να έχουν μεγάλο όφελος, αλλά και να μειώσουν την απόδοση των υπολοίπων παιχτών, με επιθέσεις παρεμβολών. Έστω ότι η μετρική θ_i εκφράζει αν ένας χρήστης είναι ιδιοτελής ή κακόβουλος, όπου

$$\begin{cases} \theta_i = 0 & \text{για τους ιδιοτελείς χρήστες} \\ \theta_i \in [-1, 0) & \text{για τους κακόβουλους χρήστες} \end{cases} \quad (4.15)$$

Προκειμένου να ενσωματωθεί ο τύπος ενός παίχτη στο παίγνιο, υιοθετείται μια τροποποιημένη εκδοχή της συνάρτησης χρησιμότητας, υπό τον συμβολισμό V_i , η οποία συνδυάζει τις αποδόσεις όλων των παιχτών:

$$V_i = U_i + \theta_i \sum_{j \neq i} U_j \quad (4.16)$$

Χάρη στην εξίσωση 4.16 καθίσταται εμφανής η διαφορά μεταξύ ιδιοτελών και κακόβουλων χρηστών. Οι μεν πρώτοι, έχοντας $\theta_i = 0$ επιδιώκουν αποκλειστικά την μεγιστοποίηση της συνάρτησης χρησιμότητάς τους U_i , ανεξαρτήτως των υπολοίπων χρηστών. Οι κακόβουλοι συμμετέχοντες από την άλλη, επηρεάζονται αρνητικά ($\theta_i < 0$) από την απόδοση των υπολοίπων παιχτών, με αποτέλεσμα να επιδιώκουν, μέσω επιθέσεων παρεμβολών, να μειώσουν την συνάρτηση χρησιμότητας των αντιπάλων τους.

Προκειμένου η συνάρτηση χρησιμότητας V_i να είναι κοίλη, πρέπει να επιλεγούν τιμές του θ_i ώστε να επαληθεύεται η εξής ανισότητα:

$$\frac{d^2 U_i}{dp_i^2} + \theta_i \sum_{j \neq i} \frac{d^2 U_j}{dp_i^2} \leq 0, \quad \theta_i \in [-1, 0] \quad \forall i \in \{N \cup M\} \quad (4.17)$$

Οι κατάλληλες τιμές για το θ_i επιλέγονται με πειραματική μελέτη και παρουσιάζονται στην ενότητα 5.

Μπεϋζιανό Παίγνιο σε Κανονική Μορφή

Η ταυτότητα του κάθε χρήστη θ_i είναι γνωστή μόνο από τον ίδιο τον χρήστη και από κανέναν άλλον. Έτσι διαμορφώνεται ένα Μπεϋζιανό παίγνιο ελλιπούς πληροφορίας, το οποίο μπορεί να αναπαρασταθεί με την πλειάδα $G = \{I, \Theta, P, \Psi, V\}$, όπου:

- Σύνολο παιχτών $I = N \cup M = 1, \dots, i, \dots, |N| + |M|$
- Σύνολο τύπων των παιχτών $\Theta = \Theta_1 \times \dots \times \Theta_i \times \dots \times \Theta_{|N|+|M|}$, όπου $\Theta_i = \{\theta_i \in [-1, 0]\}$
- Σύνολο δράσεων των παιχτών $P = P_1 \times \dots \times P_i \times \dots \times P_{|N|+|M|}$, όπου $P_i = \{p_i \in [0, P_{max}]\}$
- Σύνολο πιθανοτήτων $\Psi = \Psi_1 \times \dots \times \Psi_i \times \dots \times \Psi_{|N|+|M|}$, όπου $\Psi_i = Pr\{\psi_i = \theta_i\}$
- Σύνολο Αποδόσεων $U = \{U_1, \dots, U_i, \dots, U_{|N|+|M|}\}$

Με βάση αυτή τη μοντελοποίηση, κάθε χρήστης i επιδιώκει να μεγιστοποιήσει την απόδοσή του, ανάλογα με τον τύπο του και την εκτίμησή του σχετικά με τον τύπο των υπολοίπων παιχτών. Το πρόβλημα βελτιστοποίησης που καλείται να λύσει ο κάθε συμμετέχοντας έχει την ακόλουθη μορφή:

$$\max_{p_i^t} \mathbb{E}[V_i(p_i^t, \mathbf{p}_{-i}^{t-1})] = U_i(p_i^t, \mathbf{p}_{-i}^{t-1}) + \theta_i \sum_{i \neq j} (U_j(p_i^t, \mathbf{p}_{-i}^{t-1}) * Pr[\theta_j = 0 | \theta_i]) \quad (4.18)$$

όπου η συνάρτηση $U_j(p_i^t, \mathbf{p}_{-i}^{t-1})$ αναφέρεται στην εκτίμηση του χρήστη i για την συνάρτηση χρησιμότητας ενός συμμετέχοντα j , στον νέο γύρο t . Σε αυτή την εκτίμηση είναι γνωστή μόνο η ισχύς εκπομπής των χρηστών (p_{-i}^{t-1}) στον προηγούμενο γύρο, καθώς και η μεταβλητή απόφασης του i (p_i^t). Η μοντελοποίηση αυτή επιτρέπει σε έναν κακόβουλο χρήστη να εκτιμήσει τι ζημιά θα προκαλέσει στην απόδοση ενός αντιπάλου δυναμικά, επιτρέποντας έτσι επιθέσεις παρεμβολών.

4.3.3 Υπολογισμός Πιθανότητας Ιδιοτελούς χρήστη

Η εξίσωση 4.18 που μελετήθηκε στην προηγούμενη υποενότητα προϋποθέτει τον υπολογισμό της δεσμευμένης πιθανότητας ενός χρήστη να είναι ιδιοτελής $Pr[\theta_j = 0 | \theta_i]$. Για να πραγματοποιήσει ένας συμμετέχοντας i μια τέτοια εκτίμηση, χρησιμοποιεί τις εξής δύο πληροφορίες που του είναι γνωστές:

- Τον τύπο που έχει ο ίδιος (θ_i)
- Τη απόδοσή του μοντέλου του στον προηγούμενο γύρο a_i^{t-1} καθώς και την απόδοση του χρήστη j (a_j^{t-1}), η οποία κοινοποιήθηκε δημόσια από τον εξυπηρετητή στο τέλος της εποχής $t-1$.

Υποθέτοντας ότι ένας κακόβουλος χρήστης i κάνει επιθέσεις δηλητηριασμού του αθροιστικού μοντέλου, είναι αναμενόμενο να λαμβάνει χαμηλότερες τιμές a_i συγκριτικά με τους υπόλοιπους συμμετέχοντες. Αν λοιπόν ένας χρήστης j ξέρει ότι είναι κακόβουλος ($\theta_j < 0$) και στον προηγούμενο γύρο αξιολογήθηκε το μοντέλο του με a_j , τότε είναι αναμενόμενο να θεωρεί πολύ πιθανό ότι ο i είναι επίσης κακόβουλος, αν $a_i < a_j$. Αντίστροφα, αν δύο ιδιοτελείς συμμετέχοντες a και b λαμβάνουν τιμές a_a, a_b , με $a_a < a_b$, τότε είναι αναμενόμενο ο a , γνωρίζοντας ότι ο ίδιος δεν έκανε επίθεση δηλητηριασμού και έλαβε a_a , να θεωρεί πως ούτε ο b έδρασε κακόβουλα, αφού το μοντέλο του κρίθηκε πιο επιτυχές από του ίδιου.

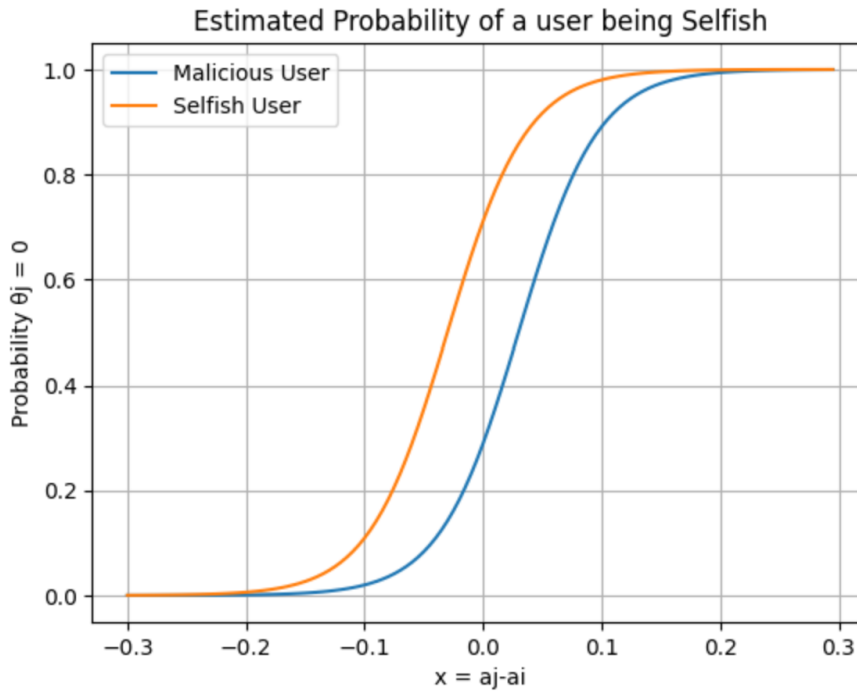
Με αφορμή αυτή την απλή παρατήρηση, μοντελοποιήθηκε η δεσμευμένη πιθανότητα ενός χρήστη να είναι κακόβουλος με την χρήση της σιγμοειδούς συνάρτησης υπερβολικής εφασπτομένης:

$$f(x, \mu) = \frac{1}{2} \left(\frac{e^{a(x+\mu)} - e^{-a(x+\mu)}}{e^{a(x+\mu)} + e^{-a(x+\mu)}} + 1 \right) \quad (4.19)$$

όπου a μία σταθερά. Η μεταβλητή μ αντικατοπτρίζει τον τύπο του χρήστη και παίρνει είτε θετικές τιμές αν $\theta_i = 0$ είτε αρνητικές για κακόβουλους χρήστες. Έτσι ο υπολογισμός της πιθανότητας είναι ο εξής:

$$Pr[\theta_j = 0 | \theta_i] = f(a_j - a_i, \pm\mu) \quad (4.20)$$

Οι τιμές που χρησιμοποιήθηκαν για την εκτίμηση των πιθανοτήτων προέκυψαν από πειραματική μέλητη. Με βάση την αρχικοποίηση $\mu = \pm 0.03$, $a = 15$, τα αποτελέσματα της συνάρτησης f φαίνονται στην γραφική παράσταση 4.2 και στον πίνακα 4.2



Σχήμα 4.2: Γραφική αναπαράσταση των πιθανοτήτων με $\mu = \pm 0.03$, $a = 15$

Πίνακας 4.2: Πίνακας τιμών εκτίμησης πιθανοτήτων με $\mu = \pm 0.03$, $a = 15$

Χρήστης i	Διαφορά $a_j - a_i$	$Pr[\theta_j = 0 \theta_i]$
Ιδιοτελής	0.05	91.68%
Ιδιοτελής	-0.02	57.44%
Κακόβουλος	0.05	64.56%
Κακόβουλος	-0.02	18.24%

Διεξαγωγή Πειραμάτων

Πρωτού παρουσιαστούν τα αποτελέσματα της εργασίας, είναι απαραίτητο να διευκρινιστούν οι παράμετροι με τις οποίες εκτελέστηκαν τα πειράματα. Συνεπώς στο κεφάλαιο 5 πραγματοποιείται αναφορά στις τιμές των σταθερών που επιλέχθηκαν, στα νευρωνικά δίκτυα που χρησιμοποιήθηκαν, στον τρόπο που αρχικοποιήθηκαν οι χρήστες στον χώρο και σε λοιπές πληροφορίες ή επιλογές που λήφθηκαν κατά την προσομοίωση των δικτύων Ομόσπονδης Μάθησης.

5.1 Επιλογή τιμών σταθερών μεταβλητών

Το τηλεπικοινωνιακό πλαίσιο NOMA που χρησιμοποιήθηκε στην παρούσα εργασία δεν επιτρέπει την ταυτόχρονη υποστήριξη περισσότερων από 5 χρηστών, συνεπώς τα πειράματα περιορίστηκαν στην ύπαρξη τεσσάρων (4) ιδιοτελών χρηστών και ενός (1) κακόβουλου συμμετέχοντα, ο οποίος μπορεί να προβεί σε επιθέσεις παρεμβολών ή δηλητηριασμού. Στον πίνακα 5.1 παρουσιάζονται οι τιμές που δόθηκαν στις σταθερές του συστήματος.

Πίνακας 5.1: Αρχικοποίηση Σταθερών

Μεταβλητή	Ερμηνεία	Αρχικοποίηση
$ N + M $	Πλήθος ιδιοτελών συν κακόβουλων χρηστών	5 (4 + 1)
P_{max}	Η μέγιστη ισχύς μετάδοσης σημάτων των χρηστών	1 Watt
(c_1, c_2, c_3)	Σταθερές της συνάρτησης χρησιμότητας (4.14)	(10, 1, 2)
T	Σταθερά κανονικοποίησης του χρόνου που απαιτείται T_i στην 4.14	0.1 sec
B	Εύρος ζώνης του συστήματος	20MHz
P_{tol}	Ελάχιστη διαφορά ισχύος για μετάδοση χωρίς παρεμβολές 4.12, 4.13)	-94dBm
I_0	Πυκνότητα ενεργειακού φάσματος AWGN με μηδενικό μέσο όρο	-174 dBm/Hz

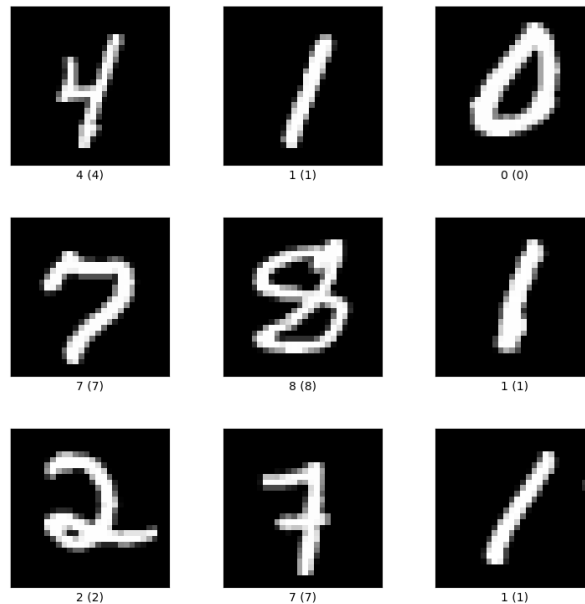
5.2 Μοντέλα Μηχανικής Μάθησης και Παράμετροι εκπαίδευσης

Η παρούσα διπλωματική εργασία βασίστηκε στον κώδικα του [50] για την δημιουργία των μοντέλων μηχανικής μάθησης και την εκπαίδευση αυτών, έτσι ώστε τα αποτελέσματα να συνάδουν με την διεθνή βιβλιογραφία.

5.2.1 MNIST Dataset

Ως σύνολο δεδομένων των χρηστών αλλά και του εξυπηρετητή χρησιμοποιήθηκε το MNIST dataset, το οποίο αποτελείται από μια μεγάλη συλλογή χειρόγραφων ψηφίων με αριθμούς από το 0 έως το 9. Τα ψηφία αυτά παρέχονται σε μορφή ασπρόμαυρων εικόνων, μεγέθους 28x28 pixels. Αποτελείται από 70.000 εικόνες, από τις οποίες 60.000 διαμοιράζονται ισοπόσως στους συμμετέχοντες της μηχανικής μάθησης για την εκπαίδευση των τοπικών μοντέλων και στον διακομιστή για την αξιολόγηση αυτών. Οι υπόλοιπες 10.000 εικόνες χρησιμοποιούνται για την αξιολόγηση των αθροιστικών μοντέλων που παράγονται στο τέλος κάθε εποχής και αποτελούν το σύνολο ελέγχου (test set).

Το MNIST θεωρείται από τα πλέον διαδεδομένα σύνολα δεδομένων και χρησιμοποιείται σαν σημείο αναφοράς για σύγκριση διαφορετικών μοντέλων μηχανικής μάθησης, στο πλαίσιο προβλημάτων ταξινόμησης. Η απλότητά του, το σχετικά μικρό μέγεθός του και η ευκολία χρήσης του, το καθιστούν ένα από τα πλέον χρησιμοποιημένα σύνολα δεδομένων για αξιολόγηση των μοντέλων Ομόσπονδης Μάθησης.



Σχήμα 5.1: Μερικές από τις εικόνες του συνόλου MNIST [51]

5.2.2 Συνελικτικό Νευρωνικό Δίκτυο

Τα τοπικά μοντέλα που εκπαιδεύουν οι χρήστες είναι συνελικτικά νευρωνικά δίκτυα (convolutional neural networks - CNN) για προβλήματα ταξινόμησης. Η αρχιτεκτονική τους αποτελείται από 2 συνελικτικά επίπεδα (convolutional layers) με 32 φίλτρα το κάθε ένα, ακολουθούμενα από συναρτήσεις ενεργοποίησης ReLu. Στο πρώτο συνελικτικό επίπεδο χρησιμοποιείται η τεχνική "Max Pooling", για την εξαγωγή της μέγιστης τιμής (αντί για τον μέσο όρο), σε ένα μπλοκ διαστάσεων 2x2. Αμέσως μετά το επίπεδο αυτό χρησιμοποιείται ένα επίπεδο dropout, για να αποφευχθεί το ενδεχόμενο υπερπροσαρμογής (overfitting), με ρυθμό dropout 0.25. Το αποτέλεσμα των δύο αυτών διαδοχικών επιπέδων επιπεδοποιείται (flatten) και δίνεται ως είσοδος σε ένα πυκνό επίπεδο με 128 νευρώνες και συνάρτηση ενεργοποίησης RELU. Έπειτα

ακολουθεί ένα ακόμα dropout επίπεδο, αυτή τη φορά με ρυθμό 0.5 και το συνελκτικό νευρωνικό δίκτυο ολοκληρώνεται με ένα επίπεδο 10 νευρώνων, από τους οποίους προκύπτει το αποτέλεσμα των 10 διαφορετικών κλάσεων που μελετώνται στο MNIST. Κάθε νευρώνας δηλαδή αντιπροσωπεύει έναν εκ των 10 διαφορετικών αριθμών από το 0 έως το 9 και έχει συνάρτηση ενεργοποίησης Softmax.

5.2.3 Παράμετροι Εκπαίδευσης

Όπως αναφέρθηκε σε προηγούμενες ενότητες, οι συμμετέχοντες στο δίκτυο Ομόσπονδης Μάθησης εκπαιδεύουν τοπικά το μοντέλο που τους έστειλε ο εξυπηρετητής, στον προηγούμενο γύρο, και αποστέλλουν το αποτέλεσμα πίσω στον διακομιστή. Προκειμένου να επισπευσθεί ο χρόνος ολοκλήρωσης ενός γύρου της Ομόσπονδης Μάθησης, οι χρήστες εκπαιδεύουν το μοντέλο τους μόνο για μια εποχή, με βήμα μάθησης 0.001.

Η Ομόσπονδη Μάθηση επιλέχθηκε να ολοκληρώνεται μετά το πέρας ενός συγκεκριμένου αριθμού γύρων. Σε περίπτωση που τα δεδομένα των χρηστών είναι σε IID Μορφή, το πλήθος εποχών Ομόσπονδης Μάθησης είναι 50. Αντίθετα αν τα δεδομένα είναι Non-IID τότε το πλήθος εποχών ορίστηκε να είναι 70, καθώς απαιτείται μεγαλύτερο χρονικό διάστημα για να συγκλίνει η απόδοση του συγκεντρωτικού νευρωνικού δικτύου.

5.3 Τα σενάρια που μελετήθηκαν

Η αρχικοποίηση των μοντέλων μηχανικής μάθησης που περιγράφηκαν στην προηγούμενη υποενότητα γίνεται με τυχαίο τρόπο. Αυτό σημαίνει πως δυο διαδοχικές εκτελέσεις του ίδιου κώδικα, με τις ίδιες παραμέτρους, δεν θα οδηγήσει στην δημιουργία δύο ίδιων νευρωνικών δικτύων, με τον ίδιο βαθμό επιτυχίας. Ως εκ τούτου είναι απαραίτητο να εκτελεστεί πολλές φορές η διαδικασία Ομόσπονδης Μάθησης, με σταθερές παραμέτρους, έτσι ώστε το τελικό αποτέλεσμα να προκύψει ως ο μέσος όρος των αποτελεσμάτων των διαφορετικών αυτών εκτελέσεων. Η ενότητα αυτή θα αναλύσει τον τρόπο με τον οποίο εκτελέστηκαν οι πειραματικές δοκιμές, έτσι ώστε να μειωθεί η "τυχειότητα" των αποτελεσμάτων.

Οι περισσότερες παράμετροι που έχουν σημειωθεί μέχρι στιγμής παραμένουν σταθερές καθ' όλη την διάρκεια της πειραματικής μελέτης. Προκειμένου όμως να μελετηθούν σε βάθος οι καινοτομίες της παρούσας εργασίας, αξιολογούνται οι παρακάτω περιπτώσεις:

1. Αν χρησιμοποιείται η συνάρτηση συνάθροισης FedAvg (4.7) ή η συνάρτηση ContrAvg (4.8).

Η FedAvg είναι η βασική συνάρτηση που χρησιμοποιείται στα δίκτυα Ομόσπονδης Μάθησης και αθροίζει τα τοπικά μοντέλα των χρηστών, με βάση το πλήθος των δεδομένων που αξιοποιήθηκε για την εκπαίδευσή τους. Η ContrAvg είναι η συνάρτηση που προτείνουμε, καθώς αθροίζει τα μοντέλα με βάση βαθμό συνεισφοράς τους, όπως αυτός υπολογίζεται με τις εξισώσεις 4.4, 4.5 και 4.6.

2. Αν πραγματοποιούνται επιθέσεις δηλητηριασμού ή όχι

Ο πιο σημαντικός τρόπος να αξιολογηθεί η επιτυχία του αλγορίθμου ContrAvg, είναι να δοκιμαστεί σε περιβάλλον όπου υπάρχουν επιθέσεις δηλητηριασμού από κακόβουλους χρήστες. Ως εκ τούτου είναι αναγκαίο να μελετηθούν οι δύο παραπάνω αλγόριθμοι συνάρτησης τόσο σε συνθήκες με δηλητηριασμένα μοντέλα, όσο και σε κανονικές. Η **επίθεση δηλητηριασμού δεδομένων μη καθαρής ετικέτας** που υλοποιήσαμε, είναι αυτή κατά την οποία ένας κακόβουλος χρήστης αλλάζει την κλάση των δεδομένων του και εκπαιδεύει το τοπικό του μοντέλο με λανθασμένα δεδομένα. Για παράδειγμα θα μπορούσε να θέσει την ετικέτα όλων των εικόνων με τον αριθμό 1, στον αριθμό 2, με

αποτέλεσμα το νευρωνικό δίκτυο που εκπαιδεύει, να μαθαίνει να ταξινομεί τις εικόνες με τον αριθμό 1 ως δείγματα του "2".

3. Αν πραγματοποιούνται επιθέσεις παρεμβολών

Οι επιθέσεις παρεμβολών αντιμετωπίζονται επιτυχώς χάρη στο μπειϋζιανό παίγνιο που διοργανώνεται μεταξύ των χρηστών. Ωστόσο η ύπαρξη ή μη επιθέσεων παρεμβολών επηρεάζει σε μεγάλο βαθμό τα αποτελέσματα του αθροιστικού μοντέλου στους πρώτους γύρους της Ομόσπονδης Μάθησης, προτού συγκλίνει το παίγνιο, καθώς πολλοί χρήστες δεν καταφέρνουν να στείλουν επιτυχώς το μοντέλο τους. Στο σενάριο που επιλέγεται να μην υπάρχει επίθεση παρεμβολών, όλοι οι χρήστες έχουν $\theta_i = 0$.

4. Αν τα δεδομένα έχουν μορφή IID ή Non-IID

Ως Ανεξάρτητα και Όμοια Κατανεμημένη μορφή δεδομένων (Independently and Identically Distributed - IID) θεωρούμε την περίπτωση όπου οι χρήστες έχουν δείγματα από όλες τις ετικέτες. Στο MNIST, για παράδειγμα, έχουμε IID αρχικοποίηση του δικτύου Ομόσπονδης Μάθησης, αν όλοι οι χρήστες έχουν εικόνες με όλους τους αριθμούς 0-9. Αντίθετα, στην Non-IID μορφή δεδομένων, ο κάθε χρήστης έχει δείγματα μόνο για περίπου 2 διαφορετικούς αριθμούς. Ουσιαστικά οι εικόνες του συνόλου δεδομένων τοποθετούνται σε αύξουσα σειρά, από το 0 έως το 9 και στην συνέχεια μοιράζονται ανάμεσα στους συμμετέχοντες. Ο διακομιστής εξαιρείται φυσικά από αυτή την διαδικασία, αφού προκειμένου να αξιολογήσει τα τοπικά μοντέλα, χρειάζεται δεδομένα από όλες τις κλάσεις.

Ανάλογα με το τι επιλογές θα πραγματοποιηθούν, στις παραπάνω περιπτώσεις, σχηματίζονται διαφορετικά σενάρια προς μελέτη. Για παράδειγμα συγκρίνουμε τους αλγόριθμους συνάθροισης ContrAvg από την FedAvg σε IID δεδομένα χρηστών με ταυτόχρονες επιθέσεις δηλητηριασμού και παρεμβολών. Ωστόσο, όπως αναφέρθηκε και στην αρχή της ενότητας, προκειμένου τα αποτελέσματα να μην επηρεάζονται από την τυχαία αρχικοποίηση των παραμέτρων των νευρωνικών δικτύων, κάθε διαφορετικό σενάριο εκτελείται 10 φορές και υπολογίζεται ο μέσος όρος του ποσοστού επιτυχίας. Κάθε τέτοια διαφορετική εκτέλεση, θα αποκαλείται ως **αρχικοποίηση** εφεξής. Οι αρχικοποιήσεις ενός σεναρίου, αν και έχουν όλες τις παραμέτρους και σταθερές ίδιες, διαφέρουν μεταξύ τους ως προς την θέση των χρηστών στον χώρο. Ο τρόπος που επιλέγονται οι θέσεις των συμμετεχόντων, παρουσιάζεται στην επόμενη ενότητα.

5.4 Τοπολογία και αρχικοποίηση χρηστών στον χώρο

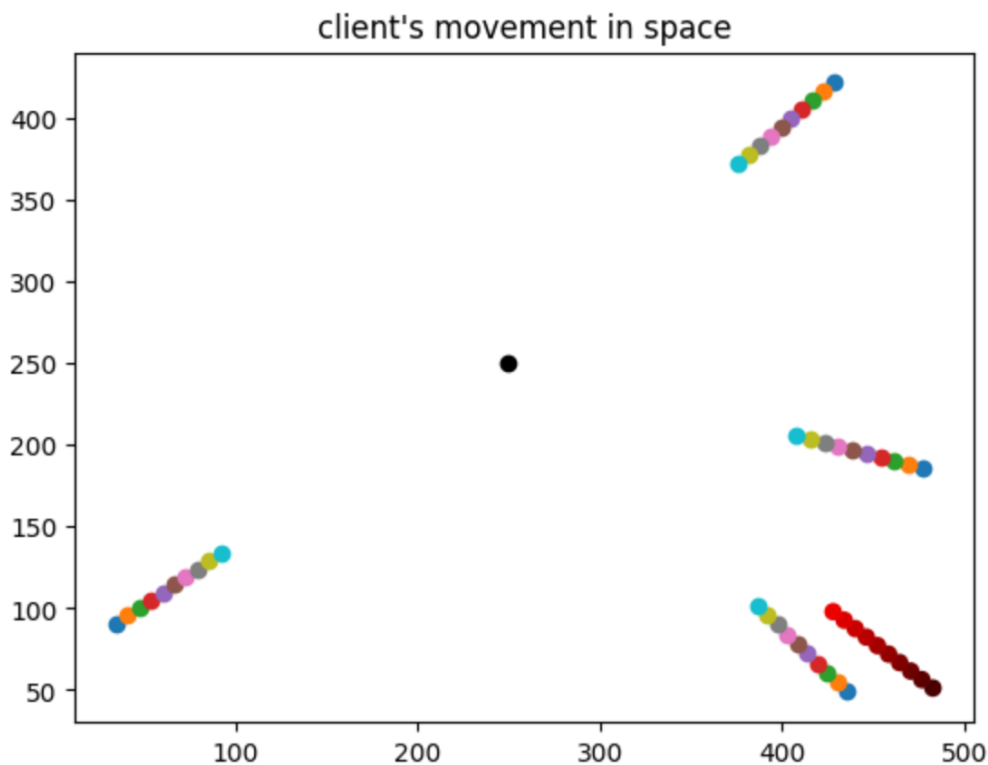
Το τηλεπικοινωνιακό πρωτόκολλο Μη Ορθογωνικής Πολλαπλής Πρόσβασης (NOMA) χρησιμοποιεί τα κέρδη καναλιού των χρηστών, προκειμένου να υπολογίσει την χρονική επιβάρυνση του συστήματος. Το κέρδος καναλιού ενός χρήστη εξαρτάται αποκλειστικά από την θέση του στον χώρο, σε σχέση με την τοποθεσία του διακομιστή με τον οποίο επικοινωνεί. Ως εκ τούτου επιλέχθηκε η διεξαγωγή του δικτύου ομόσπονδης μάθησης σε ένα καρτεσιανό επίπεδο (x, y) όπου $x \in [0, 500]$, $y \in [0, 500]$. Ο εξυπηρετητής που αναλαμβάνει την αξιολόγηση των μοντέλων και την παραγωγή του συγκεντρωτικού μοντέλου τοποθετήθηκε στην μέση του τετράγωνου χώρου, δηλαδή στην θέση $(x_{server}, y_{server}) = (250, 250)$.

Ωστόσο η τοποθέτηση των χρηστών στο επίπεδο δεν μπορεί να είναι τυχαία, εξαιτίας του τρόπου λειτουργίας του πρωτοκόλλου Μη Ορθογωνικής Πρόσβασης (NOMA). Όπως εξηγήθηκε στο θεωρητικό υπόβαθρο 2.2 και στην ενότητα 4.2, μέσω της εξίσωσης 4.9, οι χρήστες ταξινομούνται σε αύξουσα σειρά κέρδους καναλιού, προτού αποκωδικοποιηθούν τα σήματά τους με την διαδικασία SIC. Υπενθυμίζεται ότι, με βάση τον τρόπο λειτουργίας της SIC, ο πρώτος χρήστης (με το μικρότερο κέρδος καναλιού) δεν

δέχεται παρεμβολές από κανέναν άλλον συμμετέχοντα, αλλά το σήμα του παρεμβάλλεται στα σήματα των υπολοίπων. Ομοίως, ο αποστολέας με το δεύτερο μικρότερο κέρδος καναλιού επηρεάζεται μόνο από το σήμα του πρώτου και κάνει παρεμβολές στην μετάδοση όλων των άλλων χρηστών. Τέλος η μετάδοση που έχει το μεγαλύτερο κέρδος καναλιού δεν επηρεάζει κανέναν από τους υπόλοιπους χρήστες, αλλά αισθάνεται παρεμβολές από όλους.

Είναι λοιπόν εμφανές, ότι προκειμένου ένας κακόβουλος χρήστης να είναι σε θέση να κάνει επίθεση παρεμβολών στους υπόλοιπους συμμετέχοντες, δεν θα πρέπει σε καμία περίπτωση να έχει το μεγαλύτερο κέρδος καναλιού. Σε μία τέτοια περίπτωση, με όση ισχύ και να στείλει το μήνυμά του, δεν θα επηρεάσει τον σηματοθορυβικό λόγο των υπολοίπων. Επιλέχθηκε λοιπόν η τοποθέτηση του κακόβουλου χρήστη, σε θέση τέτοια ώστε να έχει το μικρότερο κέρδος καναλιού και να μπορεί να επιτεθεί μέσω παρεμβολών, σε όλους τους υπόλοιπους χρήστες.

Ωστόσο η παραπάνω επιλογή δεν είναι αρκετή ώστε να εξασφαλιστούν επιτυχείς επιθέσεις παρεμβολών από έναν κακόβουλο χρήστη. Σε περιπτώσεις όπου οι ιδιοτελείς χρήστες έχουν σημαντικά μεγαλύτερο κέρδος καναλιού από τον κακόβουλο χρήστη, τότε είναι λογικό πως η ισχύς εκπομπής που θα απαιτείται για παρεμβολές θα είναι πολύ μεγαλύτερη από το 1 Watt που έχει στην διάθεσή του ο κακόβουλος χρήστης. Μια πιο διαισθητική εξήγηση είναι ότι αν όλοι οι χρήστες βρίσκονται κοντά στον εξυπηρετητή (έχουν δηλαδή μεγάλο κέρδος καναλιού), είναι πολύ δύσκολο για έναν συμμετέχοντα που βρίσκεται πολύ μακριά να κάνει επίθεση παρεμβολών, αφού το σήμα των υπολοίπων φτάνει αρκετά καθαρά στον διακομιστή.



Σχήμα 5.2: Διαφορετικές αρχικοποιήσεις των χρηστών στον χώρο

Προκειμένου λοιπόν να μελετηθούν οι χειρότερες δυνατές περιπτώσεις, δηλαδή αυτές στις οποίες ο κακόβουλος χρήστης μπορεί να κάνει επίθεση παρεμβολών που θα επηρεάσει όλους τους συμμετέχοντες, επιλέξαμε μια ψευδοτυχαία αρχικοποίηση χρηστών στον χώρο, που πληρούσε τις παραπάνω προϋποθέσεις. Η αρχικοποίηση που μελετήθηκε τοποθετεί τον κακόβουλο χρήστη στην θέση (489, 46) και τους

υπόλοιπους ιδιοτελείς συμμετέχοντες στις θέσεις (441, 43), (485, 184), (434, 427) και (28, 86). Ξεκινώντας από αυτές τις θέσεις, οι χρήστες μετακινήθηκαν προς την κατεύθυνση του διακομιστή κατά 8 μέτρα, προκειμένου να προκύψουν 9 επιπλέον διαφορετικές αρχικοποιήσεις, στις οποίες τα κέρδη καναλιού ευνοούν επιθέσεις παρεμβολών. Οι 10 διαφορετικές αρχικοποιήσεις των χρηστών που μελετήθηκαν παρουσιάζονται στο διάγραμμα 5.2. Οι αποχρώσεις του κόκκινου αναπαριστούν την θέση του κακόβουλου χρήστη και η μαύρη κουκκίδα στο κέντρο του χώρου είναι η θέση του εξυπηρετητή.

Κάθε διαφορετική αρχικοποίηση αναπαριστάται με έναν διαφορετικό αριθμό από το 0 έως το 9. **Αρχικοποίηση ή βήμα (step) "0"**, σημαίνει ότι οι χρήστες βρίσκονται στις αρχικές θέσεις και δεν έχουν μετακινηθεί προς το κέντρο. Στο διάγραμμα 5.2 αυτή η αρχικοποίηση αποτυπώνεται με σκούρα μπλε κουκκίδα. Ως βήμα "1" χαρακτηρίζεται η αμέσως επόμενη κατάσταση, όπου οι χρήστες μετακινήθηκαν μία φορά (8 μέτρα) προς το κέντρο, που φαίνεται με πορτοκαλί κουκκίδα. Ομοίως, στο βήμα "2", οι χρήστες μετακινούνται προς την τοποθεσία του διακομιστή κατά 8 μέτρα, ξεκινώντας από την θέση που είχαν στην "1^η" αρχικοποίηση. Η διαδοχική αυτή μετακίνηση επαναλαμβάνεται μέχρι το βήμα 9.

Υπολογισμός Κέρδους Καναλιού

Η απώλεια διαδρομής (path loss) ενός σήματος, σε αστικές περιοχές, από έναν χρήστη i προς τον εξυπηρετητή, υπολογίζεται μέσω του ακόλουθου τύπου:

$$PathLoss_i \text{ (in dB)} = 128.1 + 37.6 \log_{10} \left(\frac{\sqrt{(x_{\text{server}} - x_i)^2 + (y_{\text{server}} - y_i)^2}}{1000} \right)$$

Ως κέρδος καναλιού ενός χρήστη i , προς τον διακομιστή, ορίζεται το αντιστρόφως ανάλογο μέγεθος της απώλειας διαδρομής, δηλαδή

$$G_i = \frac{1}{PathLoss_i}$$

5.5 Αρχικοποίηση Μεταβλητών

Οι εξισώσεις χρησιμότητας των χρηστών υπολογίζονται σε κάθε γύρο της ομόσπονδης μάθησης ανάλογα με τα αποτελέσματα της προηγούμενης εποχής. Έτσι, όπως είναι λογικό, απαιτείται αρχικοποίηση των απαραίτητων μεταβλητών, έτσι ώστε να είναι εφικτός ο υπολογισμός των συναρτήσεων χρησιμότητας στην πρώτη εποχή, που δεν υπάρχουν αποτελέσματα προηγούμενου γύρου.

Πίνακας 5.2: Πίνακας Αρχικοποίησης Μεταβλητών

Μεταβλητή	Ερμηνεία	Αρχικοποίηση
$a_i^{t=-1}$	Ο βαθμός συνεισφοράς των χρηστών (ενότητα 4.5)	[0.2, 0.2, 0.2, 0.2, 0.2]
$p_i^{t=-1}$	Η ισχύς μετάδοσης των χρηστών	[0.5, 0.5, 0.5, 0.5, 0.5]

Οι τιμές του βαθμού συνεισφοράς και της ισχύος, όπως παρουσιάζονται στον πίνακα 5.2, είναι συμφωνημένες και ευρέως γνωστές κατά την εκκίνηση της διαδικασίας Ομόσπονδης Μάθησης. Όλοι οι συμμετέχοντες, δηλαδή, χρησιμοποιούν αυτές τις τιμές, κατά τον πρώτο υπολογισμό της συνάρτησης χρησιμότητας στον γύρο 0.

Στο κεφάλαιο αυτό θα μελετηθούν τα αποτελέσματα της μοντελοποίησης της παρούσας εργασίας. Υπενθυμίζεται ότι κάθε σενάριο που μελετάται, περιλαμβάνει 10 διαφορετικές αρχικοποιήσεις χρηστών στον χώρο, από τις οποίες παρουσιάζεται ο μέσος όρος. Ταυτόχρονα πραγματοποιείται και ξεχωριστή μελέτη για το κατά πόσο επηρεάζει μία συγκεκριμένη αρχικοποίηση ενός σεναρίου κάποιες μεταβλητές, όπως την ισχύ μετάδοσης ενός χρήστη i (p_i).

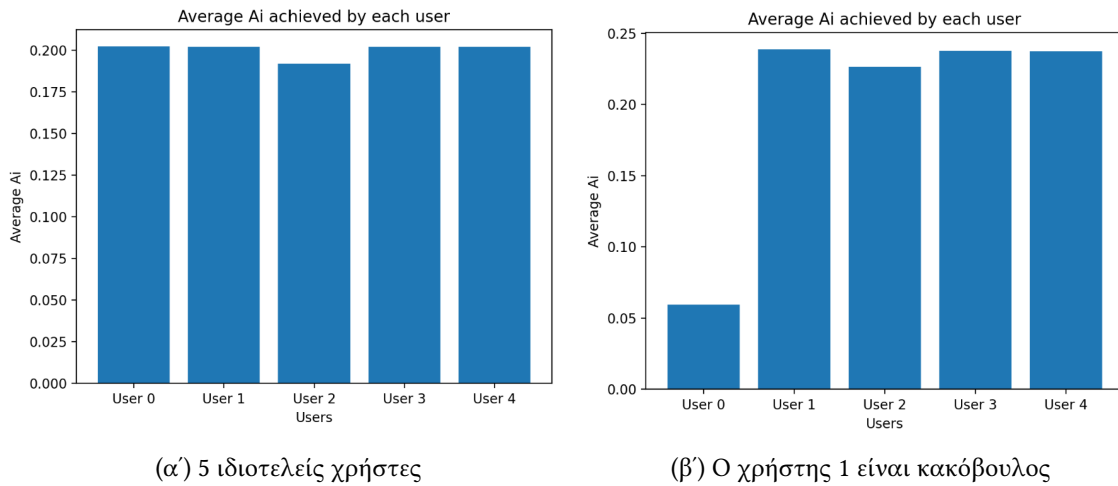
Το κεφάλαιο αυτό χωρίζεται σε τέσσερις βασικές υποενότητες. Στην πρώτη παρουσιάζονται τα αποτελέσματα των αλγορίθμων αξιολόγησης των μοντέλων των χρηστών. Έπειτα μελετάται η σύγκλιση του Μπεϋζιανού παιγνίου που λαμβάνει χώρα μεταξύ των χρηστών και παρουσιάζονται τα αποτελέσματα διαφόρων μετρικών, όπως ο χρόνος και η ενέργεια που απαιτείται σε κάθε γύρο. Τέλος γίνεται σύγκριση του ποσοστού επιτυχίας του συγκεντρωτικού μοντέλου μηχανικής μάθησης, όταν χρησιμοποιείται ο αλγόριθμος συνάθροισης ContrAvg (4.8), που προτείνει η παρούσα διπλωματική, σε αντιδιαστολή με τον ευρέως χρησιμοποιούμενο FedAvg (4.7).

6.1 Εντοπισμός Κακόβουλου Χρήστη

Για τον εντοπισμό του κακόβουλου χρήστη χρησιμοποιούνται οι αλγόριθμοι υπολογισμού συνεισφοράς που μελετήθηκαν στην ενότητα 4.1.2. Ο εκτενής αλγόριθμος υπολογισμού συνεισφοράς (Extended Contribution Calculation) έχει πολυπλοκότητα $O(2^{|N|+|M|})$, γεγονός που καθυστερούσε σημαντικά την διαδικασία Ομόσπονδης Μάθησης. Ως εκ τούτου **επιλέχθηκε η χρήση του προσεγγιστικού αλγορίθμου** (Estimated Contribution Calculation) κατά την εκτέλεση όλων των πειραμάτων που παρουσιάζονται σε αυτή την ενότητα. Ο αλγόριθμος αυτός αποτελεί καινοτομία της παρούσας διπλωματικής εργασίας και έχει χρονική πολυπλοκότητα ίση με $O(|N| + |M|)$, μειώνοντας έτσι σημαντικά τον χρόνο που απαιτείται για να ολοκληρωθεί ένας γύρος ομόσπονδης μάθησης, ακόμα και όταν το πλήθος συμμετεχόντων ($|N| + |M|$) είναι μόλις 5.

Ο προσεγγιστικός αλγόριθμος συγκρίνει το μοντέλο που έχει αποστείλει ο κάθε χρήστης με αυτά των υπολοίπων, καθώς και με το συγκεντρωτικό μοντέλο του προηγούμενου γύρου, προκειμένου να υπολογιστεί η μετρική a_i . Ο κάθε χρήστης λαμβάνει διαφορετικό a_i στο τέλος κάθε γύρου, ωστόσο το άθροισμα όλων των a_i πρέπει να αθροίζει στο 1. Φυσικά όσο μεγαλύτερη τιμή λαμβάνει ένας χρήστης, τόσο πιο αξιόλογο θεωρείται το μοντέλο του, με αποτέλεσμα να συμβάλλει περισσότερο στην διαμόρφωση του τελικού αθροιστικού μοντέλου.

Λαμβάνοντας υπόψιν -συνεπώς- τον τρόπο σχηματισμού του a_i , είναι αναμενόμενο ότι όταν όλοι οι χρήστες είναι ιδιοτελείς, θα λαμβάνουν παρόμοιο βαθμό a_i . Δεδομένου ότι στην πειραματική μελέτη

Σχήμα 6.1: Βαθμός συνεισφοράς a_i σε IID δεδομένα

επιλέχθηκε το πλήθος συμμετεχόντων να είναι ίσο με 5, αναμένεται ότι ο κάθε χρήστης θα επιτυγχάνει $a_i \approx \frac{1}{5} = 0.2$, αν κανένας δεν πραγματοποιεί επιθέσεις δηλητηριασμού. Από την άλλη, σε περιπτώσεις που ο χρήστης 1 είναι κακόβουλος, το ζητούμενο είναι να λαμβάνει τιμή a_i σημαντικά μικρότερη από αυτή που θα λάμβανε αν ήταν ιδιοτελής χρήστης. Στα σχήματα 6.1α', 6.1β', 6.2α', 6.2β' και στον πίνακα 6.1 παρουσιάζεται η μέση τιμή της μετρικής a_i που έλαβε ο εκάστοτε συμμετέχοντας σε όλες τις εποχές και των 10 διαφορετικών αρχικοποιήσεων στον χώρο.

6.1.1 Η μετρική a_i σε IID δεδομένα

Όταν τα δεδομένα που έχουν στην διάθεσή τους οι χρήστες έχουν IID μορφή, τότε ο κάθε χρήστης έχει δείγματα από όλες τις κλάσεις του συνόλου δεδομένων. Στην περίπτωση του συνόλου MNIST, που χρησιμοποιήθηκε στην παρούσα εργασία, κάθε χρήστης έχει εικόνες που αναπαριστούν τους αριθμούς 0 έως 9, οπότε τα τοπικά μοντέλα εκπαιδεύονται με δείγματα από όλες τις διαφορετικές κλάσεις δεδομένων. Αυτό έχει ως αποτέλεσμα τα δεδομένα των χρηστών να είναι "παρόμοια", για αυτό και ο βαθμός a_i των χρηστών ακολουθεί ομοιόμορφη κατανομή γύρω από το 0.2, όπως ήταν αναμενόμενο (σχήμα 6.1α').

Πράγματι, τα αποτελέσματα της προσομοίωσης επιβεβαιώνουν την παραπάνω λογική, όπως φαίνεται στο ραβδόγραμμα του σχήματος 6.1α'. Στο σχήμα αυτό, ο βαθμός a_i των χρηστών φαίνεται να ακολουθεί ομοιόμορφη κατανομή γύρω από το 0.2, με απόκλιση μικρότερη του 0.008. Οι ακριβείς τιμές των βαθμών a_i , φαίνονται στον πίνακα 6.1.

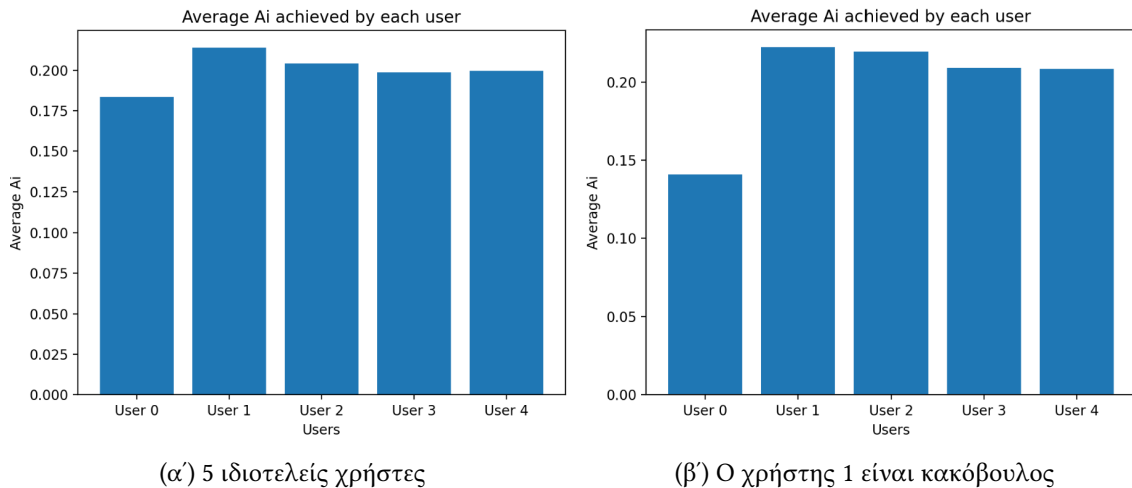
Όσον αφορά στο σενάριο όπου ο χρήστης 0 είναι κακόβουλος και επιχειρεί επιθέσεις δηλητηριασμού του αθροιστικού μοντέλου, οι επιθέσεις του γίνονται αντιληπτές από τον διακομιστή. Μάλιστα ο κακόβουλος χρήστης λαμβάνει τιμή $a_0 \approx 0.06$ την στιγμή που οι υπόλοιποι συμμετέχοντες έχουν κατά μέσο όρο $a_i \approx 0.23$. Η διαφορά αυτή στις τιμές που λαμβάνουν τα μοντέλα δεν είναι απλά πολύ μεγάλη: εξασφαλίζει ότι το αθροιστικό μοντέλο θα σχηματίζεται κατά 94% από τα μοντέλα των ιδιοτελών χρηστών και μόλις κατά 6% από το δηλητηριασμένο μοντέλο. Αυτό φυσικά ισχύει όταν χρησιμοποιείται η συνάρτηση συνάθροισης ContrAvg που αθροίζει τα τοπικά μοντέλα ανάλογα με την μετρική a_i που έλαβαν. Στην πράξη, η διαφορά μεταξύ των δύο μεθόδων συνάθροισης (FedAvg και ContrAvg) παρουσιάζει πολύ σημαντικές διαφορές ως προς το ποσοστό επιτυχίας, γεγονός που οφείλεται στον ορθό εντοπισμό του κακόβουλου χρήστη. Περισσότερες λεπτομέρειες σχετικά με τα ποσοστά επιτυχίας παρουσιάζονται σε επόμενες ενότητες.

Πίνακας 6.1: Τιμές του μέσου βαθμού a_i

Δεδομένα	Χρήστης 0	Χρήστης 1	Χρήστης 2	Χρήστης 3	Χρήστης 4
IID, χωρίς κακόβουλο χρήστη	0.2021	0.2019	0.1925	0.2020	0.20195
IID, με κακόβουλο τον χρήστη 0	0.0596	0.2386	0.2266	0.2378	0.23735
NON IID, χωρίς κακόβουλο χρήστη	0.18365	0.2138	0.20415	0.1987	0.1997
NON IID, με κακόβουλο τον χρήστη 0	0.1409	0.2221	0.2195	0.2091	0.2083

6.1.2 Η μετρική a_i σε NON IID δεδομένα

Η μελέτη των περιπτώσεων όπου οι χρήστες έχουν NON IID δεδομένα είναι αρκετά διαφορετική. Αρχικά οι χρήστες δεν έχουν δεδομένα από όλες τις κλάσεις. Στην περίπτωση του MNIST, το πλήθος των κλάσεων είναι 10 (αριθμοί από το 0 έως το 9) και μοιράζεται σε πέντε διαφορετικούς χρήστες. Αυτό έχει ως αποτέλεσμα ο κάθε συμμετέχοντας να έχει πρόσβαση σε εικόνες που αναπαριστούν 2 ή το πολύ 3 αριθμούς, αφού το σύνολο δεδομένων μοιράζεται ισοπόσως σε όλους τους χρήστες.

Σχήμα 6.2: Βαθμός συνεισφοράς a_i σε NON IID δεδομένα

Όπως καθίσταται εμφανές από σχήμα 6.2α', στην περίπτωση που όλοι οι συμμετέχοντες είναι ιδιοτελείς και δεν υπάρχουν δηλητηριασμένα μοντέλα, ο μέσος βαθμός a_i που λαμβάνουν οι χρήστες ακολουθεί -όπως και προηγουμένως- ομοιόμορφη κατανομή, αλλά με ελαφρώς μεγαλύτερη διασπορά. Η μεγαλύτερη αυτή διασπορά οφείλεται πιθανότατα στο ότι τα δεδομένα των χρηστών διαφέρουν παρασάγγας μεταξύ τους. Σε κάθε γύρο, οι χρήστες εκπαιδεύουν το τοπικό μοντέλο τους έτσι ώστε να αναγνωρίζει μόνο τους δύο (ή το πολύ τρεις) αριθμούς που έχουν στην διάθεσή τους ως εικόνες. Έτσι τα μοντέλα που προκύπτουν είναι πολύ διαφορετικά και η σύγκρισή τους καθίσταται δυσκολότερη. Όσο εύκολο είναι να συγκριθούν δύο παρόμοια μοντέλα, τόσο δύσκολο είναι να συγκριθούν δύο αρκετά διαφορετικά μοντέλα που επιχειρούν να κατηγοριοποιήσουν διαφορετικούς αριθμούς.

Λαμβάνοντας αυτό υπόψιν, είναι εμφανές ότι στο διάγραμμα 6.2α' όλοι οι χρήστες λαμβάνουν παρόμοια τιμή a_i , με κέντρο το 0.2 αλλά και ότι στην περίπτωση που υπάρχει κακόβουλος χρήστης (6.2β') αυτός αναγνωρίζεται με επιτυχία και σημειώνει $a_0 \approx 0.14$. Η τιμή αυτή είναι μικρότερη κατά περίπου 0.07 από την αμέσως μικρότερη και βοηθάει σημαντικά στον περιορισμό των επιπτώσεων μιας επίθεσης δηλητηριασμού.

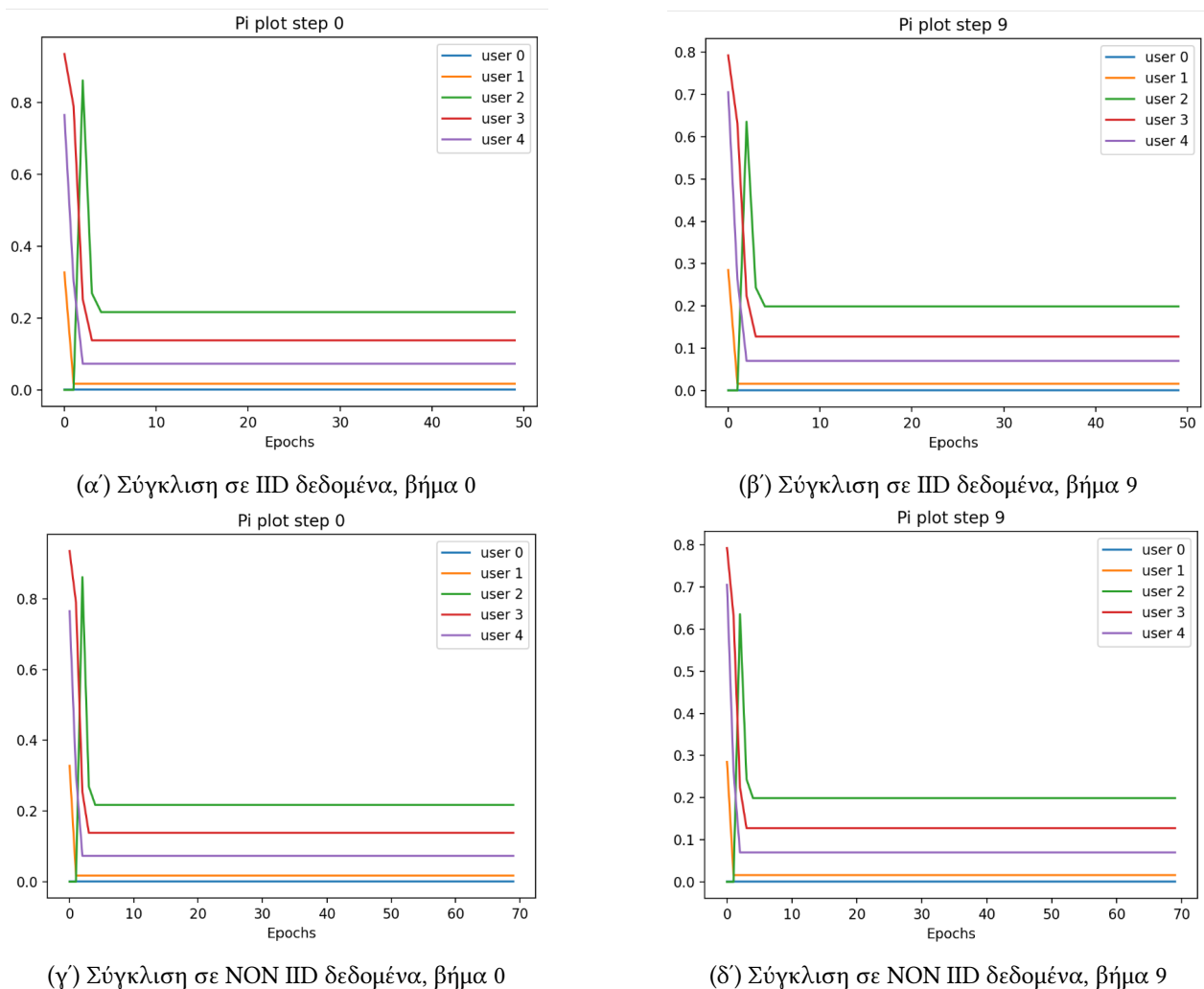
6.2 Σύγκλιση Παιγνίου

Σκοπός αυτής της ενότητας είναι η παρουσίαση των αποτελεσμάτων που αφορούν στην σύγκλιση του Μπεϋζιανού παιγνίου, που διαδραματίζεται μεταξύ των συμμετεχόντων. Σκοπός του παιγνίου είναι η επιλογή της βέλτιστης ισχύος μετάδοσης για τον κάθε χρήστη, έτσι ώστε να εξασφαλιστεί επιτυχής μετάδοση. Η εξίσωση χρησιμότητας των χρηστών (4.14) λαμβάνει υπόψιν της την ισχύ με την οποία μετέδωσαν οι χρήστες τα σήματά τους στην προηγούμενη εποχή. Τέλος, ανάλογα με το αν ένας χρήστης είναι ιδιοτελής ή κακόβουλος, καλείται να λύσει το πρόβλημα βελτιστοποίησης 4.18, επιλέγοντας την τιμή ισχύος που μεγιστοποιεί την απόδοσή του.

Στην ενότητα αυτή πραγματοποιείται διαφορετική μελέτη για κάθε αρχικοποίηση των χρηστών στον χώρο, αφού αυτή επηρεάζει την ελάχιστη απαιτούμενη ενέργεια για επιτυχή μετάδοση δεδομένων.

6.2.1 Αρχικοποιήσεις χωρίς κακόβουλους χρήστες

Τα σενάρια όπου δεν υπάρχουν κακόβουλοι χρήστες παρουσιάζουν σημαντικές ομοιότητες μεταξύ τους, είτε οι χρήστες έχουν IID είτε NON IID δεδομένα. Επειδή δεν υπάρχουν επιθέσεις παρεμβολών, το παίγνιο συγκλίνει μετά από το πολύ 5 γύρους ομόσπονδης μάθησης.



Σχήμα 6.3: Η σύγκλιση της ισχύος μετάδοσης p_i σε σενάρια με 5 ιδιοτελείς χρήστες

Με βάση το σχήμα 6.3 μπορεί εύκολα να διαπιστώσει κανείς ότι όλες οι περιπτώσεις που παρουσιάζ-

στηκαν είναι παρόμοιες. Τα διαγράμματα είναι σχεδόν τα ίδια, αν εξαιρέσει κανείς τις αρχικές τιμές p_i με τις οποίες μεταδίδουν οι συμμετέχοντες το τοπικό τους μοντέλο. Για παράδειγμα παρατηρείται ότι στο βήμα 0, είτε έχουμε IID είτε NON IID δεδομένα, οι αρχική ισχύς μετάδοσης των χρηστών 2 και 3 είναι σημαντικά πάνω από 0.8 Watts, ενώ στην αρχικοποίηση με βήμα 9 είναι αρκετά χαμηλότερη. Αυτό φυσικά οφείλεται στο ότι σε κάθε βήμα, οι χρήστες πλησιάζουν όλο και περισσότερο στον εξυπηρετητή (που βρίσκεται στο κέντρο της αρχικοποίησης), με αποτέλεσμα να μειώνεται η ισχύς μετάδοσης που απαιτείται για ορθή αποκωδικοποίηση του σήματός τους.

Αξίζει να σημειωθεί, ότι ο χρήστης 0, αν και φαίνεται να έχει μηδενική ισχύ μετάδοσης, στην πραγματικότητα έχει απλώς πάρα πολύ μικρή, της τάξης των 0.00057 Watt. Αυτό είναι απολύτως λογικό, αφού αφενός ο χρήστης 0 έχει το μικρότερο κέρδος καναλιού, με αποτέλεσμα να μην αισθάνεται παρεμβολές από κανέναν άλλον χρήστη και αφετέρου, επειδή δεν είναι κακόβουλος σε αυτές τις αρχικοποιήσεις, δεν έχει κανένα κίνητρο να μεταδώσει με σημαντικά μεγαλύτερη ισχύ από όση από όση απαιτείται.

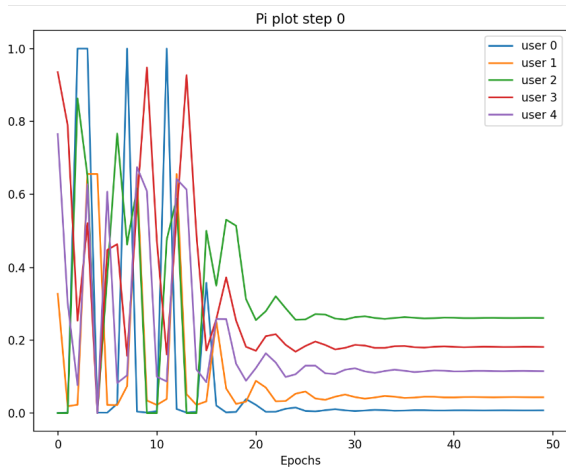
Η ίδια λογική με τα κέρδη καναλιού μπορεί να εφαρμοστεί και για τους υπόλοιπους χρήστες. Ο συμμετέχοντας 1, για παράδειγμα, έχει το δεύτερο μικρότερο κέρδος καναλιού μετά τον χρήστη 0, με αποτέλεσμα να αισθάνεται παρεμβολές μόνο από αυτόν. Ως εκ τούτου, χρειάζεται να μεταδώσει το σήμα του με περισσότερη ισχύ, ώστε να φτάσει επιτυχώς στον διακομιστή. Αμέσως μετά τον χρήστη 1, το μικρότερο κέρδος καναλιού έχει ο 4, ακολουθούμενος από τους χρήστες 3 και 2. Δεν είναι -συνεπώς- παράλογο, όταν κανένας δεν επιχειρεί επίθεση παρεμβολών, η τελική ισχύς μετάδοσης των συμμετεχόντων να είναι ταξινομημένη ανάλογα με το κέρδος καναλιού τους. Από τον χρήστη 0 που δεν αισθάνεται παρεμβολές, άρα δεν έχει μεγάλες απαιτήσεις σε ισχύ, μέχρι τον χρήστη 2, που έχοντας το υψηλότερο κέρδος, αναγκάζεται να μεταδώσει με την μεγαλύτερη ισχύ, ώστε να αντιπαραβάλλει την επίδραση των υπόλοιπων σημάτων.

6.2.2 Αρχικοποιήσεις με κακόβουλο χρήστη σε IID δεδομένα

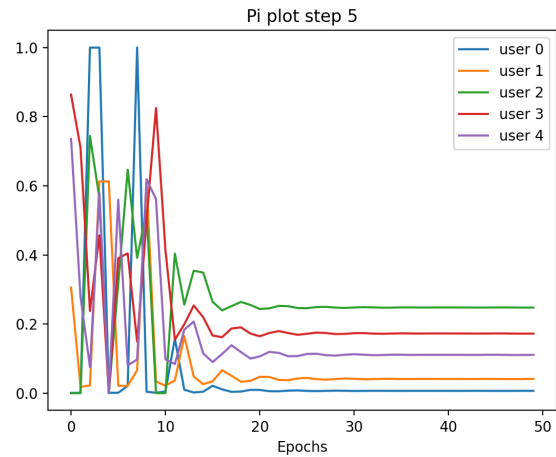
Τα αποτελέσματα είναι πολύ διαφορετικά όταν στο δίκτυο Ομόσπονδης Μάθησης συμμετέχει ένας κακόβουλος χρήστης που επιχειρεί αφενός να δηλητηριάσει το συγκεντρωτικό μοντέλο και αφετέρου να προκαλέσει παρεμβολές στην επικοινωνία των συμμετεχόντων με τον διακομιστή. Οι εξισώσεις 4.12 και 4.13 υπολογίζουν την ελάχιστη ισχύ με την οποία πρέπει να μεταδώσει ένας χρήστης το σήμα του, προκειμένου αυτό να αποκωδικοποιηθεί ορθά από τον εξυπηρετητή. Ο υπολογισμός αυτός είναι όμως προσεγγιστικός: ο κάθε χρήστης εκτιμά σε κάθε γύρο πόση ισχύ θα πρέπει -κατ' ελάχιστον- να καταναλώσει για την αποστολή του μηνύματός του, βασιζόμενος στην ισχύ μετάδοσης των υπολοίπων χρηστών στην προηγούμενη εποχή της Ομόσπονδης Μάθησης.

Την εκτίμηση αυτή μπορεί να αξιοποιήσει ένας κακόβουλος χρήστης για να πραγματοποιήσει επίθεση παρεμβολών. Επειδή οι αρχικοποιήσεις που μελετώνται έχουν επιλεγεί έτσι ώστε σε όλες ο κακόβουλος χρήστης να έχει το μικρότερο κέρδος καναλιού, όταν αυτός στείλει το σήμα του με πολύ μεγαλύτερη ισχύ από όση έστειλε στον προηγούμενο γύρο, τότε αισθάνονται τις παρεμβολές όλοι οι υπόλοιποι συμμετέχοντες. Έτσι, επειδή δεν είναι εφικτό να προβλεφθεί η επίθεση αυτή του κακόβουλου χρήστη, έχει υπολογιστεί λαθεμένα η ελάχιστη ισχύ μετάδοσης από τους ιδιοτελείς χρήστες, με αποτέλεσμα να μεταδίδουν -πιθανώς- με ισχύ μικρότερη από αυτή που είναι αναγκαία. Το αποτέλεσμα είναι ότι τα σήματα των χρηστών φτάνουν στον εξυπηρετητή με τόσες παρεμβολές, που καθίσταται αδύνατη η αποκωδικοποίησή τους, με αποτέλεσμα να αποκλείονται από τον εκάστοτε γύρο της Ομόσπονδης Μάθησης. Ο αριθμός των επιθέσεων που πραγματοποιούνται εξαρτάται άμεσα από την αρχικοποίηση των χρηστών στον χώρο, όπως μαρτυρούν οι εικόνες του σχήματος 6.4.

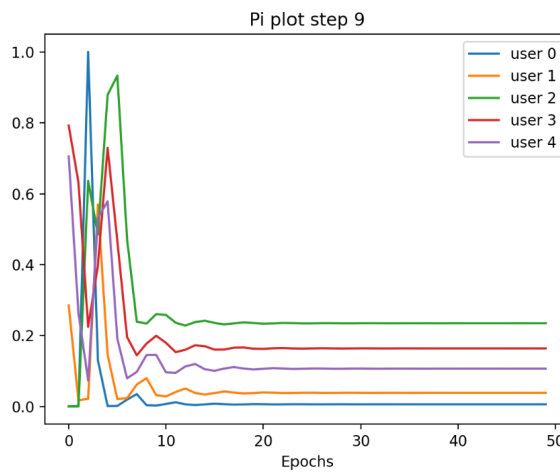
Οι τρεις διαφορετικές αρχικοποιήσεις του σχήματος 6.4 διαφέρουν μεταξύ τους σημαντικά στις πρώ-



(α') Σύγκλιση σε IID δεδομένα, βήμα 0



(β') Σύγκλιση σε IID δεδομένα, βήμα 9



(γ') Σύγκλιση σε IID δεδομένα, βήμα 9

Σχήμα 6.4: Η σύγκλιση της ισχύος μετάδοσης p_i σε σενάρια όπου ο χρήστης 1 είναι κακόβουλος, πραγματοποιεί επιθέσεις παρεμβολών, και το σύνολο δεδομένων των χρηστών έχει IID μορφή

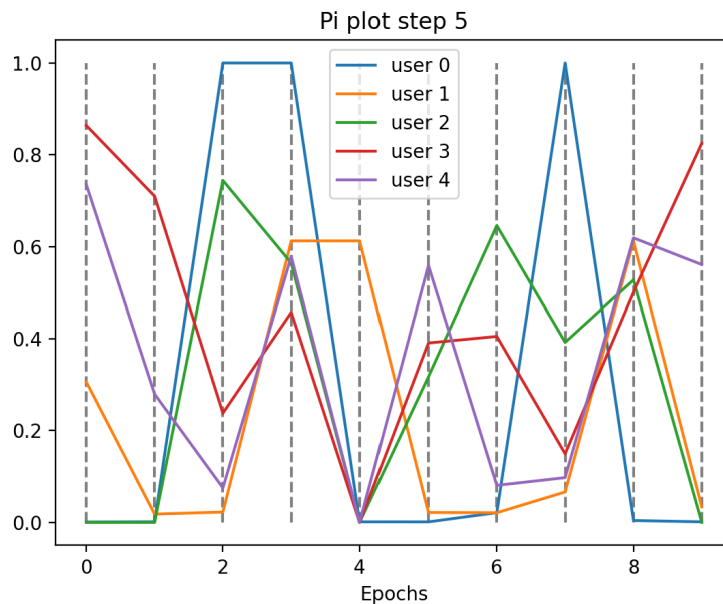
τες εποχές της ομόσπονδης Μάθησης. Μια πρώτη ματιά είναι αρκετή για να διαπιστωθεί πως όσο πιο κοντά πλησιάζουν οι χρήστες στον εξυπηρετητή, όσο μεγαλύτερο δηλαδή είναι το βήμα, τόσο πιο γρήγορα συγκλίνει το παίγνιο. Ξεκινώντας από τις τουλάχιστον 35 εποχές που απαιτούνται για σύγκλιση στην αρχικοποίηση με βήμα 0, ο αριθμός μειώνεται στις 27 όταν το βήμα είναι 5 και καταλήγει περίπου στις 20 για το βήμα 9. Οι διαφορές αυτές στην σύγκλιση εξαρτώνται από το πλήθος επιθέσεων παρεμβολών που πραγματοποιεί ο κακόβουλος χρήστης, οι οποίες με την σειρά τους εξαρτώνται από τα κέρδη καναλιού. Διαισθητικά, θα μπορούσε κανείς να ισχυριστεί, ότι όσο πιο κοντά στον εξυπηρετητή βρίσκονται οι χρήστες (όσο μεγαλύτερο είναι δηλαδή το βήμα) τόσο πιο "καθαρά" φτάνουν τα σήματα σε αυτόν. Έτσι προκειμένου να πραγματοποιηθεί επίθεση παρεμβολών απαιτείται περισσότερη ισχύς μετάδοσης από μέρος του κακόβουλου χρήστη, κάτι που αφενός είναι ασύμφορο και αφετέρου μπορεί να μην είναι εφικτό, λόγω του περιορισμού της ισχύος μετάδοσης σε τιμές κάτω του 1 Watt.

Η ισχύς μετάδοσης του κακόβουλου χρήστη αποτυπώνεται στα διαγράμματα του σχήματος 6.4 με μπλε χρώμα. Κάθε κορύφωση στην συνάρτηση ισχύος φανερώνει μια επίθεση παρεμβολών. Όταν η αρχικοποίηση είναι με βήμα 0, πραγματοποιούνται περίπου 6 επιθέσεις παρεμβολών πριν την σύγκλιση. Ο αριθμός αυτός μειώνεται σε 4 στο βήμα 5 και σε μόλις 2 επιθέσεις στο βήμα 9. Παρόλο που παρουσιάζονται τα αποτελέσματα μόνο των πρώτων 10 βημάτων, από 0 έως 9, παρατηρήθηκε ότι αν οι χρήστες

εξακολουθούσαν να κινούνται προς το κέντρο, από μια απόσταση και μετά ο κακόβουλος χρήστης δεν θα κατάφερνε να πραγματοποιήσει καθόλου επιθέσεις και το παίγνιο θα συνέκλινε πολύ γρήγορα. Τέτοιες περιπτώσεις ωστόσο δεν παρουσιάζουν ερευνητικό ενδιαφέρον, για αυτό και δεν μελετήθηκαν περαιτέρω.

Είναι επίσης σημαντικό να παρατηρηθεί, ότι παρά τις επιθέσεις που πραγματοποιούνται στους πρώτους γύρους της Ομόσπονδης μάθησης, αφού επέλθει η σύγκλιση του παιγνίου, τα αποτελέσματα είναι πανομοιότυπα με αυτά που παρουσιάστηκαν όταν δεν υπήρχε κακόβουλος συμμετέχοντας. Όπως αναφέρθηκε και στην προηγούμενη υποενότητα (6.2.1), ο χρήστης 0 έχει το μικρότερο κέρδος καναλιού, ακολουθούμενος από τους χρήστες 1, 4, 3 και 2 σε αύξουσα σειρά. Αυτή η σειρά είναι που καθορίζει την τελική ισχύ με την οποία καταλήγουν να μεταδίδουν τα σήματά τους οι χρήστες, ανεξάρτητα από το αν πραγματοποιούνται επιθέσεις. Είτε δηλαδή ο χρήστης 0 είναι ιδιοτελής, όπως στην ενότητα 6.2.1, είτε είναι κακόβουλος, είναι εμφανές ότι τον συμφέρει να μεταδίδει με την μικρότερη ισχύ τα σήματά του, επειδή έχει το μικρότερο κέρδος. Ομοίως ο χρήστης 2 που έχει το μεγαλύτερο κέρδος καναλιού υποχρεώνεται σε όλα τα σενάρια να επιλέγει την μεγαλύτερη ισχύ μετάδοσης, ανεξάρτητα από την παρουσία κακόβουλου χρήστη. Η παρατήρηση ότι όσο αυξάνεται το κέρδος καναλιού, τόσο αυξάνεται η τελική ισχύς ενός χρήστη επιβεβαιώνεται σε όλα τα σενάρια, ανεξάρτητα από τις επιθέσεις παρεμβολών.

Η διεξαγωγή του Μπεύζιανού παιγνίου επιτρέπει λοιπόν την επιτυχή αντιμετώπιση των επιθέσεων παρεμβολών, αφού τα τελικά αποτελέσματα είτε υπάρχει κακόβουλος συμμετέχοντας, είτε όχι είναι σχεδόν τα ίδια.



Σχήμα 6.5: Η ισχύς μετάδοσης των χρηστών στο βήμα 5 κατά τις πρώτες 10 εποχές, όπου ο χρήστης 0 είναι κακόβουλος και τα δεδομένα των χρηστών έχουν IID μορφή

Αξίζει να γίνει μια εκτενέστερη αναφορά στο πως η ισχύς μετάδοσης του κακόβουλου χρήστη επηρεάζει την ισχύ των υπολοίπων χρηστών. Το σχήμα 6.5 αποτελεί μια μεγέθυνση της εικόνας 6.4β', προκειμένου να εξηγηθεί η αλληλεπίδραση των χρηστών στους πρώτους γύρους της Ομόσπονδης Μάθησης. Η αρχικοποίηση του παιγνίου θέτει ως αρχική τιμή ισχύος όλων των παιχτών ίση με 0.5 Watt. Αυτή είναι και η τιμή που λαμβάνουν υπόψιν τους οι συμμετέχοντες, προκειμένου να επιλέξουν την ισχύ τους στην εποχή 0.

Δεδομένου ότι όλοι οι χρήστες ξεκινούν με $p_i^{t=-1} = 0.5$ Watt, ο κακόβουλος χρήστης 0 εκτιμά ότι αν οι υπόλοιποι μεταδώσουν με την ίδια ισχύ και στον επόμενο γύρο, δεν θα καταφέρει να κάνει επίθεση

παρεμβολής, εξαιτίας του περιορισμού $p_i \leq P_{max} = 1\text{Watt}$. Ταυτόχρονα επειδή έχει το μικρότερο κέρδος καναλιού, δεν αισθάνεται παρεμβολές από κανέναν άλλον συμμετέχοντα, με αποτέλεσμα η συμφέρουσα λύση για αυτόν είναι να επιλέξει μια πολύ μικρή τιμή ισχύος. Οι χρήστες 1,3 και 4 στέλνουν με σημαντικά μεγαλύτερη ισχύ, προκειμένου να μεταδοθεί το σήμα τους σωστά, δεδομένου ότι όλοι στέλνουν με $p_i^{t=-1} = 0.5\text{ Watt}$. Ο χρήστης 2 ωστόσο, που έχοντας το μεγαλύτερο κέρδος καναλιού αισθάνεται παρεμβολές από όλους τους υπόλοιπους χρήστες, με αποτέλεσμα η (εκτιμώμενη) ελάχιστη απαιτούμενη ενέργεια για ορθή μετάδοση να υπερβαίνει την P_{max} . Έτσι στις εποχές 0 και 1, ο χρήστης 2 δεν μεταδίδει καθόλου το σήμα του, πεπεισμένος ότι όση ισχύ και να χρησιμοποιήσει, δεν θα έχει επιτυχή μετάδοση λόγω παρεμβολών.

Στην εποχή 1, η κατάσταση παραμένει ίδια για τους χρήστες 0 και 2. Οι συμμετέχοντες 1, 3 και 4 ωστόσο, μειώνουν σημαντικά την ισχύ μετάδοσής τους, αφού παρατήρησαν ότι στον γύρο 0 ότι οι παρεμβολές στο σήμα τους ήταν μικρότερες από αυτές που είχαν προβλέψει. Το γεγονός ότι στην εποχή 1 μειώθηκε η ισχύς μετάδοσης των περισσότερων χρηστών, έδωσε την δυνατότητα στον χρήστη 2 να μεταδώσει το σήμα του στον γύρο 2. Ωστόσο δεν είναι μόνο ο χρήστης 2 που αξιοποίησε τις χαμηλές τιμές μετάδοσης ισχύος των συμμετεχόντων 1, 3 και 4. Ο κακόβουλος χρήστης 0 επέλεξε να πραγματοποιήσει επίθεση παρεμβολών στην ίδια εποχή ($p_0 = 1\text{ Watt}$), με αποτέλεσμα να αποκλείει όλους σχεδόν τους χρήστες (με εξαίρεση τον 2) από την διαδικασία της Ομόσπονδης μάθησης.

Στην εποχή 3 αντιλαμβάνονται οι συμμετέχοντες 1, 3 και 4 ότι το σήμα τους απορρίφθηκε από τον διακομιστή λόγω παρεμβολών και αναγκάζονται να αυξήσουν σημαντικά την ισχύ τους, ώστε να εξασφαλίσουν επιτυχή μετάδοση. Ο χρήστης 2 από την άλλη, εκτιμώντας ότι οι υπόλοιποι χρήστες (πλην του 0) θα συνεχίσουν να επιλέγουν την χαμηλή τιμή ισχύος που είχαν στον γύρο 2, "επαναπαύεται" και μειώνει περαιτέρω την ισχύ του. Η νέα μειωμένη ισχύς του χρήστη 2, σε συνδυασμό με τις "αναπάντεχες" αυξήσεις των παρεμβολών που του προκαλούν οι χρήστες 1,3 και 4 έχουν ως αποτέλεσμα να αποκλειστεί ο χρήστης από την τρίτη εποχή της Ομόσπονδης Μάθησης. Αυτή η παρατήρηση παρουσιάζει έντονο ενδιαφέρον, καθώς ο αποκλεισμός του χρήστη 2, δεν οφείλεται σε αιφνιδιαστική επίθεση παρεμβολών, αλλά στην απρόσμενη αύξηση της ισχύος των υπολοίπων χρηστών. Με άλλα λόγια, ο κακόβουλος συμμετέχοντας 0 και η αυξημένη ισχύς μετάδοσής του, δεν είναι ο μοναδικός τρόπος για να αποκλειστεί ένας χρήστης από την Ομόσπονδη μάθηση.

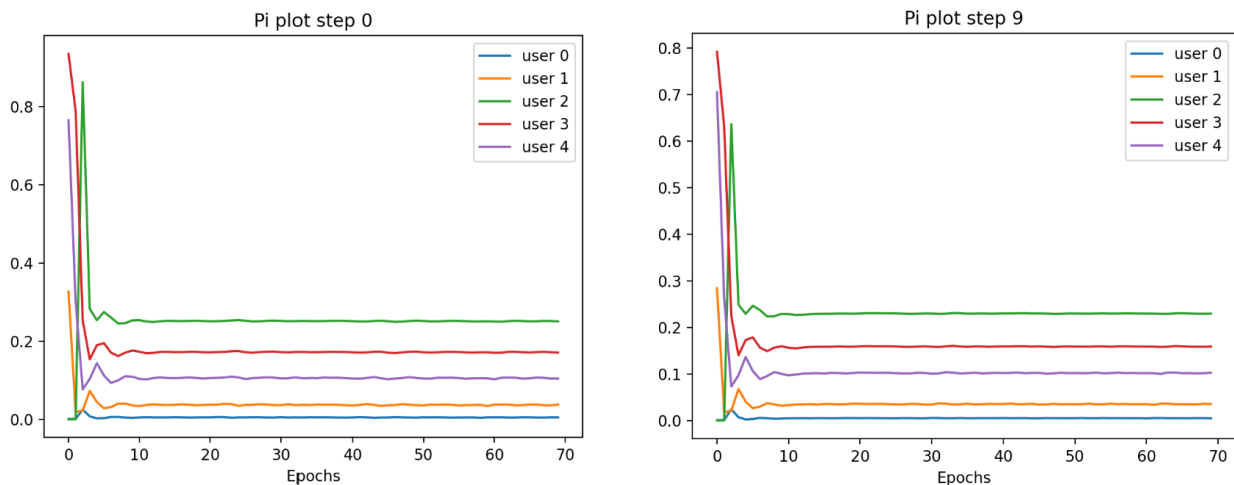
Τα παραπάνω φαινόμενα επαναλαμβάνονται αρκετές φορές μέχρι να επέλθει σύγκλιση. Ο μεν κακόβουλος χρήστης 0 αναμένει να μειωθεί σημαντικά η ισχύς μετάδοσης των υπολοίπων χρηστών, ώστε να προβεί σε επίθεση, ενώ οι δε ιδιοτελείς συμμετέχοντες ανταποκρίνονται στις επιθέσεις, με καθυστέρηση μιας εποχής. Για αυτό άλλωστε και σε όλα τα διαγράμματα, κάθε φορά που πραγματοποιείται επίθεση (κορυφώνεται απότομα η μπλε συνάρτηση του σχήματος 6.5) τότε στον επόμενο γύρο ακολουθούν απότομες αυξήσεις ισχύος από μέρους των υπολοίπων χρηστών.

6.2.3 Αρχικοποιήσεις με κακόβουλο χρήστη σε NON IID δεδομένα

Το σενάριο όπου οι χρήστες έχουν Non IID σύνολο από δεδομένα είναι αρκετά διαφορετικό το αναμενόμενο. Φαινομενικά, το σύνολο δεδομένων των χρηστών δεν επηρεάζει την σύγκλιση του παιγνίου, ωστόσο η πραγματικότητα, με βάση το διάγραμμα 6.6, είναι αρκετά διαφορετική. Σε αντίθεση με το IID σενάριο που μελετήθηκε προηγουμένως, όταν τα δεδομένα των χρηστών έχουν NON IID μορφή, ο κακόβουλος χρήστης φαίνεται να χάνει το κίνητρο του για επίθεση. Επίσης ίσως παρατηρήσει κάποιος ότι μετά την 10η εποχή, όταν έχει επέλθει η σύγκλιση, εμφανίζονται μικρές αυξομειώσεις στην ισχύ μετάδοσης των χρηστών, που αντιστοιχούν στην λογική της σύγκλισης ενός παιγνίου. Και οι δύο αυτές "ανωμαλίες" οφείλονται σε έναν κοινό παρονομαστή.

Η συνάρτηση χρησιμότητας ενός συμμετέχοντα 4.14 είναι αυτή που διαμορφώνει την εξίσωση 4.18 και καθορίζει την ισχύ μετάδοσης που επιλέγεται σε κάθε εποχή, από τον εκάστοτε χρήστη. Η συνάρτηση χρησιμότητας, ωστόσο, πέρα από την μεταβλητή απόφασης -την ισχύ- επηρεάζεται και από μια άλλη -εντελώς ανεξάρτητη- μεταβλητή, τον βαθμό συνεισφοράς των χρηστών a_i . Ο βαθμός συνεισφοράς αλλάζει από εποχή σε εποχή, ανάλογα με το πόσο αξιόλογο τοπικό μοντέλο έστειλε ο κάθε χρήστης στον διακομιστή. Στην περίπτωση όμως που οι χρήστες χρησιμοποιούν NON IID δεδομένα, ο βαθμός συνεισφοράς είναι πολύ δυσκολότερο να υπολογισθεί σωστά, κάτι αποδείχθηκε στα διαγράμματα 6.2 της προηγούμενης ενότητας.

Αφενός, όταν τα δεδομένα είναι NON IID, οι τιμές που λαμβάνει η μετρική a_i έχουν εντονότερες διαβαθμίσεις, με αποτέλεσμα να μεταβάλλουν -ελαφρώς- σε κάθε γύρο την εξίσωση χρησιμότητας και συνεπώς την ισχύ την μεγιστοποιεί. Η αυξομείωση αυτή του βαθμού συνεισφοράς a_i δεν αναιρεί την σύγκλιση του παιγνίου, αφού είναι πλήρως ανεξάρτητη από αυτό. Το ίδιο φαινόμενο παρατηρείται και σε IID μορφή δεδομένων, απλώς σε τέτοια σενάρια, οι μεταβολές του a_i είναι ελαφρώς μικρότερες, με αποτέλεσμα να μην αποτυπώνονται με ακρίβεια στο διάγραμμα ισχύος.



(α') Σύγκλιση σε IID δεδομένα, βήμα 0

(β') Σύγκλιση σε IID δεδομένα, βήμα 9

Σχήμα 6.6: Η σύγκλιση της ισχύος μετάδοσης p_i σε σενάρια όπου ο χρήστης 1 είναι κακόβουλος, πραγματοποιεί επιθέσεις παρεμβολών, και το σύνολο δεδομένων των χρηστών έχει NON IID μορφή

Αφετέρου, όπως παρουσιάστηκε στο κεφάλαιο 6.1, ο κακόβουλος χρήστης λαμβάνει -κατά μέσο όρο- μεγαλύτερη τιμή a_i όταν τα δεδομένα έχουν NON IID μορφή. Αντίστοιχα, επειδή το άθροισμα όλων των a_i ισούται με 1, ο βαθμός συνεισφοράς των υπολοίπων χρηστών μειώνεται όταν ο κακόβουλος χρήστης επιτυγχάνει καλύτερες αποδόσεις. Η σημαντική διαφορά με τα σενάρια που χρησιμοποιούν IID δεδομένα, είναι η πρόβλεψη που πραγματοποιεί ο κακόβουλος χρήστης, αναφορικά με το αν οι υπόλοιποι συμμετέχοντες είναι ιδιοτελείς ή όχι. Ο υπολογισμός αυτός, με την χρήστη σιγμοειδούς συνάρτησης (4.19, 4.20), αποτυπώθηκε στο διάγραμμα του σχήματος 4.2.

Υπενθυμίζεται ότι προκειμένου ο κακόβουλος χρήστης 0 να πραγματοποιήσει πρόβλεψη σχετικά με το αν οι υπόλοιποι συμμετέχοντες είναι κακόβουλοι, συγκρίνει τον βαθμό συνεισφοράς που έλαβε το μοντέλο του, με τον βαθμό των υπολοίπων χρηστών. Αν κάποιο άλλο μοντέλο έχει λάβει τιμή $a_i < a_0$, τότε ο χρήστης 0 μπορεί να αποφανθεί με μεγάλη βεβαιότητα, ότι και ο χρήστης i είναι κακόβουλος. Ωστόσο, όταν ένας χρήστης j έχει λάβει βαθμό $a_j < a_0$, τότε ο κακόβουλος χρήστης δεν μπορεί να εκτιμήσει με βεβαιότητα το κατά πόσο ο j είναι ιδιοτελής χρήστης ή απλώς ένας κακόβουλος χρήστης που έτυχε να μεταδώσει ένα ελαφρώς καλύτερο τοπικό μοντέλο.

Σε σενάρια με IID δεδομένα, ο κακόβουλος χρήστης 0 λάμβανε τιμές $a_0 \approx 0.06$ και οι υπόλοιποι συμμετέχοντες τιμές $a_i > 0.22$. Η διαφορά συνεπώς της συνεισφοράς των κανονικών τοπικών μοντέλων από το δηλητηριασμένο ήταν της τάξης 0.16, οδηγώντας έτσι τον κακόβουλο χρήστη να θεωρεί -με μεγάλη βεβαιότητα- ότι οι υπόλοιποι συμμετέχοντες είναι ιδιοτελείς. Όσον αφορά στα Non IID σενάρια όμως, η διαφορά αυτή μειώνεται από 0.16 σε 0.06. Αυτό έχει ως αποτέλεσμα την αβεβαιότητα του χρήστη 0, περί την ταυτότητα των υπολοίπων συμμετεχόντων.

Το αποτέλεσμα της αβεβαιότητας αυτής αντανακλάται στο πρόβλημα βελτιστοποίησης 4.18 που καλείται να λύσει ο κακόβουλος χρήστης, σύμφωνα με το οποίο επιδιώκει αφενός την αύξηση του κέρδους του και αφετέρου την μείωση της συνάρτησης χρησιμότητας των υπολοίπων χρηστών, με γνώμονα την πιθανότητα αυτοί να είναι ιδιοτελείς. Όταν λοιπόν, λόγω αβεβαιότητας, μειώνεται η πιθανότητα των υπολοίπων χρηστών να είναι ιδιοτελείς, τότε ο κακόβουλος χρήστης χάνει το κίνητρό του για επίθεση. Άλλωστε μια επίθεση παρεμβολών απαιτεί αυξημένη ισχύ, κάτι που μειώνει την συνάρτηση χρησιμότητάς του.

Παρά το μειωμένο κίνητρο και την ελλιπή προσπάθεια του κακόβουλου συμμετέχοντα να προβεί σε επιθέσεις, τα αποτελέσματα μετά την σύγκλιση συνάδουν με αυτά των προηγούμενων σεναρίων. Είτε πρόκειται για IID δεδομένα ή NON IID, είτε οι συμμετέχοντες είναι όλοι ιδιοτελείς είτε υπάρχουν επιθέσεις, μετά την σύγκλιση του παίγνιου, η ισχύς μετάδοσης των χρηστών καταλήγει να εξαρτάται άμεσα από το κέρδος καναλιού τους. Όσο μικρότερο κέρδος καναλιού έχει ένας χρήστης, τόσο μικρότερη είναι η τελική ισχύς με την οποία επιλέγει να μεταδώσει το μήνυμά του.

6.3 Η απόδοση του Μπεϋζιανού Παιγνίου

Σκοπός της παρούσας ενότητας είναι η παρουσίαση διαφόρων μετρικών που σχετίζονται με το Μπεϋζιανό Παιγνίο που λαμβάνει χώρα μεταξύ των χρηστών. Η μελέτη αυτή αφορά διάφορες όψεις της επικοινωνίας των χρηστών με τον εξυπηρετητή, όπως η ισχύς, η ενέργεια και ο χρόνος που απαιτούνται για την μετάδοση των μηνυμάτων. Επίσης θα παρουσιαστούν οι τιμές που λαμβάνει η συνάρτηση χρησιμότητας, που αποδεικνύουν έμπρακτα, ότι ένας κακόβουλος χρήστης θα ωφελούνταν περισσότερο αν δεν έκανε επιθέσεις δηλητηριασμού και παρεμβολών.

Όπως και στις προηγούμενες ενότητες του παρόντος κεφαλαίου, τα αποτελέσματα που παρουσιάζονται θα αποτελούν τον μέσο όρο όλων των διαφορετικών αρχικοποιήσεων ενός σεναρίου. Ωστόσο, προκειμένου να υπάρξει καλύτερη ερμηνεία των αποτελεσμάτων, κρίνεται σκόπιμο να αλλάξει ο συμβολισμός των χρηστών. Όπως είχε αναφερθεί στην ενότητα 4.2, το κέρδος καναλιού είναι η μετρική που επηρεάζει τον ρυθμό μεταφοράς στο τηλεπικοινωνιακό μοντέλο NOMA. Το κέρδος καναλιού ενός συμμετέχοντα υπολογίζεται με βάση την θέση του στον χώρο, αναφορικά με την τοποθεσία του εξυπηρετητή και είναι αυτό που καθορίζει την σειρά με την οποία αποκωδικοποιούνται τα σήματα των χρηστών από τον διακομιστή. Όπως αποδείχθηκε στην προηγούμενη ενότητα, μετά την σύγκλιση του παιγνίου, όσο μικρότερο κέρδος καναλιού έχει ένας χρήστης, τόσο μικρότερη είναι η τελική ισχύς μετάδοσης που επιλέγει.

Πίνακας 6.2: Αντιστοίχιση των ταξινομημένων κερδών καναλιού $G_0 \leq G_1 \leq G_2 \leq G_3 \leq G_4$ στους συμμετέχοντες του παιγνίου

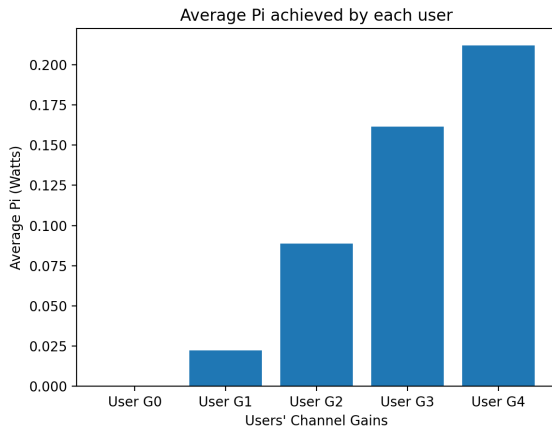
Αντιστοίχιση	Χρήστης 0	Χρήστης 1	Χρήστης 2	Χρήστης 3	Χρήστης 4
Κέρδος καναλιού σε αύξουσα σειρά	G_0	G_1	G_4	G_3	G_2

Ως εκ τούτου, αν ταξινομήσουμε τα κέρδη καναλιού των χρηστών ως $G_0 \leq G_1 \leq G_2 \leq G_3 \leq G_4$, όπου G_i δεν είναι το κέρδος καναλιού του χρήστη i , αλλά το $i + 1$ μικρότερο κέρδος καναλιού μεταξύ των παιχτών, τότε η αντιστοίχιση μεταξύ κερδών και χρηστών είναι αυτή που φαίνεται στον πίνακα 6.2. Όπως εξηγήθηκε και στην ενότητα 5.4, ο χρήστης 0, που σε μερικά σενάρια είναι κακόβουλος, έχει το μικρότερο κέρδος καναλιού, δηλαδή το G_0 . Το αμέσως μικρότερο κέρδος G_1 , ανήκει στον χρήστη 1 και έπειτα ο συμμετέχοντας 4 έχει το κέρδος G_2 . Τέλος, το δεύτερο μεγαλύτερο κέρδος G_3 ανήκει στον χρήστη 3 και το G_4 στον χρήστη 2. Για το υπόλοιπο της ενότητας αυτής, θα αναφέρονται οι χρήστες μέσω του κέρδους καναλιού τους. Για παράδειγμα, ως χρήστης G_4 , θα αναφέρεται ο χρήστης 2, που έχει το μεγαλύτερο κέρδος καναλιού.

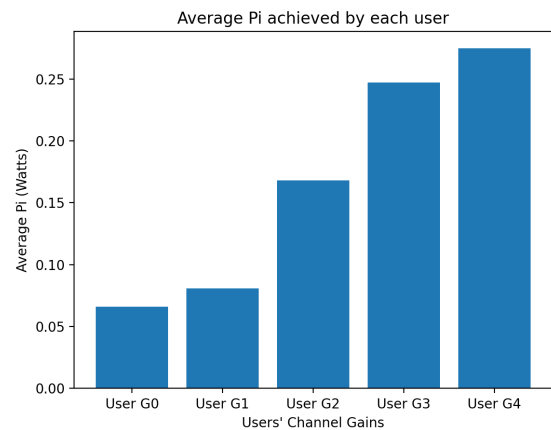
6.3.1 Η μέση ισχύς μετάδοσης

Μια πρώτη μελέτη της ισχύος μετάδοσης των χρηστών πραγματοποιήθηκε στην ενότητα 6.2, όπου παρουσιάστηκαν οι τιμές ισχύος των χρηστών πριν και μετά την σύγκλιση του Μπεϋζιανού παιγνίου. Η παρούσα ενότητα έχει ως σκοπό να παρουσιάσει τα αποτελέσματα όλων των εποχών και αρχικοποιήσεων ταυτόχρονα, ώστε να προκύψει ένα γενικό συμπέρασμα, για το πως η τοποθεσία των χρηστών στον χώρο επηρεάζει την ισχύ μετάδοσής τους. Ως εκ τούτου, στο διάγραμμα 6.7 παρουσιάζεται ο μέσος όρος της ισχύος που επέλεξαν οι συμμετέχοντες στις 10 διαφορετικές αρχικοποιήσεις που μελετήθηκαν. Σε αυτές, αν και αλλάζει η θέση των χρηστών στον χώρο, δεν μεταβάλλεται η σχέση των κερδών καναλιού. Ο χρήστης 0 για παράδειγμα, που έχει το μικρότερο κέρδος καναλιού (G_0) στην πρώτη αρχικοποίηση, εξακολουθεί να έχει το μικρότερο και στην τελευταία αρχικοποίηση. Τα αποτελέσματα παρουσιάζονται

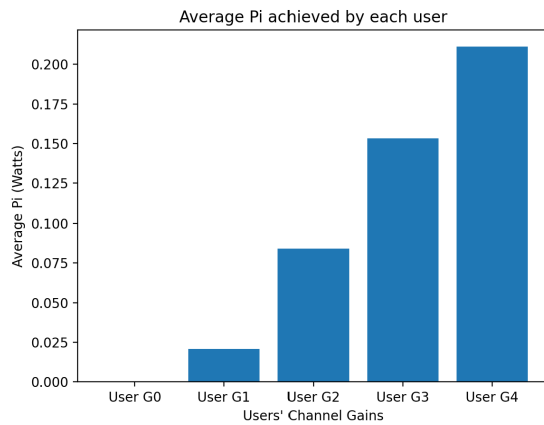
με τους χρήστες να είναι ταξινομημένοι σε αύξουσα σειρά ως προς το κέρδος καναλιού τους, με την αντιστοίχιση που αναφέρεται στον πίνακα 6.2.



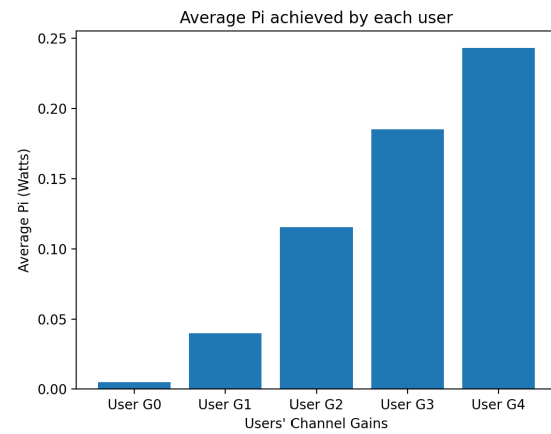
(α') Σενάριο με IID δεδομένα χρηστών, χωρίς κακόβουλο χρήστη



(β') Σενάριο με IID δεδομένα χρηστών, με κακόβουλο χρήστη G_0



(γ') Σενάριο με NON IID δεδομένα χρηστών, χωρίς κακόβουλο χρήστη



(δ') Σενάριο με NON IID δεδομένα χρηστών, με κακόβουλο χρήστη G_0

Σχήμα 6.7: Η μέση ισχύς του κάθε χρήστη, ανάλογα αν τα δεδομένα είναι IID ή NON IID και το κατά πόσο υπάρχει κακόβουλος συμμετέχωντες που πραγματοποιεί επιθέσεις παρεμβολών και δηλητηριασμού

Η διαφορά μεταξύ των αρχικοποιήσεων με και χωρίς κακόβουλο χρήστη είναι εμφανής και αντικατοπτρίζει το αποτέλεσμα των επιθέσεων παρεμβολών. Ο χρήστης G_0 όταν είναι ιδιοτελής (σχήματα 6.7α', 6.7γ') μεταδίδει το σήμα του με σχεδόν μηδενική ισχύ, αφού όντας ο χρήστης με το μικρότερο κέρδος καναλιού, δεν αισθάνεται παρεμβολές από κανέναν άλλον συμμετέχοντα. Όταν όμως αυτός δρα κακοβούλως και πραγματοποιεί επιθέσεις παρεμβολών, η μέση ισχύς του αυξάνεται σημαντικά. Η αύξηση αυτή επηρεάζει την μέση ισχύ όλων των υπολοίπων χρηστών, αφού προκειμένου να αντιμετωπίσουν τις παρεμβολές από τον χρήστη G_0 υποχρεώνονται να αυξήσουν και οι ίδιοι το σήμα τους. Αυτός είναι ο λόγος που στα σχήματα 6.7β', 6.7δ' η μέση ισχύς του κάθε χρήστη είναι σημαντικά αυξημένη συγκριτικά με σχήματα χωρίς επιθέσεις.

Στην υποενότητα 6.2 παρατηρήθηκε ότι στα σενάρια με NON IID data, ο κακόβουλος χρήστης δεν είχε μεγάλο κίνητρο για επιθέσεις. Μάλιστα φάνηκε να πραγματοποιεί με δυσκολία μια-δύο επιθέσεις, με αποτέλεσμα το παίγνιο να συγκλίνει μέσα σε 5 εποχές. Η διαφορά μεταξύ των IID και NON IID αρχικοποιήσεων είναι ευδιάκριτη στα σχήματα 6.7β' και 6.7δ'. Από την μία, η μέση ισχύς του κακόβουλου χρήστη

είναι πάνω από 0.05 Watts, σε IID δεδομένα, ενώ από την άλλη δεν φτάνει ούτε τα 0.01 Watts σε NON IID αρχικοποιήσεις. Ενώ οι επιθέσεις του κακόβουλου χρήστη είναι σημαντικά εντονότερες στο IID σενάριο, δεν μπορεί να αμφισβητήσει κάποιος ότι και για NON IID δεδομένα πραγματοποιούνται επιθέσεις, απλώς λιγότερο συχνά. Για αυτό άλλωστε παρατηρείται και η αύξηση ισχύος από την εικόνα 6.7γ' στην εικόνα 6.7δ'.

Τέλος ανεξάρτητα από το πόσο έντονες ή όχι επιθέσεις παρεμβολών υπάρχουν, παρατηρείται ότι η μέση ισχύς που επιλέγει ο κάθε χρήστης εξαρτάται από το κέρδος καναλιού του. Αυτό εξηγείται σύμφωνα με το τηλεπικοινωνιακό μοντέλο NOMA, που χρησιμοποιεί την τεχνική SIC (ενότητα 2.2) για την διαδοχική ακύρωση παρεμβολών. Με βάση αυτή, όσο μικρότερο κέρδος καναλιού έχει κάποιος, τόσο λιγότεροι συμμετέχοντες του προκαλούν παρεμβολές με τα σήματά τους, με αποτέλεσμα να μην υπάρχει ανάγκη για χρήση μεγάλης ισχύος μετάδοσης. Από την άλλη, ένας συμμετέχοντας με μεγάλο κέρδος καναλιού πρέπει να βεβαιωθεί ότι το σήμα του θα μπορέσει να αποκωδικοποιηθεί ορθά από τον εξυπηρετητή, παρά τις παρεμβολές από τους υπολοίπους, με αποτέλεσμα να καταφεύγει σε μεγαλύτερη ισχύ μετάδοσης.

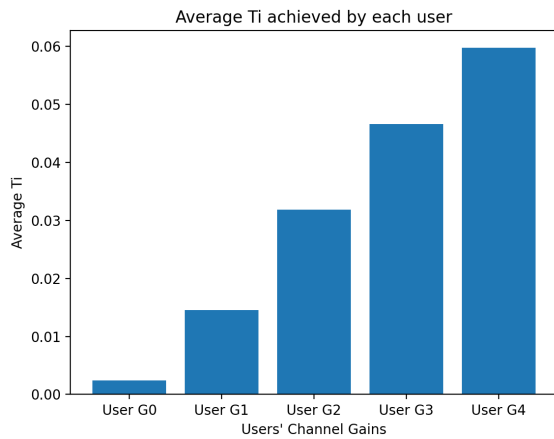
6.3.2 Ο μέσος χρόνος που απαιτείται

Σύμφωνα με την μελέτη που πραγματοποιήθηκε στην ενότητα 4 και την εξίσωση 4.10, ο χρόνος που απαιτείται για την μετάδοση των δεδομένων των χρηστών εξαρτάται από το πλήθος τους και είναι αντιστρόφως ανάλογος από τον ρυθμό μετάδοσης. Ενώ τα δεδομένα των χρηστών είναι σταθερά σε κάθε εποχή της Ομόσπονδης μάθησης και ίσα για όλους τους συμμετέχοντες, ο ρυθμός μετάδοσης (εξίσωση 4.9) είναι μεταβλητός και εξαρτάται από την ισχύ που επιλέγουν και οι υπόλοιποι χρήστες. Μάλιστα, ο ρυθμός μετάδοσης ενός συμμετέχοντα i αυξάνεται όταν αυτός χρησιμοποιεί μεγαλύτερη ισχύ, αλλά μειώνεται όταν αυξάνεται η ισχύς μετάδοσης των χρηστών με μικρότερο κέρδος καναλιού. Ως εκ τούτου, για να μειωθεί ο χρόνος που απαιτείται για την μετάδοση των δεδομένων ενός χρήστη, θα πρέπει να αυξηθεί ο ρυθμός μετάδοσης, δηλαδή είτε ο χρήστης να αυξήσει την ισχύ του, είτε οι υπόλοιποι συμμετέχοντες (που έχουν μικρότερο κέρδος καναλιού από αυτόν) θα πρέπει να την μειώσουν.

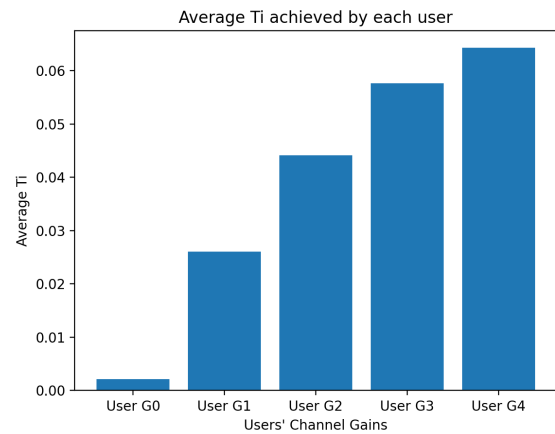
Στο σχήμα 6.8 παρουσιάζεται ο μέσος ρυθμός μετάδοσης, όπως αυτός εκτιμήθηκε από τον εκάστοτε χρήστη σε κάθε εποχή των 10 διαφορετικών αρχικοποιήσεων που εκτελέστηκαν. Με μια πρώτη σύγκριση των εικόνων 6.8α' με 6.8γ' και 6.8β' με 6.8δ' είναι εμφανές ότι οι διαφορές μεταξύ IID και NON IID δεδομένων είναι ελάχιστες. Αυτό είναι αναμενόμενο, καθώς και οι διαφορές μεταξύ των ισχύων μετάδοσης (σχήμα 6.7) ήταν εξίσου δύσκολο να επισημανθούν, ιδίως όταν επρόκειτο για σενάρια χωρίς κακόβουλους χρήστες.

Παρόλα αυτά οι διαφορές μεταξύ των σεναρίων με και χωρίς κακόβουλους χρήστες είναι εντονότερες. Αρχικά ο χρόνος που απαιτείται για την αποστολή του τοπικού μοντέλου του χρήστη με κέρδος G_0 είναι μικρότερος σε σενάρια με επιθέσεις, αν και δεν είναι ιδιαίτερα ορατή η διαφορά στο σχήμα 6.8. Το γεγονός αυτό είναι πλήρως αναμενόμενο, αφού ο χρήστης με G_0 έχει το μικρότερο κέρδος καναλιού και ως κακόβουλος χρήστης, αυξάνει την ισχύ του όταν πραγματοποιεί επιθέσεις. Ως εκ τούτου, ο ρυθμός μετάδοσής του αυξάνεται σημαντικά, με αποτέλεσμα να μειώνεται ο χρόνος που απαιτείται για μετάδοση. Ωστόσο η παραπάνω λογική δεν επαναλαμβάνεται στους υπόλοιπους συμμετέχοντες της ομόσπονδης μάθησης. Ο χρήστης G_0 κατάφερε να αυξήσει τον ρυθμό μετάδοσής του σε σενάρια με επιθέσεις καθώς δεν αισθανόταν παρεμβολές από τα σήματα των υπολοίπων. Οι υπόλοιποι χρήστες, ωστόσο, επηρεάζονται από την αυξημένη ισχύ μετάδοσης του χρήστη G_0 , καθώς και από την ισχύ των υπολοίπων χρηστών με χαμηλότερο κέρδος καναλιού. Η αύξηση που παρατηρείται στον χρόνο μετάδοσης των χρηστών G_1 , G_2 , G_3 και G_4 σε σενάρια με επιθέσεις, έναντι σεναρίων με 5 ιδιοτελείς χρήστες, είναι ενδεικτική της μείωσης του ρυθμού μετάδοσης των χρηστών αυτών. Η μείωση αυτή οφείλεται στο γεγονός ότι παρόλο που όλοι

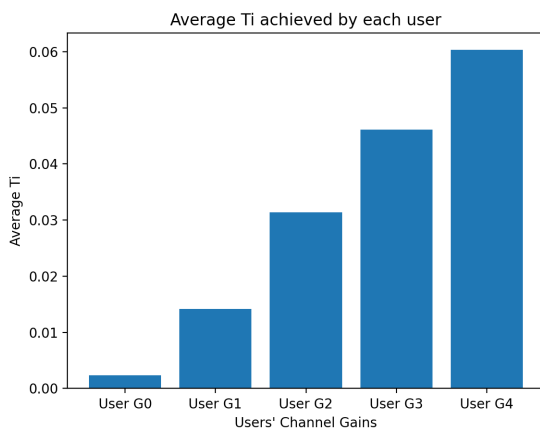
οι χρήστες G_1 , G_2 , G_3 και G_4 αυξάνουν την ισχύ τους σε σενάρια με επιθέσεις (όπως έγινε εμφανές στα σχήματα 6.7) η αύξηση των παρεμβολών που αισθάνονται υπερκεράζει την όποια αύξηση ισχύος.



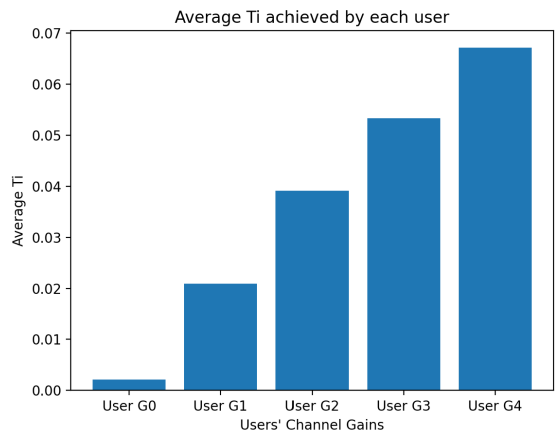
(α') Σενάριο με IID δεδομένα χρηστών, χωρίς κακόβουλο χρήστη



(β') Σενάριο με IID δεδομένα χρηστών, με κακόβουλο χρήστη G_0



(γ') Σενάριο με NON IID δεδομένα χρηστών, χωρίς κακόβουλο χρήστη



(δ') Σενάριο με NON IID δεδομένα χρηστών, με κακόβουλο χρήστη G_0

Σχήμα 6.8: Ο μέσος χρόνος (σε seconds) που απαιτείται για την μετάδοση του κάθε χρήστη, ανάλογα αν τα δεδομένα είναι IID ή NON IID και το κατά πόσο υπάρχει κακόβουλος συμμετέχοντας που πραγματοποιεί επιθέσεις παρεμβολών και δηλητηριασμού

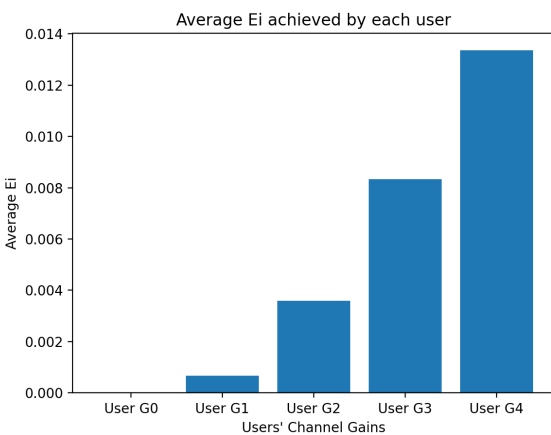
Αξίζει επίσης να επισημανθεί ότι όσο μικρότερο κέρδος καναλιού έχει ένας χρήστης, τόσο λιγότερος χρόνος απαιτείται για την μετάδοση του σήματός του, είτε πρόκειται για σενάρια επιθέσεων είτε για σενάρια με 5 ιδιοτελείς συμμετέχοντας. Πρακτικά, ένας χρήστης με μικρό κέρδος καναλιού νιώθει λιγότερες παρεμβολές συγκριτικά με χρήστες που έχουν μεγαλύτερο κέρδος. Οπότε, αυτός ο χρήστης μπορεί χρησιμοποιήσει μικρότερη ισχύ για την μετάδοσή του και ταυτόχρονα καταφέρνει να διατηρεί τον ρυθμό μετάδοσής του υψηλό, εξαιτίας των χαμηλών παρεμβολών, και συνεπακόλουθα το χρόνο μετάδοσής του χαμηλό.

6.3.3 Η μέση ενέργεια που απαιτείται

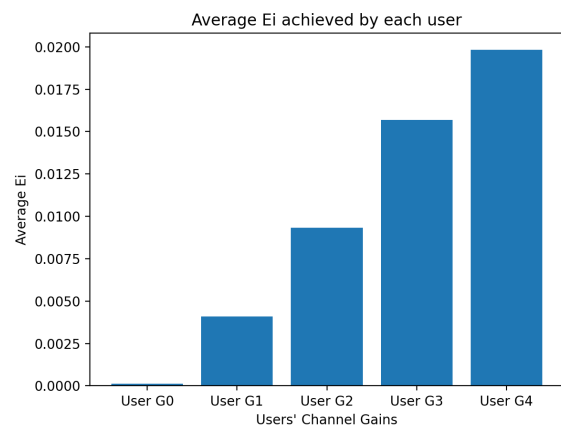
Στην ενότητα 4 παρουσιάστηκε η ενέργεια που απαιτείται από τον εκάστοτε χρήστη για την μετάδοση των δεδομένων του στον εξυπηρετητή, μέσω της εξίσωσης 4.11. Η ενέργεια αυτή αποτελεί το γινόμενο του χρόνου μετάδοσης με την ισχύ που επιλέγει ο συμμετέχοντας σε κάθε γύρο της Ομόσπονδης Μάθησης,

δηλαδή το γινόμενο των αποτελεσμάτων που παρουσιάστηκαν στα σχήματα 6.7 και 6.8.

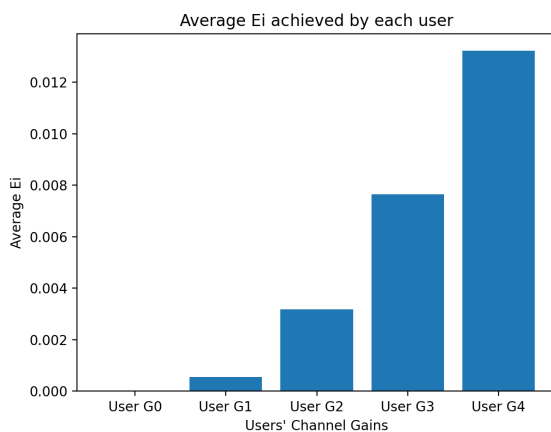
Όπως είναι αναμενόμενο, η μέση ενέργεια που απαιτείται για την μετάδοση είναι άμεσα εξαρτώμενη από το κέρδος καναλιού ενός χρήστη. Άλλωστε, το πόρισμα των προηγούμενων δύο υποενοτήτων κατέστησε σαφές πως όσο μικρότερο κέρδος καναλιού έχει ένας συμμετέχοντας, τόσο υψηλότερο ρυθμό μετάδοσης επιτυγχάνει, συνεπώς τόσο λιγότερος χρόνος απαιτείται για την αποστολή του τοπικού του μοντέλου. Ταυτόχρονα αφού επέλθει η σύγκλιση του Μπεϋζιανού παιχνιδιού οι χρήστες με μικρό κέρδος καναλιού χρησιμοποιούν λιγότερη ισχύ κατά την μετάδοσή τους, ανεξαρτήτως από την ύπαρξη επιθέσεων και την χρήση IID ή NON IID δεδομένων. Δεν είναι -συνεπώς- παράλογο, ότι οι χρήστες με το μεγαλύτερο κέρδος καναλιού, που έχουν τον μεγαλύτερο χρόνο μετάδοσης και χρησιμοποιούν την μεγαλύτερες τιμές ισχύος, θα έχουν και το μεγαλύτερο γινόμενο χρόνου με ισχύ, δηλαδή την μεγαλύτερη απαιτούμενη ενέργεια. Αντίστοιχα οι χρήστες με μικρό κέρδος καναλιού, που δεν αισθάνονται τόσες παρεμβολές από τους υπόλοιπους συμμετέχοντες, επιτυγχάνουν ταχύτερη μετάδοση των δεδομένων τους, με μικρότερη ισχύ, σπαταλώντας έτσι σημαντικά λιγότερη ενέργεια στην κάθε τους μετάδοση.



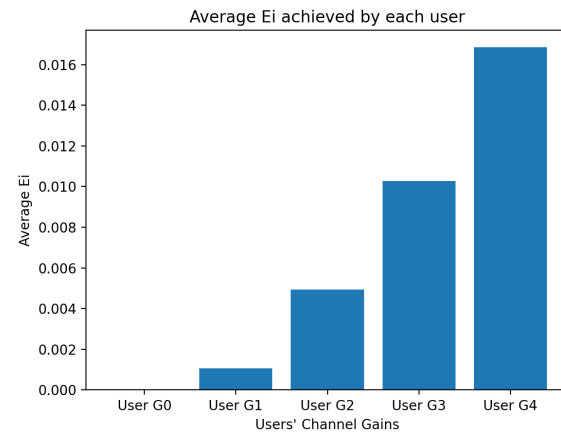
(α) Σενάριο με IID δεδομένα χρηστών, χωρίς κακόβουλο χρήστη



(β) Σενάριο με IID δεδομένα χρηστών, με κακόβουλο χρήστη G_0



(γ) Σενάριο με NON IID δεδομένα χρηστών, χωρίς κακόβουλο χρήστη



(δ) Σενάριο με NON IID δεδομένα χρηστών, με κακόβουλο χρήστη G_0

Σχήμα 6.9: Η μέση ενέργεια που απαιτείται για την μετάδοση του κάθε χρήστη, ανάλογα αν τα δεδομένα είναι IID ή NON IID και το κατά πόσο υπάρχει κακόβουλος συμμετέχοντας που πραγματοποιεί επιθέσεις παρεμβολών και δηλητηριασμού

6.4 Το ποσοστό επιτυχίας του Συγκεντρωτικού Μοντέλου

Μετά την Ολοκλήρωση όλων των εποχών, ο εξυπηρετητής έχει σχηματίσει το τελικό Συγκεντρωτικό Μοντέλο, που είναι άλλωστε και ο στόχος της Ομόσπονδης Μάθησης. Χρησιμοποιώντας ως μέτρο σύγκρισης τα σενάρια όπου δεν πραγματοποιούνται καθόλου επιθέσεις δηλητηριασμού και παρεμβολών, μελετάται ο βαθμός επιτυχίας της προτεινόμενης μεθόδου στην αντιμετώπιση τέτοιων απειλών.

Τα σενάρια που μελετήθηκαν και θα παρουσιαστούν στις επόμενες υποενότητες είναι τα εξής:

1. Χωρίς επίθεση δηλητηριασμού και παρεμβολών, με Contribution Averaging (ContrAvg)

Στο σενάριο αυτό δεν υπάρχει κακόβουλος χρήστης που πραγματοποιεί επιθέσεις και υποβαθμίζει την ποιότητα του τελικού μοντέλου. Η συνάρτηση που χρησιμοποιείται σε κάθε εποχή για την παραγωγή του συγκεντρωτικού μοντέλου είναι η ContrAvg, που προτείνει η παρούσα διπλωματική εργασία. Το σενάριο αυτό απεικονίζεται με μπλε χρώμα στα επόμενα σχήματα.

2. Χωρίς επίθεση δηλητηριασμού και παρεμβολών, με Federated Averaging (FedAvg)

Στο σενάριο αυτό δεν υπάρχει κακόβουλος χρήστης που πραγματοποιεί επιθέσεις και υποβαθμίζει την ποιότητα του τελικού μοντέλου. Η συνάρτηση που χρησιμοποιείται για την παραγωγή του συγκεντρωτικού μοντέλου είναι η FedAvg, η συνάρτηση δηλαδή που προτείνουν οι δημιουργοί της Ομόσπονδης Μάθησης. Το σενάριο αυτό θα απεικονίζεται με πορτοκαλί χρώμα.

3. Με επίθεση δηλητηριασμού και παρεμβολών, με Contribution Averaging (ContrAvg)

Στο σενάριο αυτό υπάρχει ένας κακόβουλος συμμετέχοντας, ο οποίος πραγματοποιεί ταυτόχρονα και επιθέσεις δηλητηριασμού και επιθέσεις παρεμβολών. Η συνάρτηση ContrAvg, που προτείνει η παρούσα εργασία, μειώνει σημαντικά τον αντίκτυπο αυτών των επιθέσεων, ελέγχοντας πόσο αξιόλογο είναι το τοπικό μοντέλο του εκάστοτε συμμετέχοντα. Το σενάριο αυτό απεικονίζεται με πράσινο χρώμα.

4. Με επίθεση δηλητηριασμού και παρεμβολών, με Federated Averaging (FedAvg)

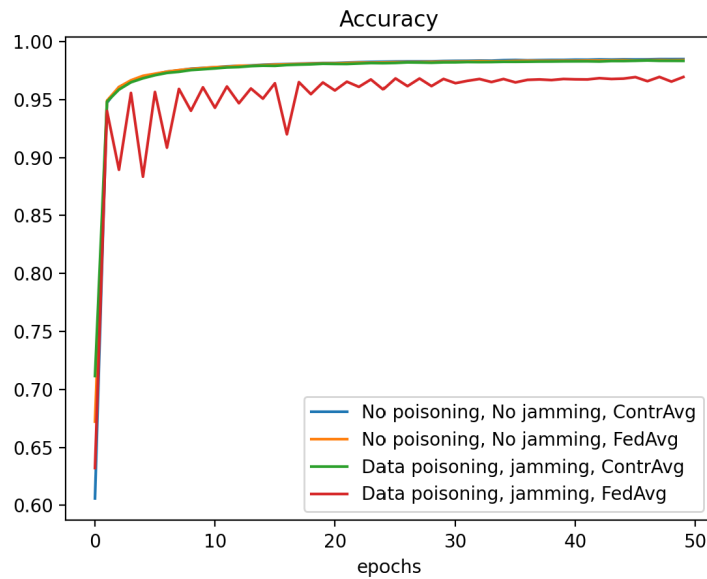
Στο σενάριο αυτό υπάρχει ένας κακόβουλος συμμετέχοντας, ο οποίος πραγματοποιεί ταυτόχρονα και επιθέσεις δηλητηριασμού και επιθέσεις παρεμβολών. Η συνάρτηση που χρησιμοποιείται για την παραγωγή του συγκεντρωτικού μοντέλου είναι η FedAvg, που προτείνουν οι δημιουργοί της Ομόσπονδης Μάθησης. Το σενάριο αυτό απεικονίζεται με κόκκινο χρώμα.

Όπως και στην προηγούμενες ενότητες του κεφαλαίου αυτού, τα αποτελέσματα κάθε σεναρίου προέρχονται από τον μέσο όρο δέκα διαφορετικών αρχικοποιήσεων. Έτσι εξασφαλίζεται η εγκυρότητα τους και εξαλείφεται η παράμετρος τυχαιότητας, που χαρακτηρίζει συχνά τα μοντέλα μηχανικής μάθησης.

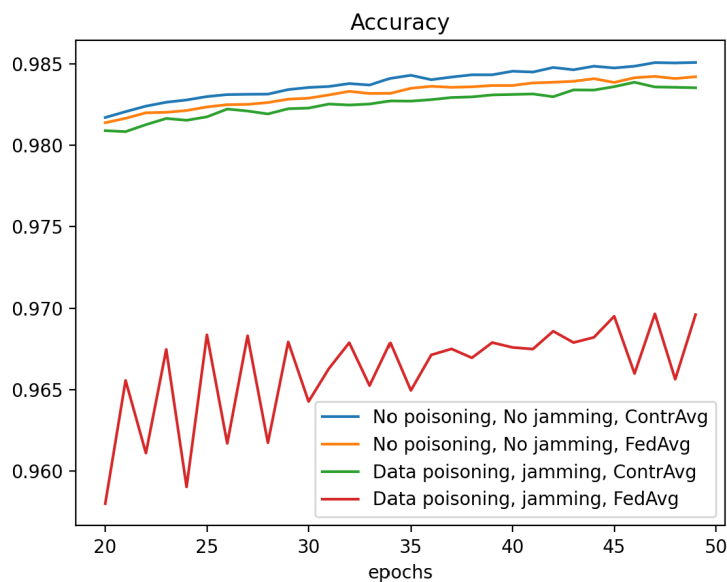
6.4.1 Μελέτη σε IID δεδομένα

Όταν τα δεδομένα των χρηστών έχουν IID Μορφή, επιτυγχάνονται πολύ υψηλά ποσοστά επιτυχίας στο δίκτυο Ομόσπονδης Μάθησης. Τα αποτελέσματα παρουσιάζονται στην εικόνα 6.10α' ή με "μεγένθυση" στην εικόνα 6.10β', όπου είναι ευκρινέστερη η διαφορά στο ποσοστό επιτυχίας μεταξύ των διαφορετικών σεναρίων.

Είναι εμφανές ότι η χρήση της συνάρτησης ContrAvg, που προτείνει η παρούσα εργασία, επιτυγχάνει καλύτερα αποτελέσματα συγκριτικά με την FedAvg τόσο σε σενάρια χωρίς επίθεση, όσο και σε σενάρια με επιθέσεις.



(α) Όλες οι εποχές, 0-50



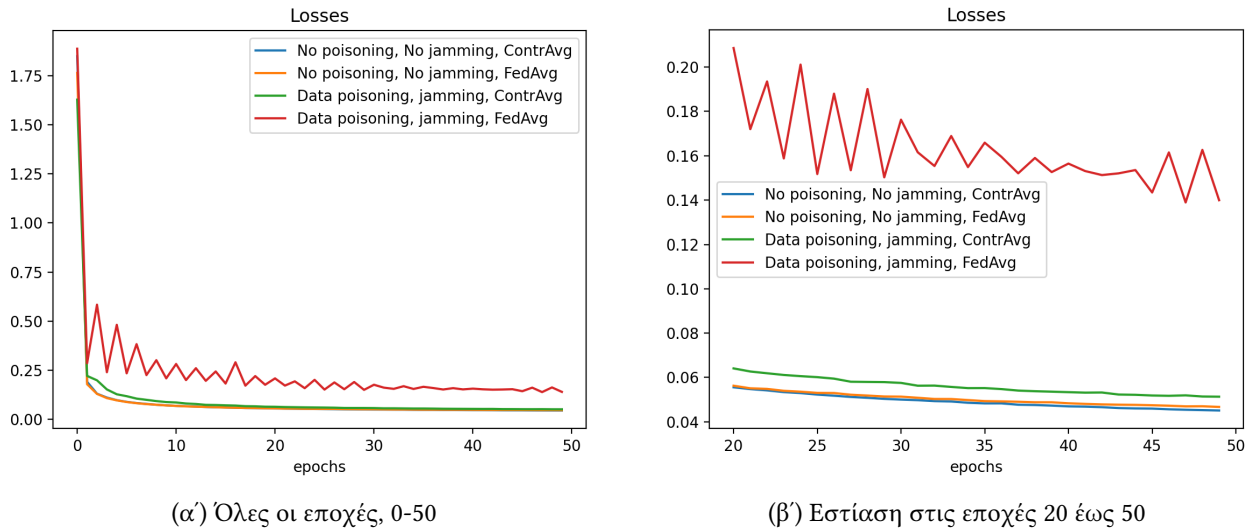
(β) Εστίαση στις εποχές 20 έως 50

Σχήμα 6.10: Ποσοστό επιτυχίας (Accuracy) 4 διαφορετικών σεναρίων, σε IID δεδομένα

Στο τέλος της εποχής 50, το ποσοστό επιτυχίας σε σενάριο χωρίς επιθέσεις έφτασε το 98.51% για τον αλγόριθμο ContrAvg, έναντι 98.42% του FedAvg. Σε περιπτώσεις που μελετήθηκαν επιθέσεις δηλητηριασμού και παρεμβολών, η διαφορά των δύο αλγορίθμων ήταν σημαντικά μεγαλύτερη, με 98.35% για τον ContrAvg και 96.96% για τον FedAvg.

Η μεγαλύτερη επιτυχία του αλγορίθμου ContrAvg που προτείνει η παρούσα εργασία δεν είναι απλώς ότι περνάει τον βασικό αλγόριθμο FedAvg, αλλά το πως επιτυγχάνει να βελτιώσει τα αποτελέσματα σε συνθήκες επιθέσεων δηλητηριασμού. Αν παρατηρήσει κανείς την γραφική παράσταση του ποσοστού επιτυχίας της FedAvg, όταν αυτή δοκιμάζεται σε σενάρια επιθέσεων, θα παρατηρήσει πολύ έντονες αυξομειώσεις από εποχή σε εποχή. Οι απότομες αυτές αλλαγές οφείλονται στο ότι υπάρχουν 5 χρήστες στο σύστημα, ένας από τους οποίους είναι κακόβουλος και δηλητηριάζει τα δεδομένα του. Δεδομένου ότι όλοι οι χρήστες έχουν ίδιο πλήθος δεδομένων, σύμφωνα με την FedAvg το αθροιστικό μοντέλο προκύ-

ππει από τον μέσο όρο των παραμέτρων όλων των τοπικών μοντέλων. Συνεπώς το αθροιστικό μοντέλο επηρεάζεται κατά $\frac{1}{5}$ ή 20% από το τοπικό μοντέλο του κακόβουλου χρήστη, το οποίο είναι δηλητηριασμένο. Το ποσοστό αυτό είναι αρκετά μεγάλο, για αυτό και δίνει την δυνατότητα στον κακόβουλο χρήστη να επιτυγχάνει μεγάλη βλάβη στο συγκεντρωτικό μοντέλο, ειδικά μόλις τείνει να βελτιωθεί το ποσοστό επιτυχίας (accuracy) του.



(α) Όλες οι εποχές, 0-50

(β) Εστίαση στις εποχές 20 έως 50

Σχήμα 6.11: Οι απώλειες (Loss) 4 διαφορετικών σεναρίων, σε IID δεδομένα

Οι απώλειες (Loss) του συγκεντρωτικού μοντέλου της Ομόσπονδης μάθησης παρουσιάζουν παρόμοια συμπεριφορά με το ποσοστό επιτυχίας (Accuracy). Είναι εμφανές ότι ο αλγόριθμος ContrAvg που προτείνει η παρούσα εργασία επιτυγχάνει πολύ μικρότερες απώλειες συγκριτικά με τον βασικό αλγόριθμο FedAvg που χρησιμοποιείται στην Ομόσπονδη Μάθηση, τόσο σε σενάρια με επιθέσεις όσο και σε σενάρια χωρίς.

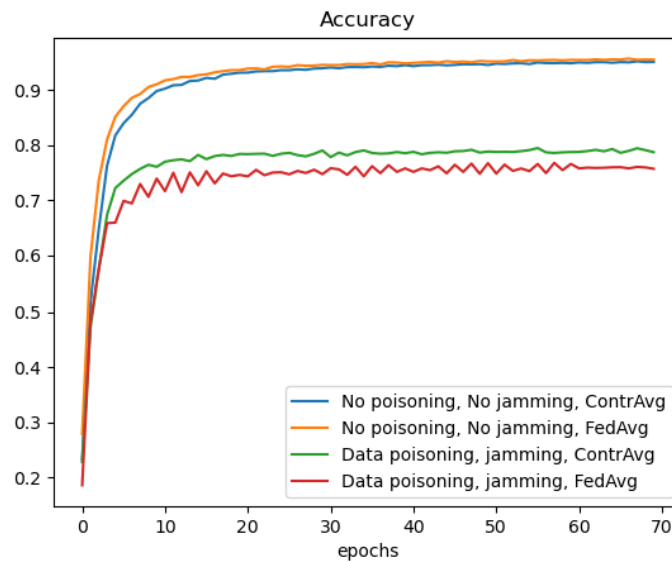
Σενάριο	Αλγόριθμος	Accuracy	Loss
Χωρίς επιθεση	FedAvg	98.42%	0.047
	ContrAvg (Ours)	98.51%	0.045
Με επιθέσεις	FedAvg	96.96%	0.14
	ContrAvg (Ours)	98.35%	0.051

Πίνακας 6.3: Τα συγκεντρωτικά αποτελέσματα τεσσάρων διαφορετικών σεναρίων σε IID δεδομένα, στο τέλος της εποχής 50

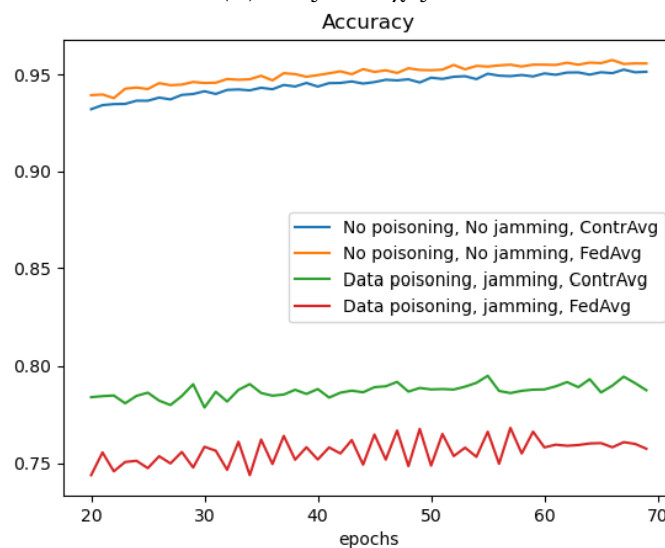
Πιο συγκεκριμένα, οι απώλειες όταν χρησιμοποιείται η συνάρτηση ContrAvg είναι 0.045 όταν υπάρχουν 5 ιδιοτελείς χρήστες και 0.051 όταν υπάρχει 1 κακόβουλος συμμετέχοντας που πραγματοποιεί επιθέσεις δηλητηριασμού και παρεμβολών. Τα αποτελέσματα αυτά είναι σημαντικά βελτιωμένα, αν συγκριθούν με αυτά του αλγορίθμου FedAvg, που αποτελεί το επίπεδο αναφοράς στην Ομόσπονδη Μάθηση, καθώς αυτός επιτυγχάνει 0.047 και 0.14 αντίστοιχα. Όπως και στο ποσοστό επιτυχίας, η διαφορά των δύο αλγορίθμων ως προς τις απώλειες είναι εντονότερη όταν αυτοί δοκιμάζονται σε περιβάλλον με κακόβουλο χρήστη. Σε τέτοιες περιπτώσεις, τα αποτελέσματα του αλγορίθμου FedAvg είναι αρκετά θορυβώδες, με έντονες αυξομειώσεις, ένα φαινόμενο που έχει εξαλείψει πλήρως ο ContrAvg που προτείνει η παρούσα εργασία.

6.4.2 Μελέτη σε NON IID δεδομένα

Η μελέτη των σεναρίων στα οποία οι χρήστες έχουν το σύνολο δεδομένων τους σε NON IID μορφή είναι αρκετά διαφορετική. Υπενθυμίζεται ότι σε τέτοιες αρχικοποιήσεις ο κάθε χρήστης έχει εικόνες που αντιστοιχούν μονάχα σε 2 αριθμούς και εκπαιδεύει το τοπικό μοντέλο του ώστε να αναγνωρίζει αποκλειστικά αυτούς. Το γεγονός αυτό καθιστά τον εντοπισμό του κακόβουλου χρήστη σημαντικά δυσκολότερο, όπως έγινε εμφανές στην ενότητα 6.1. Πιο συγκεκριμένα, όταν τα δεδομένα έχουν NON IID μορφή, ο κακόβουλος χρήστης θ επιτυγχάνει $A_\theta \approx 0.15$ (έναντι 0.06 σε IID), ενώ οι υπόλοιποι λαμβάνουν τιμές a_i ελαφρώς μεγαλύτερες του 0.2. Επειδή η μετρική a_i είναι αυτή που διαμορφώνει την συνάρτηση ContrAvg (4.8), το συγκεντρωτικό μοντέλο καταλήγει να επηρεάζεται σημαντικά περισσότερο από τις επιθέσεις δηλητηριασμού του κακόβουλου χρήστη, όταν τα δεδομένα είναι σε NON IID μορφή, αντί για IID.



(α') Όλες οι εποχές, 0-70



(β') Εστίαση στις εποχές 20 έως 70

Σχήμα 6.12: Ποσοστό επιτυχίας (Accuracy) 4 διαφορετικών σεναρίων, σε NON IID δεδομένα

Στο σχήμα 6.12 παρουσιάζονται τα ποσοστά επιτυχίας των αλγορίθμων FedAvg και ContrAvg, τόσο σε σεναρία με επιθέσεις όσο και σε χωρίς. Παρατηρείται ότι όταν οι χρήστες έχουν NON IID δεδομένα και

δεν πραγματοποιούνται καθόλου επιθέσεις, ο αλγόριθμος FedAvg έχει ελαφρώς καλύτερα αποτελέσματα από τον ContrAvg που προτείνει η παρούσα διπλωματική, αν και η διαφορά τους είναι μόλις 0.4%. Στις αρ-
χικοποιήσεις όμως που πραγματοποιούνται επιθέσεις, ο αλγόριθμος ContrAvg έχει σημαντικά καλύτερα αποτελέσματα και λιγότερη αστάθεια, συγκριτικά με τον FedAvg.

Για να γίνει κατανοητή η διαφορά στα αποτελέσματα των δύο αλγορίθμων, όταν υπόκεινται σε επιθέσεις από έναν κακόβουλο συμμετέχοντα, πρέπει πρώτα να πραγματοποιηθεί μια αναφορά στο είδος επιθέσεων δηλητηριασμού που μελετήθηκαν. Οι επιθέσεις αυτές ανήκουν στην κατηγορία δηλητηριασμού ετικετών (label poisoning), στην οποία ο κακόβουλος χρήστης αλλάζει τις επισημάνσεις των δεδομένων του, έτσι ώστε να είναι εσφαλμένες (ενότητα 2.1.2). Στα πειράματα της παρούσας εργασίας, ο κακόβουλος χρήστης αύξανε τον αριθμό της ετικέτας κατά 1. Ως αποτέλεσμα, οι εικόνες με τον αριθμό 0 να έχουν ως επισήμανση τον αριθμό 1, οι εικόνες με τον αριθμό 1 να αντιστοιχίζονται στον αριθμό 2 και ούτω καθεξής. Όταν όμως τα δεδομένα των χρηστών είναι σε NON IID μορφή, τότε ο κακόβουλος χρήστης έχει προσεγγιστικά εικόνες μόνο από 2 αριθμούς, τους οποίους δεν έχει κανένas από τους υπόλοιπους συμμετέχοντες. Εξαιτίας της επίθεσης δηλητηριασμού όμως, οι ετικέτες των δυο αριθμών αυτών αλλοιώνονται πλήρως. Ως εκ τούτου, το συγκεντρωτικό μοντέλο δεν μπορεί να εκπαιδευτεί πάνω σε αυτούς τους δύο αριθμούς, αφού κανένas χρήστης δεν εκπαιδεύει το τοπικό μοντέλο του ώστε να τους αναγνωρίζει.

Αν θεωρηθεί ότι το σύνολο αξιολόγησης (test set), στο οποίο εξετάζεται η επιτυχία του συγκεντρωτικού μοντέλου, περιλαμβάνονται όλοι οι αριθμοί από το 0 έως το 9 με ομοιόμορφη κατανομή, τότε ο κάθε αριθμός συναντάται -προσεγγιστικά- στο $\frac{1}{10}$ των εικόνων του συνόλου αξιολόγησης. Εξαιτίας όμως της επίθεσης δηλητηριασμού, δύο (2) από τις δέκα (10) διαθέσιμες ετικέτες αγνοούνται πλήρως, με αποτέλεσμα το συγκεντρωτικό μοντέλο να μην έχει εκπαιδευτεί στην αναγνώριση του 20% των ετικετών του συνόλου αξιολόγησης. Ως αποτέλεσμα το άνω φράγμα του ποσοστού επιτυχίας οποιουδήποτε αλγορίθμου να τίθεται περίπου στο 80%, αφού το υπόλοιπο 20% αποτελείται από εικόνες στις οποίες δεν εκπαιδεύτηκε ποτέ το συγκεντρωτικό μοντέλο.

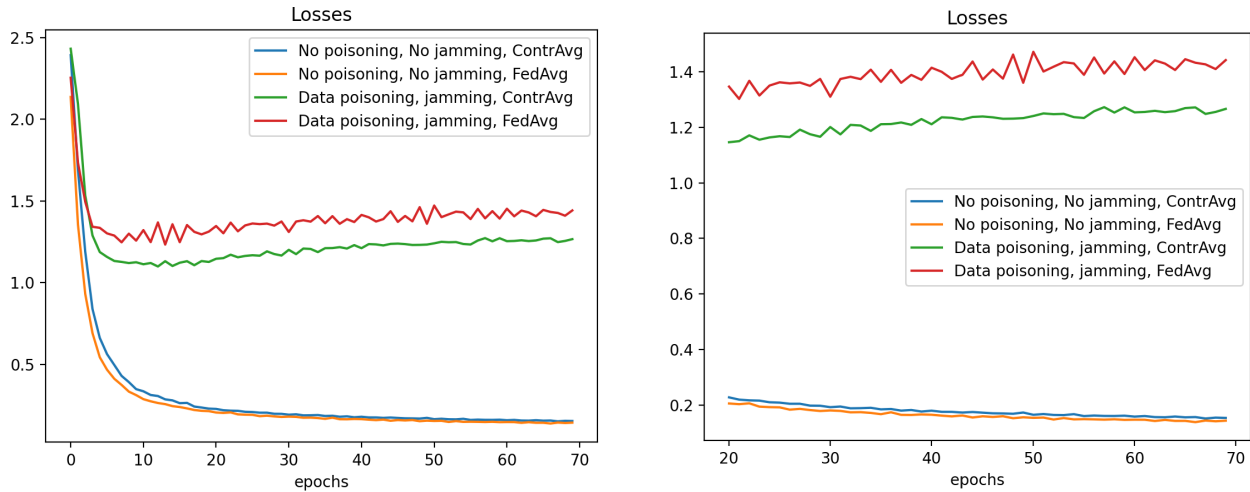
Η παραπάνω παρατήρηση καθιστά εμφανή την διαφορά στα αποτελέσματα των αλγορίθμων FedAvg και ContrAvg. Όχι μόνο επιτυγχάνει κατά 3% καλύτερα αποτελέσματα ο αλγόριθμος που προτείνει η παρούσα εργασία, αλλά τα αποτελέσματα αυτά είναι μόλις 1.24% μακριά από το ανώτατο όριο (80%) που οποιοσδήποτε αλγόριθμος θα μπορούσε να επιτύχει. Ταυτόχρονα, ο αλγόριθμος FedAvg αφενός παρουσιάζει μεγαλύτερη αστάθεια από εποχή σε εποχή και αφετέρου επιτυγχάνει μόλις 75.74% ποσοστό επιτυχίας, έναντι 78.74% του ContrAvg.

Σενάριο	Αλγόριθμος	Accuracy	Loss
Χωρίς επίθεση	FedAvg	95.56%	0.143
	ContrAvg (Ours)	95.13%	0.153
Με επιθέσεις	FedAvg	75.74%	1.442
	ContrAvg	78.76%	1.266

Πίνακας 6.4: Τα συγκεντρωτικά αποτελέσματα τεσσάρων διαφορετικών σεναρίων σε NON IID δεδομένα, στο τέλος της εποχής 50

Οι απώλειες του αθροιστικού μοντέλου ακολουθούν παρόμοια τάση με το ποσοστό επιτυχίας. Όταν οι χρήστες έχουν NON IID δεδομένα και δεν πραγματοποιούνται επιθέσεις δηλητηριασμού ή παρεμβολών, ο αλγόριθμος FedAvg έχει ελαφρώς καλύτερα αποτελέσματα απωλειών (0.143) από τον ContrAvg (0.153) που προτείνει η παρούσα διπλωματική. Η διαφορά αυτή βέβαια είναι μάλλον αμελητέα, μπροστά στην βελτίωση που επιφέρει ο αλγόριθμος ContrAvg σε σεναρία επιθέσεων. Όπως διαπιστώθηκε και κατά την

μελέτη του ποσοστού ακρίβειας, η γραφική συνάρτηση των απωλειών του αλγορίθμου FedAvg παρουσιάζει έντονες αυξομειώσεις, εξαιτίας των επιθέσεων δηλητηριασμού. Το φαινόμενο αυτό έχει καταφέρει να περιορίσει σημαντικά ο αλγόριθμος ContrAvg, που όχι μόνο εμφανίζει λιγότερες απώλειες, αλλά και λιγότερη αστάθεια όταν υπάρχει κακόβουλος χρήστης στο σύστημα Ομόσπονδης Μάθησης.



(α) Όλες οι εποχές, 0-70

(β) Εστίαση στις εποχές 20 έως 70

Σχήμα 6.13: Οι απώλειες (Loss) 4 διαφορετικών σεναρίων, σε IID δεδομένα

Με μια πιο προσεκτική ματιά θα μπορούσε κάποιος να παρατηρήσει ότι σε σενάρια με επιθέσεις και NON IID δεδομένα, οι απώλειες φαίνεται να αυξάνονται ελαφρώς όσο περνούν οι εποχές. Το φαινόμενο αυτό παρουσιάζει έντονο ενδιαφέρον και για να ερμηνευτεί είναι αναγκαία η κατανόηση της μετρικής των απωλειών.

Στο τέλος κάθε εποχής της ομόσπονδης μάθησης, το συγκεντρωτικό μοντέλο που προκύπτει από τα τοπικά μοντέλα, αξιολογείται ως προς το ποσοστό επιτυχίας του και τις απώλειες του. Ενώ το μεν ποσοστό επιτυχίας ασχολείται αποκλειστικά με το σε πόσες εικόνες αναγνωρίζει το συγκεντρωτικό μοντέλο σωστά τον αριθμό, η μετρική των απωλειών λειτουργεί ελαφρώς διαφορετικά. Μια απλοϊκή εξήγηση των απωλειών είναι ότι κάθε φορά που το μοντέλο αποφασίζει τι νούμερο αναπαριστάται στην κάθε εικόνα, δίνει και ένα ποσοστό για την σιγουριά του στην εκτίμηση αυτή. Οι απώλειες εκφράζουν όχι μόνο αν το μοντέλο αναγνώρισε σωστά τον αριθμό στην εικόνα, αλλά και το κατά πόσο ήταν σίγουρο για την εκτίμηση αυτή. Ως εκ τούτου, οι απώλειες αυξάνονται, αν το νευρωνικό δίκτυο δώσει εσφαλμένη απάντηση με μεγάλη βεβαιότητα.

Η αύξηση που παρατηρείται στις απώλειες του σχήματος 6.13 μπορεί να ερμηνευτεί εύκολα, αν αναλογιστεί κανείς τον αντίκτυπο που έχει ο κακόβουλος χρήστης στις εκτιμήσεις του συγκεντρωτικού μοντέλου. Οι επιθέσεις δηλητηριασμού που πραγματοποιούνται στις συγκεκριμένες αρχικοποιήσεις αλλάζουν την ετικέτα των εικόνων με τον αριθμό 0 στον αριθμό 1 και των εικόνων με το 1 στον αριθμό 2. Όσο περνούν οι εποχές, το συγκεντρωτικό μοντέλο μαθαίνει με όλο και μεγαλύτερη βεβαιότητα, ότι οι εικόνες που δείχνουν τον αριθμό 0 είναι ο αριθμός 1, κάτι που φυσικά οδηγεί σε λανθασμένες προβλέψεις. Ο συνδυασμός μίας -επανελημμένως- λανθασμένης πρόβλεψης με μια διαρκώς μεγαλύτερη σιγουριά, οδηγούν σε αύξηση των απωλειών.

Ένα τέτοιο φαινόμενο μπορεί να αποτραπεί με πολλούς τρόπους. Αρχικά αν ο κακόβουλος χρήστης άλλαξε την επίθεση δηλητηριασμού και αντί να αυξάνει τον αριθμό της κάθε ετικέτας κατά 1, τον άλλαζε τυχαία, τότε το συγκεντρωτικό μοντέλο δεν θα μάθαινε με μεγάλη βεβαιότητα να κάνει λανθασμένες προ-

βλέψεις. Ταυτόχρονα θα μπορούσαν να εφαρμοστούν επιπλέον κριτήρια στην Ομόσπονδη Μάθηση, που θα απέρριπταν το μοντέλο ενός χρήστη, αν για παράδειγμα δεν επιτύγχανε κάποιο συγκεκριμένο βαθμό συνεισφοράς a_i . Έτσι, για παράδειγμα, αν ένα τοπικό μοντέλο είχε $a_i \leq 0.17$, τότε δεν θα θεωρούνταν επιτυχές και θα απορρίπτονταν ως επίθεση δηλητηριασμού. Τέτοιες ενέργειες και επιλογές πιθανώς να αύξαναν σημαντικά το ποσοστό επιτυχίας και να περιορίζαν τις απώλειες, ωστόσο δεν αποτελούν τον σκοπό της παρούσας διπλωματικής εργασίας.

Επίλογος και Μελλοντικές Επεκτάσεις

Ο βασικός σκοπός της παρούσας πτυχιακής εργασίας είναι να βελτιώσει την ασφάλεια ασύρματων δικτύων Ομόσπονδης Μάθησης, θωρακίζοντάς τα απέναντι σε επιθέσεις παρεμβολών και δηλητηριασμού. Σε κάθε εποχή της Ομόσπονδης Μάθησης, οι χρήστες επιλέγουν, μέσω του Μπεϋζιανού Παιγνίου, την τιμή ισχύος εκπομπής που εξυπηρετεί βέλτιστα τις ανάγκες τους, ώστε να μεταδώσουν στον εξυπηρετητή το τοπικό τους μοντέλο. Όταν ο διακομιστής λάβει τα τοπικά μοντέλα όλων των συμμετεχόντων, προβαίνει στον υπολογισμό του βαθμού συνεισφοράς τους, μέσα από μία τροποποιημένη εκδοχή της τιμής Shapley. Ο βαθμός συνεισφοράς αυτός χρησιμοποιείται στην συνέχεια από τον αλγόριθμο συνάθροισης (ContrAvg) που προτείνουμε, έτσι ώστε τα τοπικά μοντέλα με τα καλύτερα αποτελέσματα να συνεισφέρουν περισσότερο στην δημιουργία του συγκεντρωτικού μοντέλου. Ο εκάστοτε γύρος της Ομόσπονδης Μάθησης ολοκληρώνεται με την κοινοποίηση του συγκεντρωτικού μοντέλου στους χρήστες, ώστε να το εκπαιδεύσουν εκ νέου στην επόμενη εποχή. Η συνεισφορά των προτάσεων της παρούσας εργασίας επιβεβαιώθηκε με την διεξαγωγή εκτενών πειραμάτων:

1. Η πειραματική μελέτη τουλάχιστον 80 διαφορετικών αρχικοποιήσεων της Ομόσπονδης Μάθησης απέδειξε ότι ο αλγόριθμος ContrAvg επιτυγχάνει πάντα καλύτερα αποτελέσματα, σε σενάρια επιθέσεων, από τον αντίστοιχο αλγόριθμο FedAvg της διεθνούς βιβλιογραφίας, αφού αυξάνει το ποσοστό επιτυχίας, μειώνει τις απώλειες και βοηθάει στην σύγκλιση των αποτελεσμάτων, χωρίς έντονες διακυμάνσεις μεταξύ των εποχών.
2. Ο βαθμός συνεισφοράς των χρηστών, που υπολογίζεται με πολυπλοκότητα $O(|N| + |M|)$, έναντι $O(2^{|N|+|M|})$ που απαιτεί η τιμή Shapley, εντοπίζει επιτυχώς (σε όλα τα σενάρια) τις επιθέσεις δηλητηριασμού, αναθέτοντας χαμηλότερες τιμές στα τοπικά μοντέλα που προέρχονται από δηλητηριασμένα δεδομένα.
3. Η διεξαγωγή ενός Μπεϋζιανού παιγνίου, επιτρέπει στους συμμετέχοντες να επιλέγουν την βέλτιστη ισχύ μετάδοσης του σήματός τους, αντιμετωπίζοντας έτσι επιθέσεις παρεμβολών από κακόβουλους χρήστες. Τα πειραματικά δεδομένα απέδειξαν την ταχεία σύγκλιση του παιγνίου, την εξάλειψη των επιθέσεων παρεμβολών και οδήγησαν τους χρήστες σε χρήση ισχύος μετάδοσης παρόμοια με αυτή που επιλέγουν όταν απουσιάζει ο κακόβουλος χρήστης. Αποδείχθηκε επίσης ότι ο κακόβουλος χρήστης ωφελείται περισσότερο, σε βάθος χρόνου, αν δεν πραγματοποιεί επιθέσεις παρεμβολών και δηλητηριασμού, αφού το σύστημα που προτείνουμε εντοπίζει και τα δύο είδη επιθέσεων, με αποτέλεσμα να μειώνεται το κέρδος του.

Τα πειράματα της παρούσας εργασίας διεξήχθησαν υπό περιορισμένους υπολογιστικούς πόρους, που δεν επέτρεψαν την εκπαίδευση δικτύων Ομόσπονδης Μάθησης με την τιμή Sharpley, που απαιτεί $O(2^{|N|+|M|})$ αξιολογήσεις σε κάθε εποχή. Συνεπώς, η σύγκριση της τιμής Sharpley με τον προσεγγιστικό βαθμό συνεισφοράς που χρησιμοποιήθηκε, τόσο ως προς τον χρόνο όσο και προς την επιτυχία τους, αφήνεται ως μελλοντική επέκταση. Επίσης ο αλγόριθμος ContrAvg δοκιμάστηκε μόνο υπό την απειλή επιθέσεων δηλητηριασμού αλλαγής ετικέτας, που αποτελεί υποκατηγορία των επιθέσεων μη καθαρής ετικέτας. Είναι συνεπώς εύλογη η αξιολόγηση του αλγορίθμου αυτού και του βαθμού συνεισφοράς, στην αντιμετώπιση διαφορετικών επιθέσεων δηλητηριασμού δεδομένων (καθαρής και μη ετικέτας) ή στον δηλητηριασμό μοντέλου.

Βιβλιογραφία

- [1] B. McMahan, E. Moore, D. Ramage, S. Hampson και B. A. γ Arcas, «Communication-efficient learning of deep networks from decentralized data», στο *Artificial intelligence and statistics*, PMLR, 2017, σσ. 1273–1282.
- [2] A. Hard, K. Rao, R. Mathews, S. Ramaswamy, F. Beaufays, S. Augenstein, H. Eichner, C. Kiddon και D. Ramage, «Federated learning for mobile keyboard prediction», *arXiv preprint arXiv:1811.03604*, 2018.
- [3] Q. Yang, Y. Liu, T. Chen και Y. Tong, «Federated machine learning: Concept and applications», *ACM Transactions on Intelligent Systems and Technology (TIST)*, τόμ. 10, αρθμ. 2, σσ. 1–19, 2019.
- [4] Y. Liu, Y. Kang, T. Zou, Y. Pu, Y. He, X. Ye, Y. Ouyang, Y.-Q. Zhang και Q. Yang, «Vertical Federated Learning: Concepts, Advances, and Challenges», *IEEE Transactions on Knowledge and Data Engineering*, 2024.
- [5] C. Ju, D. Gao, R. Mane, B. Tan, Y. Liu και C. Guan, «Federated Transfer Learning for EEG Signal Classification», στο *2020 42nd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*, 2020, σσ. 3040–3045. doi: 10.1109/EMBC44109.2020.9175344.
- [6] T. Song, Y. Tong και S. Wei, «Profit allocation for federated learning», στο *2019 IEEE International Conference on Big Data (Big Data)*, IEEE, 2019, σσ. 2577–2586.
- [7] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings κ.ά., «Advances and open problems in federated learning», *Foundations and trends® in machine learning*, τόμ. 14, αρθμ. 1–2, σσ. 1–210, 2021.
- [8] V. Shejwalkar, A. Houmansadr, P. Kairouz και D. Ramage, «Back to the Drawing Board: A Critical Evaluation of Poisoning Attacks on Production Federated Learning», στο *2022 IEEE Symposium on Security and Privacy (SP)*, 2022, σσ. 1354–1371. doi: 10.1109/SP46214.2022.9833647.
- [9] Y.-A. Xie, J. Kang, D. Niyato, N. T. T. Van, N. C. Luong, Z. Liu και H. Yu, «Securing federated learning: A covert communication-based approach», *IEEE Network*, 2022.
- [10] L. Lyu, H. Yu και Q. Yang, «Threats to federated learning: A survey», *arXiv preprint arXiv:2003.02133*, 2020.
- [11] H. Lee, J. Kim, R. Hussain, S. Cho και J. Son, «On defensive neural networks against inference attack in federated learning», στο *ICC 2021-IEEE International Conference on Communications*, IEEE, 2021, σσ. 1–6.

- [12] L. Zhu, Z. Liu και S. Han, «Deep leakage from gradients», *Advances in neural information processing systems*, τόμ. 32, 2019.
- [13] B. Zhao, K. R. Mopuri και H. Bilen, «idlg: Improved deep leakage from gradients», *arXiv preprint arXiv:2001.02610*, 2020.
- [14] L. Melis, C. Song, E. De Cristofaro και V. Shmatikov, «Exploiting unintended feature leakage in collaborative learning», στο *2019 IEEE symposium on security and privacy (SP)*, IEEE, 2019, σσ. 691–706.
- [15] B. Hitaj, G. Ateniese και F. Perez-Cruz, «Deep models under the GAN: information leakage from collaborative deep learning», στο *Proceedings of the 2017 ACM SIGSAC conference on computer and communications security*, 2017, σσ. 603–618.
- [16] C. Dwork, A. Roth κ.ά., «The algorithmic foundations of differential privacy», *Foundations and Trends® in Theoretical Computer Science*, τόμ. 9, αρθμ. 3--4, σσ. 211–407, 2014.
- [17] M. Benmalek, M. A. Benrekia και Y. Challal, «Security of federated learning: Attacks, defensive mechanisms, and challenges», *Revue des Sciences et Technologies de l'Information-Série RIA: Revue d'Intelligence Artificielle*, τόμ. 36, αρθμ. 1, σσ. 49–59, 2022.
- [18] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal και K. Seth, «Practical secure aggregation for privacy-preserving machine learning», στο *proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, 2017, σσ. 1175–1191.
- [19] A. Shafahi, W. R. Huang, M. Najibi, O. Suci, C. Studer, T. Dumitras και T. Goldstein, «Poison frogs! targeted clean-label poisoning attacks on neural networks», *Advances in neural information processing systems*, τόμ. 31, 2018.
- [20] W. Y. B. Lim, N. C. Luong, D. T. Hoang, Y. Jiao, Y.-C. Liang, Q. Yang, D. Niyato και C. Miao, «Federated learning in mobile edge networks: A comprehensive survey», *IEEE Communications Surveys & Tutorials*, τόμ. 22, αρθμ. 3, σσ. 2031–2063, 2020.
- [21] N. Bouacida και P. Mohapatra, «Vulnerabilities in federated learning», *IEEE Access*, τόμ. 9, σσ. 63 229–63 249, 2021.
- [22] M. S. Ali, H. Tabassum και E. Hossain, «Dynamic user clustering and power allocation for uplink and downlink non-orthogonal multiple access (NOMA) systems», *IEEE access*, τόμ. 4, σσ. 6325–6343, 2016.
- [23] A. Benjebbovu, A. Li, Y. Saito, Y. Kishiyama, A. Harada και T. Nakamura, «System-level performance of downlink NOMA for future LTE enhancements», στο *2013 IEEE globecom workshops (GC Wkshps)*, IEEE, 2013, σσ. 66–70.
- [24] Y. Endo, Y. Kishiyama και K. Higuchi, «Uplink non-orthogonal access with MMSE-SIC in the presence of inter-cell interference», στο *2012 international symposium on wireless communication systems (ISWCS)*, IEEE, 2012, σσ. 261–265.
- [25] N. Zhang, J. Wang, G. Kang και Y. Liu, «Uplink nonorthogonal multiple access in 5G systems», *IEEE Communications Letters*, τόμ. 20, αρθμ. 3, σσ. 458–461, 2016.
- [26] K. Higuchi και A. Benjebbour, «Non-orthogonal multiple access (NOMA) with successive interference cancellation for future radio access», *IEICE Transactions on Communications*, τόμ. 98, αρθμ. 3, σσ. 403–414, 2015.

- [27] R. Razavi, M. Dianati και M. A. Imran, «Non-orthogonal multiple access (NOMA) for future radio access», *5G Mobile Communications*, σσ. 135–163, 2017.
- [28] M. D. Davis και S. J. Brams, *Game Theory*, <https://www.britannica.com/science/game-theory>, Encyclopedia Britannica, 10 May. 2024, Accessed 20 June 2024, 2024.
- [29] H. Peters, *Game theory: A Multi-leveled approach*. Springer, 2015.
- [30] J. C. Harsanyi, «Games with incomplete information», *The American Economic Review*, τόμ. 85, αρθμ. 3, σσ. 291–303, 1995.
- [31] J. F. Nash κ.ά., «Non-cooperative games», 1950.
- [32] M. Hamidi, H. Liao και F. Szidarovszky, «Non-cooperative and cooperative game-theoretic models for usage-based lease contracts», *European Journal of Operational Research*, τόμ. 255, αρθμ. 1, σσ. 163–174, 2016.
- [33] M. Shor, *Symmetric Game*, Dictionary of Game Theory Terms, Game Theory .net, Web accessed: 20/06/24. διεύθυν.: <https://www.gametheory.net/dictionary/Games/SymmetricGame.html>.
- [34] S.-F. Cheng, D. M. Reeves, Y. Vorobeychik και M. P. Wellman, «Notes on equilibria in symmetric games», 2004.
- [35] M. J. Osborne και A. Rubinstein, *A course in game theory*. MIT press, 1994.
- [36] N. Nisan, T. Roughgarden, É. Tardos και V. V. Vazirani, *Algorithmic Game Theory*. Cambridge, UK: Cambridge University Press, 2007, ISBN: 9780521872829.
- [37] J. F. Nash, «Equilibrium Points in N-Person Games», *Proceedings of the National Academy of Sciences*, τόμ. 36, αρθμ. 1, σσ. 48–49, 1950. DOI: 10.1073/pnas.36.1.48.
- [38] L. S. Shapley κ.ά., «A value for n-person games», 1953.
- [39] R. Gupta και J. Gupta, «Federated learning using game strategies: State-of-the-art and future trends», *Computer Networks*, τόμ. 225, σ. 109 650, 2023.
- [40] M. Hu, W. Yang, Z. Luo, X. Liu, Y. Zhou, X. Chen και D. Wu, «AutoFL: a Bayesian game approach for autonomous client participation in federated edge learning», *IEEE Transactions on Mobile Computing*, 2022.
- [41] Y. Zou, S. Feng, D. Niyato, Y. Jiao, S. Gong και W. Cheng, «Mobile device training strategies in federated learning: An evolutionary game approach», στο *2019 International Conference on Internet of Things (iThings) and IEEE Green Computing and Communications (GreenCom) and IEEE Cyber, Physical and Social Computing (CPSCom) and IEEE Smart Data (SmartData)*, IEEE, 2019, σσ. 874–879.
- [42] J. Kang, Z. Xiong, D. Niyato, H. Yu, Y.-C. Liang και D. I. Kim, «Incentive design for efficient federated learning in mobile networks: A contract theory approach», στο *2019 IEEE VTS Asia Pacific Wireless Communications Symposium (APWCS)*, IEEE, 2019, σσ. 1–5.
- [43] L. Nagalapatti και R. Narayanam, «Game of gradients: Mitigating irrelevant clients in federated learning», στο *Proceedings of the AAAI Conference on Artificial Intelligence*, τόμ. 35, 2021, σσ. 9046–9054.
- [44] Z. Tang, F. Shao, L. Chen, Y. Ye, C. Wu και J. Xiao, «Optimizing federated learning on non-IID data using local Shapley value», στο *Artificial Intelligence: First CAAI International Conference, CICA 2021, Hangzhou, China, June 5–6, 2021, Proceedings, Part II 1*, Springer, 2021, σσ. 164–175.

- [45] Y. Shi και Y. E. Sagduyu, «Jamming attacks on federated learning in wireless networks», *arXiv preprint arXiv:2201.05172*, 2022.
- [46] Y. Shi και Y. E. Sagduyu, «How to launch jamming attacks on federated learning in NextG wireless networks», στο *2022 IEEE Globecom Workshops (GC Wkshps)*, IEEE, 2022, σσ. 945–950.
- [47] R. Ruby, H. Yang και K. Wu, «Anti-jamming strategy for federated learning in internet of medical things: A game approach», *IEEE Journal of Biomedical and Health Informatics*, τόμ. 27, αρθμ. 2, σσ. 888–899, 2022.
- [48] R. Gupta και J. Gupta, «Federated learning using game strategies: State-of-the-art and future trends», *Computer Networks*, τόμ. 225, σ. 109 650, 2023.
- [49] E. Tahanian, M. Amouei, H. Fateh και M. Rezvani, «A game-theoretic approach for robust federated learning», *International Journal of Engineering*, τόμ. 34, αρθμ. 4, σσ. 832–842, 2021.
- [50] Z. Wang, *Traditional Federated Learning*, <https://github.com/wzljerry/Hierarchical-Federated-Learning/blob/main/FL.ipynb>, 2022.
- [51] Y. LeCun, C. Cortes και C. Burges, «MNIST handwritten digit database», *AT&T Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, τόμ. 2, 2010.

Απόδοση

Ομόσπονδη Μάθηση
 Διακομιστής, Εξυπηρετητής
 Συσκευές άκρων
 Συγκεντρωτικό Μοντέλο
 Ακρίβεια, Ποσοστό Επιτυχίας
 Μη Ορθογωνική Πολλαπλή Πρόσβαση
 Ρυθμός μεταφοράς
 Κέρδος καναλιού
 Χρονική επιβάρυνση
 Ενεργειακή επιβάρυνση
 Παρεμβολή
 Χρησιμότητα
 Συνελκτικό Νευρωνικό Δίκτυο
 Ετικέτα, επισήμανση
 Σύνολο δεδομένων
 Σύνολο αξιολόγησης
 Σηματοθορυβικός λόγος
 Σταθμός βάσης
 Κατώφλι
 Ανερχόμενη ζεύξη
 κατερχόμενη ζεύξη
 Διαδοχική Ακύρωση Παρεμβολών
 Λευκός Προσθετικός Θόρυβος Γκάους
 Ισοροπία Ναs

Ξενόγλωσσος όρος

Federated Learning
 Server
 Edge Devices
 Aggregated Model
 Accuracy
 Non Orthogonal Multiple Access (NOMA)
 data rate
 channel gain
 time overhead
 energy overhead
 jamming
 utility
 convolutional neural network (CNN)
 label
 dataset
 test set
 Signal Noise Ratio (SNR)
 Base Station (BS)
 threshold
 uplink
 downlink
 Successive Interference Cancellation (SIC)
 Additive White Gaussian Noise (AWGN)
 Nash Equilibrium

Στην ενότητα 4 παρουσιάστηκε η συνάρτηση χρησιμότητας των χρηστών 4.14. Η συνάρτηση αυτή οφείλει να είναι κοίλη (να έχει μοναδικό μέγιστο) όταν ο εκάστοτε χρήστης μπορεί να μεταδώσει το μήνυμά του, δηλαδή για $P_{thres} < P_{max}$. Στο παράρτημα αυτό πραγματοποιείται η απόδειξη κοιλότητας.

Ισχυρισμός 1. Η συνάρτηση U_i είναι κοίλη ως προς την p_i για κάθε \mathbf{p}_{-i}

Απόδειξη. Έστω

$$f(p_i) = U_i(p_i^t, \mathbf{p}_{-i}^{t-1}) = c_1 a_i^{t-1} + c_2 \frac{T - T_i(p_i, \mathbf{p}_{-i}^{t-1})}{T} - c_3 p_i^t \quad (7.1)$$

για $P_{thres} < P_{max}$, όπου το διάνυσμα \mathbf{p}_{-i}^{t-1} είναι σταθερό και ανεξάρτητο της μεταβλητής απόφασης p_i του γύρου t . Οι σταθερές a_i^{t-1} , c_1 , c_2 , c_3 , T είναι επίσης ανεξάρτητες του p_i και θετικές, οπότε με παραγωγήσι ως προς p_i έχουμε:

$$\begin{aligned} \frac{df}{dp_i} &= \frac{d(c_1 a_i^{t-1} + c_2 \frac{T - T_i(p_i, \mathbf{p}_{-i}^{t-1})}{T} - c_3 p_i^t)}{dp_i} \\ &= 0 + \frac{c_2}{T} \frac{d(T - T_i(p_i, \mathbf{p}_{-i}^{t-1}))}{dp_i} - c_3 \\ &= \frac{c_2}{T} \frac{d(T_i(p_i, \mathbf{p}_{-i}^{t-1}))}{dp_i} - c_3 \end{aligned}$$

Συνεπώς, προκειμένου να αποδειχθεί η κοιλότητα της συνάρτησης $f(p_i)$, αρκεί να αποδειχθεί ότι ο όρος $-T_i$ είναι κοίλος ως προς την μεταβλητή p_i . Με βάση τις εξισώσεις 4.9, 4.10 και αντικαθιστώντας ως $I_i = \sum_{j=1}^{i-1} (G_j P_j)$, προκύπτει:

$$g(x) = -T_i(p_i, \mathbf{p}_{-i}^{t-1}) = \frac{-Z(\mathbf{W}_i)}{R_i} = \frac{-Z(\mathbf{W}_i)}{B \log_2(1 + \frac{G_i P_i}{I_i + I_0})}$$

Συνεπώς η παράγωγος της συνάρτησης $g(x)$ ως προς p_i ισούται με:

$$\frac{dg}{dx} = \frac{Z(\mathbf{W}_i) \log(2)}{B \log^2(1 + \frac{G_i P_i}{I_i + I_0})} \frac{1}{(1 + \frac{G_i P_i}{I_i + I_0})} \frac{G_i}{I_i + I_0}$$

και η δεύτερη παράγωγος:

$$\frac{d^2 g}{dx^2} = -2 \frac{Z(\mathbf{W}_i) \log(2)}{B \log^3(1 + \frac{G_i P_i}{I_i + I_0})} \frac{1}{(1 + \frac{G_i P_i}{I_i + I_0})^2} \frac{G_i^2}{(I_i + I_0)^2} - \frac{Z(\mathbf{W}_i) \log(2)}{B \log^2(1 + \frac{G_i P_i}{I_i + I_0})} \frac{1}{(1 + \frac{G_i P_i}{I_i + I_0})^2} \frac{G_i^2}{(I_i + I_0)^2}$$

Δεδομένου ότι όλες οι σταθερές είναι θετικές, ως φυσικά μεγέθη που εκφράζουν κέρδος καναλιού, πλήθος δεδομένων και ισχύ και αφού το όρισμα όλων των λογαρίθμων είναι μεγαλύτερο του 1, τότε είναι εμφανές ότι:

$$\frac{d^2 g}{dx^2} < 0, \quad \forall p_i \in [0, P_{max}]$$

Συνεπώς η $g(x)$ είναι κοίλη σε όλο το πεδίο ορισμού της, με αποτέλεσμα και η συνάρτηση χρησιμότητας $U_i(p_i^t, \mathbf{p}_{-i}^{t-1})$ να είναι κοίλη, όταν ένας χρήστης έχει επιτυχή μετάδοση. ■

