



**ΠΑΝΕΠΙΣΤΗΜΙΟ ΠΕΙΡΑΙΩΣ
ΤΜΗΜΑ ΨΗΦΙΑΚΩΝ ΣΥΣΤΗΜΑΤΩΝ**

Π.Μ.Σ. «Τεχνοοικονομική Διοίκηση και Ασφάλεια Ψηφιακών Συστημάτων»

**ΜΕΤΑΠΤΥΧΙΑΚΗ ΔΙΠΛΩΜΑΤΙΚΗ ΕΡΓΑΣΙΑ
k-Ανωνυμία, Μια προσέγγιση στην προστασία των δεδομένων και της
ιδιωτικότητας**

Όνομα: Θρασύβουλος

Επώνυμο: Πιστοφίδης

A.M. ΜΤΕ 1532

ΕΠΙΒΛΕΠΩΝ ΚΑΘΗΓΗΤΗΣ:

Λαμπρινουδάκης Κωνσταντίνος

Τριμελής Εξεταστική Επιτροπή:

Κωνσταντίνος Λαμπρινουδάκης, Καθηγητής

Χρήστος Ξενάκης, Καθηγητής

Χριστόφορος Νταντογιάν, Καθηγητής

Πειραιάς Φεβρουάριος 2018

Ευχαριστίες

Κατ' αρχήν, θα ήθελα να εκφράσω τις ευχαριστίες μου στον επιβλέποντα καθηγητή μου Δρ. Κωνσταντίνο Λαμπρινουδάκη, τόσο για την ανάθεση της εν λόγω εργασίας όσο για την υποστήριξη, υπομονή και καθοδήγηση του κατά την διάρκεια της σύνταξης της. Η συνεργασία μας ήταν άψογη τόσο κατά την διάρκεια του Μεταπτυχιακού όσο και στο τέλος. Εν συνεχεία, θα ήθελα να ευχαριστήσω τον Δρ. Χρήστο Ξενάκη, Δρ. Χριστόφορο Νταντογιάν και τον Δρ. Παναγιώτη Ριζομυλιώτη για τις παραγωγικές συζητήσεις και το ενδιαφέρον τους.

Άξια αναφοράς, και η άριστη συνεργασία μου με την γραμματεία κα. Σκούντζου Σοφία.

Τέλος, ένα ευχαριστώ στην οικογένεια μου που συνεχίζει να με στηρίζει σε όλες μου τις επιλογές και αποφάσεις.

Περίληψη

Με την καλπάζουσα αύξηση της τεχνολογίας, αυξάνεται η σχέση και η εξάρτηση μας απο αυτή. Το διαδίκτυο, ως ο απόλυτος πλέον εκφραστής αυτής, γίνεται ο δρόμος της ανταλλαγής, επεξεργασίας και διαχείρισης της πληροφοροφορίας τόσο σε ατομικό επίπεδο ως χρήστες, όσο σε μαζικό ως εταιρείες, οργανισμοί και ιδρύματα. Η πληροφορία διαχέεται ακατάπαυστα και αποθηκεύεται παντού. Τόσο οι καταχωρήσεις όσο και τα μεταδεδομένα που αφήνουμε πίσω. Έτσι η πληροφορία, και κατ επέκταση τα Big Data αποκτούν το χαρακτηρισμό του "μαύρου χρυσού" του 21 Αιώνα.

Σε ένα τέτοιο αμυδρό περιβάλλον χρειάζονται κανόνες και μέτρα ασφάλειας, προκειμένου να διασφαλιστεί αρχικώς η ιδιωτικότητα των ατόμων, και κατ επέκταση η διασφάλιση των εταιρειών απο επιθέσεις. Σε Ευρωπαϊκό επίπεδο η θέσπιση του GDPR έχει θετικό αντίκτυπο τόσο στην ενημέρωση - συνειδητοποίηση μας, όσον αφορά το εμπόριο των δεδομένων αλλά και την συμμόρφωση των εταιρειών, οργανισμών ως προς την διαχείριση τους. Ένας απο τους τρόπους που μπορούν να διασφαλίσουν τα παράπανω είναι η ανωνυμία.

Στην εν λόγω διπλωματική θα ασχοληθούμε αναφορικά με τους αλγόριθμους ανωνυμίας, και θα επικεντρωθούμε σε έναν απο τους πρώτους -αν όχι, πρώτο- την k-ανωνυμία όπως επίσης και θα αναλύσουμε τον τρόπο λειτουργίας ενός απο πιο σημαντικά πρόγράμματα εφαρμογής της, σε μεγάλες βάσεις δεδομένων, το ARX.

Λέξεις κλειδιά: k-ανωνυμία, I-ποικιλομορφία, t-εγγύτητα, ιδιωτικότητα δεδομένων, ανωνυμία, απόρρητο, ανωνυμοποίηση, γενίκευση, απόκρυψη, ψευδοαναγνωριστικά

Πίνακας Περιεχομένων

Κεφάλαιο 1. Εισαγωγή	7
1.1 Τα είδη πληροφορίας	7
1.2 Βάσεις Δεδομένων : Ο Τρόπος παρουσίασης και διαχείρισης της πληροφορίας	8
1.2.1 Σχεσιακές βάσεις δεδομένων.	8
1.3 Απειλές στην ιδιωτικότητα	9
Κεφάλαιο 2. Ο ορισμός της k-ανωνυμίας	11
2.1 Quasi identifiers – Ψευδοαναγνωριστικά	13
Κεφάλαιο 3. Μηχανισμοί που υποστηρίζουν την k-ανωνυμία	16
3.1. Γενίκευση (Generalization)	16
3.2. Απόκρυψη (Suppression)	16
3.2.1. Ολική Γενίκευση - Global Recoding	17
3.2.2 Τοπική Γενίκευση - Local Recoding	17
Κεφάλαιο 4. Αλγόριθμοι k-Ανωνυμίας	19
4.1 Αλγόριθμος Incognito	19
4.2 Αλγόριθμος Mondrian	19
4.3 Μη-ομοιογενής γενίκευση	20
Κεφάλαιο 5. Επιθέσεις στην k-ανωνυμία	22
Κεφάλαιο 6. Σύνθετες επιθέσεις στην k-ανωνυμία και τρόποι επίλυσης	25
6.1 Εισαγωγή	25
6.2. Επίθεση της Ομοιογένειας	25
6.3. Επίθεση με πρότερη γνώση	27
6.4. I – Ποικιλομορφία (I-diversity)	28
6.4.2 Περιορισμοί I – ποικιλομορφίας	28
6.5 t- εγγύτητα	30
6.6 δ – Παρουσία (δ-Presence)	33
6.7. Ανωνυμοποίηση βάση κινδύνου	33
Κεφάλαιο 7. Εργαλεία Ανωνυμοποίησης	34
7.1 UTD	34
7.2 CAT	34
7.3 TIAMAT	34
7.4 SECRETA	35
7.5 ARX	35
7.5.1 Η Ανάλυση κινδύνου και η Ανωνυμοποίηση βάση κινδύνου με την χρήση του ARX	38
7.5.2 Αξιολόγηση χρησιμότητας	39
7.5.3 Επιπρόσθετα Χαρακτηριστικά	39
7.5.4 Χρήση του ARX	40
7.5.4.1 Η διαδικασία της ανωνυμοποίησης	40
7.5.4.2 Παραμετροποιώντας την διαδικασία ανωνυμοποίησης	41
7.5.4.3 Διαχείριση του χώρου των αποτελεσμάτων	44
7.5.4.4 Αξιολόγηση Χρηστικότητας Δεδομένων	45

7.5.4.5 Ανάλυση κινδύνων επαναπροσδιορισμού	46
7.6 Τελικές παρατηρήσεις	47
Κεφάλαιο 8. Επίλογος	48
Βιβλιογραφία	49

Κεφάλαιο 1. Εισαγωγή

Στην ψηφιακή εποχή, η ραγδαία αύξηση της δυνατότητας συλλογής, διαχείρισης και ανταλλαγής της πληροφορίας αυξάνει τις ανησυχίες περί προστασίας της ιδιωτικότητας. Οι υπολογιστές μπορούν πλέον να χρησιμοποιηθούν για να συλλέξουν και επεξεργαστούν συνήθειες, τρόπο ζωής, τοποθεσίας, όπως της συσχέτισης όλων αυτών από δεδομένα που συλλέγονται σε καθημερινή βάση. Αυτό δημιουργεί έναν αμφιλεγόμενο συμβιβασμό μεταξύ της αξίας (τόσο για την κοινωνία όσο για τα άτομα της) της διαθέσιμης, διαχέουσας γνώσης, με τον κίνδυνο που θα προκαλέσει για τα άτομα, η αποκάλυψη ή κακή χρήση των ιδιωτικών τους δεδομένων. Μια λύση σε αυτό το πρόβλημα είναι η ανωνυμία.

Ο όρος ανωνυμία αναφέρεται στην απουσία των αναγνωριστικών πληροφοριών ενός ατόμου. Ωστόσο, μια απλή κατάργηση των πιο ξεκάθαρων – σαφών αναγνωριστικών των ατόμων μπορεί να μην αποτελεί επαρκή προστασία.

Ο κύριος λόγος είναι ότι όταν οι πληροφορίες βγούν στην επιφάνεια, και συνδυαστούν με τις ήδη υπάρχουσες δημοσιευμένες, μπορούν να αποκαλύψουν την ταυτότητα ενός ατόμου.

Ένα χαρακτηριστικό και αρκετά γνωστό παράδειγμα είναι με τον Crowdsourcing διαγωνισμό του Netflix. Το 2012, το Netflix δημοσίευσε μια λίστα με τους χρήστες και τις βαθμολογίες τους. Ο κόσμος μπορούσε απλά να κατεβάσει αυτή την βάση και να ψάξει για μοτίβα - patterns. Τα δεδομένα περιλάμβαναν ψεύτικα αναγνωριστικά δεδομένα (customer ID), μαζί με τις ταινίες, τις βαθμολογίες τους και την ημέρα της βαθμολογίας. Υποστηρίζεται ότι από την στιγμή που τα αναγνωριστικά των πελατών έχουν αφαιρεθεί, οι εν λόγω πληροφορίες δεν παραβιάζουν το απόρρητο των χρηστών. Ωστόσο, οι Narayanan και Shmatikov (2008) έδειξαν πως οι πελάτες μπορούν να ταυτοποιηθούν όταν αυτή η -Netflix- λίστα συνδυαστεί με κάποια βοηθητικά δεδομένα (όπως π.χ από το IMDB). [1]

Ένα άλλο τρανταχτό παράδειγμα είναι αυτό της AOL με ουσιαστική πληροφορία αυτή την φορά τα log των αναζητήσεων των μηχανών αναζήτησης.

Τα logs αναζητήσεων των μηχανών αναζήτησης είναι τεράστια πηγή πληροφορίας τόσο για τους ερευνητές όσο και για τις εταιρείες marketing, όμως παράλληλα η δημοσίευσή τους ενδέχεται να εκθέσει την ιδιωτικότητα των χρηστών εκ των οποίων τα logs δημιουργούνται. Υπάρχει το λιγότερο μια τρανταχτή περίπτωση τέτοιας εξέθεσης των προσωπικών logs αναζητήσεων χρηστών όπου εφαρμόστηκε πολύ απλή ανωνυμοποίηση, η οποία είχε σαν αποτέλεσμα να αποκαλύψει αρκετές πληροφορίες για την αναγνώρισή τους. Αυτή η αποκάλυψη έγινε από την AOL σε μια προσπάθεια να βοηθήσει την “Ερευνητική κοινότητα ανάκτησης πληροφοριών” που έφερε σαν αποτέλεσμα να προκαλέσει ισχυρό πλήγμα τόσο στην πολιτική απορρήτου των χρηστών της AOL όσο και στην ίδια την AOL, με αγωγές και ενστάσεις εναντίον της.

Ιδανικά τα logs αναζήτησης θα πρέπει να ανωνυμοποιούνται πριν κοινοποιηθούν δημοσίως. Το πρόβλημα είναι ότι η επιτευξη του επιθυμητού επιπέδου ιδιωτικότητας στα logs καθίσταται δύσκολη, παρουσιάζοντας παράλληλα αντιστάθμισμα μεταξύ χρηστικότητας και ιδιωτικότητας των δεδομένων. Υπάρχουν αρκετές προσεγγίσεις για την ανωνυμοποίηση τέτοιων δεδομένων, αλλά συνήθως περιορίζονται στην διαγραφή συγκεκριμένων queries ή logs. Επιπλέον, συνηθισμένες τεχνικές που χρησιμοποιούνται σε στατιστικούς ελέγχους αποκάλυψης δεν χρησιμοποιήθηκαν ποτέ, παρα μόνο μέχρι πρόσφατα για τέτοιου είδους περιπτώσεις.

Ιδια λογική ισχύει και σε κάθε είδους δεδομένα, όπως τα χρονικής σήμανσης (timestamped) που παρατηρούνται σε περιπτώσεις όπως αποθήκευση τραπεζικών συναλλαγών, ιατρικές εξετάσεις, λογιστικά βιβλία αρχεία ασφάλισης. Δεδομένα δηλαδή προς διακίνηση, όπου ενδέχεται να ζητηθούν προς δημόσια κοινοποίηση

Τα προηγούμενα παραδείγματα εγείρουν την ερώτηση: τι είδους πληροφορία επιθυμούμε να προστατεύσουμε όταν μιλάμε για προστασία της ιδιωτικότητας και πως κατατάσσονται οι απειλές της

1.1 Τα είδη πληροφορίας

Με την αναφορά των ειδών πληροφορίας που επιθυμούμε να προστατεύσουμε αναφερόμαστε συνήθως στα δεδομένα χρηστών/πελατών. Αυτά δομούνται σε σχέση με τον τύπο τους, την ευαισθησία και τις συνέπειες τους στην ασφάλεια, ιδιωτικότητα και ανωνυμία. Κατηγοριοποιούνται σε γενικότερο πλαίσιο στους εξής ακόλουθους 3 τύπους:

- i) Δεδομένα αναγνώρισης χρήστη: Πρόκειται για τα δεδομένα και τις παραμέτρους που προσδιορίζουν τον χρήστη, όπως το όνομα/επώνυμο, το όνομα εισόδου (αναγνωριστικό/username), ο κωδικός, το e-mail του και η IP του
- ii) Μεταδεδομένα: Πρόκειται για δεδομένα που ορίζουν κάποια ιδιότητα του χρήστη, όπως τοποθεσία, ημερομηνία γεννήσεως, φύλο, κριτήρια αναζήτησης, και log files
- iii) Δεδομένα περιεχομένου: Πρόκειται για δεδομένα που περιέχονται σε μια συναλλαγή, όπως πληροφορίες πληρωμής, το περιεχόμενο ενός e-mail και μιας ιστοσελίδας

Στο ειδικότερο πλαίσιο της βασικής ανωνυμοποίησης, και ομαδοποίησης τους σε πίνακες, η διασύνδεση θα γίνει κριτήριο τού τι θα ορίζεται ως ευαίσθητο, μη ευαίσθητο και ψευδοαναγνωρίσιμο.

1.2 Βάσεις Δεδομένων : Ο Τρόπος παρουσίασης και διαχείρισης της πληροφορίας

1.2.1 Σχισιακές βάσεις δεδομένων.

Με τον όρο σχεσιακή βάση δεδομένων εννοείται μία συλλογή δεδομένων οργανωμένη σε συσχετισμένους πίνακες που παρέχει ταυτόχρονα ένα μηχανισμό για ανάγνωση, εγγραφή, τροποποίηση ή και πιο πολύπλοκες διαδικασίες πάνω στα δεδομένα. [14]

Κάθε γραμμή ονομάζεται tuple και περιλαμβάνει ένα σύνολο πληροφοριών που σχετίζονται με ένα άτομο. Οι στήλες χωρίζουν τα δεδομένα σε σημασιολογικές κατηγορίες που ονομάζονται γνωρίσματα. Ένα σύνολο δεδομένων – dataset αναφέρεται σε ένα μονο tuple σε έναν συγκεκριμένο πίνακα.

Για να είμαστε πιο συγκεκριμένοι βάση αναφοράς της εκθεσης της Sweeney ένας πίνακας T αναφέρεται σαν $T(A_1, \dots, A_n)$ με τα γνωρίσματα του $\{A_1, \dots, A_n\}$. Ένας διατεταγμένος n -tuple $[d_1, d_2, \dots, d_n]$ περιλαμβάνει τις τιμές που συσχετίζονται με τα γνωρίσματα του πίνακα. Για κάθε $j=1, 2, \dots, n$ η τιμή του d_j ανατίθεται στο γνώρισμα A_j .

Κάθε $T[A_1, \dots, A_j]$ σημαίνει την προβολή του T , συμπεριλαμβάνει μόνο τις τιμές A_1, \dots, A_j

1.3 Απειλές στην ιδιωτικότητα

Απειλή στην ιδιωτικότητα, κατα γενική προσέγγιση, έχουμε όταν ο επιτιθέμενος είναι ικανός να συσχετίσει την ταυτότητα ενός χρήστη με την πληροφορία που ο χρήστης θεωρεί ιδιωτική.

Ακολουθούν, οι τρεις πιο κοινοί τύποι απειλών της ιδιωτικότητας[2]:

1. Αποκάλυψη της ιδιότητας μέλους: σημαίνει ότι η διασύνδεση των δεδομένων επιτρέπει σε έναν επιτιθέμενο να καθορίσει το κατα πόσο τα δεδομένα ενός συγκεκριμένου ατόμου περιέχονται σε ένα σύνολο δεδομένων ή όχι. Ενώ αυτό άμεσα δεν αποκαλύπτει καμία πληροφορία από το ίδιο το σύνολο δεδομένων, μπορεί να επιτρέψει στον επιτιθέμενο να συναγάγει πορίσματα και μεταδεδομένα.

Στο παράδειγμα μας (πίνακας 1), τα δεδομένα προέρχονται από μητρώο καρκινοπαθών και μπορεί να συναχθεί ότι ένα άτομο έχει ή είχε καρκίνο. Αυτό αφορά τα έμμεσα ευαίσθητα χαρακτηριστικά (εννοώντας χαρακτηριστικά ενός ατόμου που δεν περιλαμβάνεται στο σύνολο δεδομένων), άλλα μοντέλα γνωστοποίησης, αφορούν απόλυτα, ευαίσθητα χαρακτηριστικά.

2. Αποκάλυψη χαρακτηριστικών: μπορεί να επιτευχθεί χωρίς απαραίτητης της σύνδεσης ενός ατόμου με ένα συγκεκριμένο στοιχείο σε ένα σύνολο δεδομένων. Προστατεύει ευαίσθητα χαρακτηριστικά, τέτοια τα οποία δεν θα ήθελε κανένα άτομο να συσχετίζεται. Ως εκ τούτου μπορεί να ενδιαφέρουν τον επιτιθέμενο, και αν αποκαλυφθούν θα προκαλέσουν κακό στα υποκείμενα των δεδομένων. Για παράδειγμα η σύνδεση σε ένα σύνολο καταχωρήσεων δεδομένων επιτρέπει την εξαγωγή πληροφοριών εάν όλα τα στοιχεία μοιράζονται μια ορισμένη τιμή ευαίσθητου χαρακτηριστικού (π.χ καρκίνος του μαστού)

3. Αποκάλυψη ταυτότητας: (ή επαναπροσδιορισμός) σημαίνει ότι ένα άτομο μπορεί να συνδεθεί με μια συγκεκριμένη καταχώρηση δεδομένων. Πρόκειται για πολύ σοβαρό είδος επίθεσης, καθώς έχει νομικές συνέπειες για τους κατόχους δεδομένων σύμφωνα με πολλούς νόμους και κανονισμούς παγκοσμίως. Από τον ορισμό προκύπτει επίσης ότι ένας επιτιθέμενος μπορεί να μάθει όλες τις ευαίσθητες πληροφορίες που εμπεριέχονται στην καταχώρηση δεδομένων για το άτομο.

Directly identifying		Quasi-identifying		Insensitive	Sensitive
Firstname	Lastname	Age	Gender	State	Diagnosis
Bradley	Rider	51	Male	NY	Colon cancer
Michael	Harlow	45	Male	MS	Hodgkin disease
Adella	Bartram	63	Female	NY	Breast cancer
Freya	King	78	Female	TX	Breast cancer
Laurena	Milton	81	Female	AL	Breast cancer

Memberships disclosure: { Firstname, Lastname, Age, Gender, State, Diagnosis }

Identity disclosure: { Firstname, Lastname, Age, Gender, State, Diagnosis }

Attribute disclosure: { Diagnosis }

Πίνακας 1 : Τύποι χαρακτηριστικών και τύποι αποκάλυψης

Με την πάροδο των χρόνων, η ερευνητική κοινότητα στην προσπάθεια της να αντιμετωπίσει τις προαναφερόμενες απειλές ανέπτυξε διάφορα μοντέλα ιδιωτικότητας, συμπεριλαμβανομένης της k-ανωνυμίας (k-anonymity : Sweeney, 2002) και της διαφορικής ιδιωτικότητας (Differential Privacy : Dwork, 2006). Στην παρούσα διπλωματική θα δοθεί έμφαση στην k-ανωνυμία, στις αδυναμίες και στις εναλλακτικές προτάσεις που ήρθαν αργότερα για να προσφέρουν επίλυση στις επιμέρους αδυναμίες της.

Κεφάλαιο 2. Ο ορισμός της k-ανωνυμίας

Κατα Sweeney:

"Given person-specific field-structured data, produce a release of the data with scientific guarantees that the individuals who are the subjects of the data cannot be re-identified while the data remain practically useful"

Η k-ανωνυμία, που προτάθηκε από την Sweeney (2002), είναι μια “αρχή προστασίας” (property) των δημοσιευμένων πληροφοριών από την επαλήθευση (reidentification). Μπορεί να χρησιμοποιηθεί για παράδειγμα, όταν μια ιδιωτική εταιρεία, όπως μια Τράπεζα που επιθυμεί να εκδόσει – μοιραστεί μια εκδοχή μιας βάσης δεδομένων πελατών που αφορά οικονομικές πληροφορίες σε κάποιους δημόσιους οργανισμούς για ερευνητικούς σκοπούς.

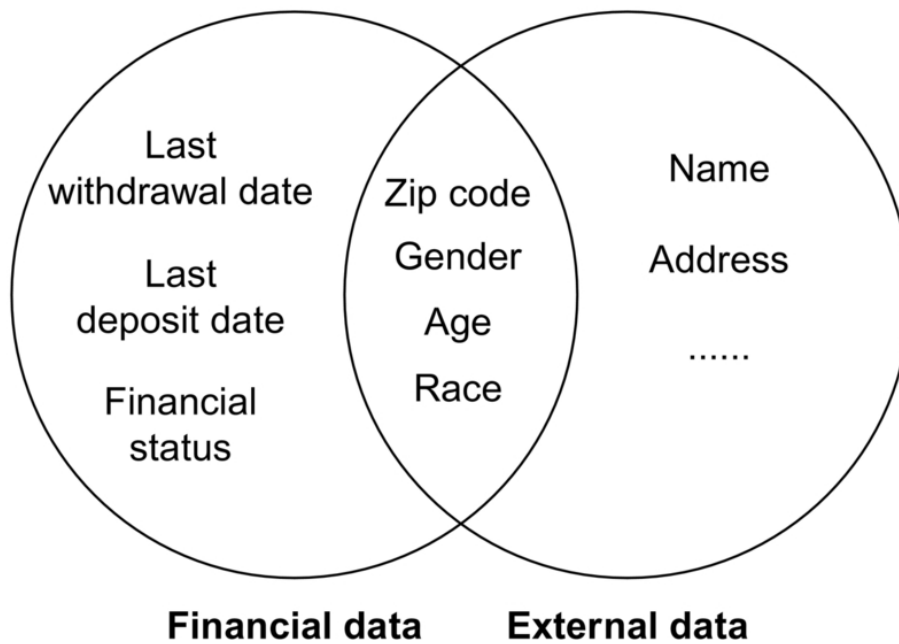
Σε μια τέτοια περίπτωση τα δεδομένα που θα δημοσιοποιηθούν θα πρέπει να έχουν την ιδιότητα της μη επαληθευσιμότητας, αναγνώρισης δηλαδή των μεμονομένων υποκειμένων τους προκειμένου να επιτευχθεί η προστασία της ιδιωτικότητας τους. Με άλλα λόγια, όλα τα δεδομένα που θα αναφέρονται στην βάση δεδομένων θα πρέπει να είναι μη – συσχετιζόμενα, ανεξάρτητα των πελατών.

Τα αρχικά δεδομένα των πελατών μιας Τράπεζας συνήθως περιέχουν πληροφορίες όπως:

Όνομα, Διεύθυνση, και Τηλεφωνικό αριθμό όπου μπορούν άμεσα να ταυτοποιήσουν τους πελάτες. Ένας πιθανός τρόπος απόκρυψης της ταυτότητας τους είναι η άμεση αφαίρεση των ευαίσθητων αυτών πληροφοριών από την βάση. Ωστόσο, αυτό δεν αποτελεί εξασφάλιση της ιδιωτικότητας των πελατών.

Οι πληροφορίες, όπως ο ταχυδρομικός κώδικας, το φύλλο, η ηλικία και η φυλή των πελατών μπορούν να τους ταυτοποιήσουν. Ο T.K (ταχυδρομικός κώδικας) προσφέρει μια κατά προσέγγιση τοποθεσία. Αναζητώντας, βάση συγκεκριμένων κριτηρίων όπως της ηλικίας, του γένους και της φυλής, είναι πολύ πιθανό να αποκαλυφθούν οι ταυτότητες των πελατών.

Ένας άλλος πιθανός τρόπος επίτευξης επαληθευσιμότητας ονομάζεται “επίθεση διασύνδεσης” (linking attack). Εκτός των γνωρίσματος όπως το όνομα και η διεύθυνση όπου μπορούν άμεσα να “σπάσουν” την ανωνυμία των δεδομένων, υπάρχουν επίσης γνωρίσματα που ονομάζονται ψευδοαναγνωριστικά – (quasi-identifier) QID που χρησιμοποιούνται για να συνδέσουν τα δημοσιευμένα δεδομένα με εξωτερικά δεδομένα. Το φύλο, η ηλικία, ο ταχυδρομικός κώδικας είναι μια τυπική **tuple** (η πλειάδα θα χαρακτηρίζεται ως tuple) και αυτή - η tuple - από τα δημοσιευμένα δεδομένα έχει αυξημένη πιθανότητα εμφάνισης εξίσου σε κάποια εξωτερικά δεδομένα. Δηλαδή, εάν υπάρχουν εξωτερικοί πίνακες όπως, λίστες καταχώρησης ψηφοφόρων, τότε συνδέοντας τα με τα QID των δημοσιευμένων δεδομένων, μπορούν να αποκαλυφθούν οι ταυτότητες των πελατών (σχεδ. 1).[1]



σχεδ.1 Επίθεση της διασύνδεσης μεταξύ των δημοσιευμένων δεδομένων και των εξωτερικών δεδομένων

Η *k*-ανωνυμία ως πρόταση απαιτεί κάθε εγγραφή των δημοσιευμένων δεδομένων, να μπορεί να αντιστοιχιστεί σε τουλάχιστον *k* εγγραφές των αρχικών δεδομένων. Με άλλα λόγια, κάθε εγγραφή από τα δημοσιευμένα δεδομένα θα έχει τουλάχιστον *k*-1 ίδιες εγγραφές στα ίδια δημοσιευμένα δεδομένα. Για παράδειγμα στον Πίνακα 1, ο (a) είναι τα αυθεντικά δεδομένα και (b) είναι τα δεδομένα που προκύπτουν από τον (a). Ο (b) έχει *k*-ανωνυμία όπου $k = 2$.

Η Latanya Sweeney [3] παρουσίασε την αρχή της *k*-ανωνυμίας και απέδειξε ότι εάν τα δημοσιευμένα δεδομένα εφαρμόζουν την ιδιότητα της *k*-ανωνυμίας, τότε μπορούν να προστατευθούν από την “Επίθεση της διασύνδεσης”. Αυτό γιατί κάθε εγγραφή των δημοσιευμένων δεδομένων θα έχει τουλάχιστον *k*-1 ίδιες εγγραφές.

Άρα, σκοπός της *k*-ανωνυμίας είναι να καταστήσει κάθε μια εγγραφή μη-διάκριτη ανάμεσα σε άλλες *k*-1.

Παράδειγμα k-ανωνυμίας (k=2)

Πίνακας 1

(a) Αρχικά Δεδομένα				
Όνομα	Γένος	Φυλή	Ηλικία	T.K
Αλίκη	Θηλυκό	Λευκή	17	21103
Su	Θηλυκό	Κίτρινη	22	21300
Denzel	Αρσενικό	Μαύρη	27	21110
Κατερίνα	Θηλυκό	Λευκή	15	21102
Rose	Θηλυκό	Μαύρη	29	21109
Andy	Αρσενικό	Κίτρινη	24	21304

(b) Διαμοιραζόμενα δεδομένα προερχόμενα απο (a)			
Γένος	Φυλή	Ηλικία	T.K
Θ ή Α	Λευκή	15 - 19	211*
Θ ή Α	Κίτρινη	20 - 24	213*
Θ ή Α	Μαύρη	25 - 29	211*
Θ ή Α	Λευκή	15 - 19	211*
Θ ή Α	Μαύρη	25 - 29	211*
Θ ή Α	Κίτρινη	20 - 24	213*

2.1 Quasi identifiers – Ψευδοαναγνωριστικά

Μια ομάδα απο γνωρίσματα που δεν είναι δομικά μοναδικά, αλλά θα μπορούσαν να είναι εμπειρικά μοναδικά, συνεπώς κατα βάση θα μπορούσαν να αναγνωρίζουν – προσδιορίζουν μια πληθυσμιακή μονάδα.

	Quasi – Identifier - Ψευδοαναγνωριστικά			Ευαίσθητα Δεδομένα
	Ηλικία	Γένος	T.K	Έσοδα
1	35	Αρσενικό	81243	300,000
2	48	Θηλυκό	83123	30,000
3	40	Αρσενικό	81205	1,000,000
4	60	Αρσενικό	73193	100,000
5	27	Θηλυκό	83123	60,000
6	60	Αρσενικό	71234	20,000
7	27	Θηλυκό	83981	25,000
8	35	Θηλυκό	83012	30,000
9	27	Αρσενικό	81021	40,000
10	46	Αρσενικό	73013	25,000
11	46	Θηλυκό	83561	70,000
12	40	Αρσενικό	81912	40,000
13	48	Αρσενικό	72231	1,500,000

Πίνακας 1a: Ιδιωτικός πίνακας P

Ο παραπάνω πίνακας είναι ιδιωτικός καθώς δεν συμπεριλαμβάνει προσωπικά δεδομένα όπως Ονοματεπώνυμο. Το εισόδημα κατατάσσεται στα ευαίσθητα δεδομένα. Δεν θα πρέπει να είναι εφικτή η αναγνώριση του εισοδήματος ενός ατόμου χρησιμοποιώντας αυτόν τον πίνακα.

Εάν συγκεκριμένα δεδομένα όπως η ηλικία, το γένος, και ο Τ.Κ είναι μοναδικοί τότε μπορεί να χρησιμοποιηθούν για την αναγνώριση ενός ατόμου και ανάλογα του εισοδήματος του. Μπορούν να χρησιμοποιηθούν για να ταυτοποιηθούν τα δεδομένα σε σχέση με έναν πίνακα που περιλαμβάνει ηλικία, γένος, Τ.Κ όπως και το όνομα. Εάν υπάρχει και το ανάλογο ενδιαφέρον από τον επιτιθέμενο σε συγκεκριμένο άτομο που ήδη γνωρίζει... και άρα ορισμένα από τα πεδία, γίνεται ακόμα πιο εύκολο.

Το ψευδοαναγνωριστικό για τον πίνακα P από το Πίνακα 1a μπορεί να έχει τον εξής τύπο $Q_P = \{age, gender, zip\}$.

Το ελάχιστο σύνολο από γνωρίσματα $Q = Q_1, \dots, Q_d$ με το οποίο ένας πίνακας T μπορεί να συζευχθεί με κάποιες εξωτερικές πληροφορίες για να αναγνωριστούν ατομικές εγγραφές ονομάζεται ψευδοαναγνωριστικό σύνολο.

	Quasi – Identifier - Ψευδοαναγνωριστικά			Ευαίσθητα Δεδομένα
	Ηλικία	Γένος	Τ.Κ	Έσοδα
1	<45	Αρσενικό	81****	40,000
2	<45	Αρσενικό	81****	40,000
3	<45	Αρσενικό	81****	300,000
4	<45	Αρσενικό	81****	1,000,000
5	≥45	Αρσενικό	7*****	20,000
6	≥45	Αρσενικό	7*****	25,000
7	≥45	Αρσενικό	7*****	100,000
8	≥45	Αρσενικό	7*****	1,500,000
9	*	Θηλυκό	83****	25,000
10	*	Θηλυκό	83****	30,000
11	*	Θηλυκό	83****	30,000
12	*	Θηλυκό	83****	60,000
13	*	Θηλυκό	83****	70,000

Πίνακας 1b: Γενικευμένος πίνακας G1 όπου βασίζεται στον P Εφαρμογή k-ανωνυμίας για k=4

Για να πετύχουμε k -ανωνυμία στον P , τα γνωρίσματα που είναι στα πεδία των ψευδοαναγνωριστικών θα πρέπει να γενικευτούν. Ένας $*$ στο T.K. θα μπορούσε να σημαίνει οποιοδήποτε στοιχείο. Στην στήλη της ηλικίας το <45 σημαίνει ότι η ηλικία είναι κάτω των 45, ενώ \geq σημαίνει ότι η ηλικία μπορεί να είναι ίση ή ανώτερη των 45 και $*$ μπορεί να σημαίνει οποιοδήποτε αριθμό. Ο G1 όπως φαίνεται και από τον πίνακα είναι μια γενικευμένη έκδοση του P που εκπληρώνει της k -ανωνυμία με $k=4$. Κάθε πεδίο με τον ίδιο quasi-identifier περιλαμβάνει τουλάχιστον 4 εισόδους. Εάν ο επιτιθέμενος ενδιαφέρεται για τα έσοδα ενός ατόμου με ένα συγκεκριμένο ψευδοαναγνωριστικό, υπάρχουν τουλάχιστον $k-1=3$ άλλοι άνθρωποι με το ίδιο ψευδοαναγνωριστικό στον πίνακα.

Εάν υπάρχουν το λιγότερο k tuples με τα ίδια ψευδοαναγνωριστικά, δεν είναι εφικτό να προσδιορίσεις ένα μόνο που βασίζεται σε αυτά. Υπάρχουν $k-1$ tuples με το ίδιο ψευδοαναγνωριστικό που δεν είναι διακριτά από το tuple που ψάχνει ο επιτιθέμενος.

Κεφάλαιο 3. Μηχανισμοί που υποστηρίζουν την k -ανωνυμία

Μετα την πρόταση της k -ανωνυμίας, υλοποιήθηκαν διάφορες προσπάθειες σχεδίασης ενός καλού αλγόριθμου που θα μετέτρεπε μια βάση σε μορφή ανάλογη της υλοποίησης του ορισμού. Οι δύο κύριες τεχνικές που χρησιμοποιήθηκαν για την επιβολή k -ανωνυμίας στα δημοσιευμένα δεδομένα είναι η Γενίκευση και η Απόκρυψη (suppression).

3.1. Γενίκευση (Generalization):

Η διαδικασία κατά την οποία οι τιμές των ψευδοαναγνωριστικών (Quasi Identifiers) αντικαθίστανται με γενικότερες ή με βάση συγκεκριμένες ιεραρχίες γενίκευσης. Στόχος της γενίκευσης είναι η διατήρηση μέρους της πληροφορίας της αρχικής τιμής χωρίς αυτή να αλλάζει και να αλλοιώνεται πλήρως.

Επιπρόσθετα, κάθε τιμή μπορεί να γενικευθεί σε πολλά στάδια και σε πιο γενικές σημασιολογικές τιμές. Τα διαφορετικά επίπεδα γενίκευσης του πεδίου τιμών ενός γνωρίσματος, στα οποία οδηγούνται οι αρχικές τιμές με κάθε γενίκευση συνήθως αποτυπώνονται με τη μορφή δένδρου, το οποίο ονομάζεται ιεραρχία γενίκευσης του πεδίου τιμών, (Domain Generalization Hierarchy).

Στην περίπτωση του Πίνακα 1, οι τιμές του φύλου, της ηλικίας και του T.K. από τον (a) αντικαθίστανται από μια γενικευμένη έκδοση (b). Η γενίκευση μπορεί να εφαρμοστεί σε επίπεδα από ένα και μόνο κελί σε ολόκληρο tuple (γραμμή) γνωρισμάτων για την επίτευξη k -ανωνυμίας.

3.2. Απόκρυψη (Suppression),

Συνίσταται στην αφαίρεση ευαίσθητων γνωρισμάτων για την μείωση της ποσότητας της γενίκευσης όταν επιτυγχάνεται k -ανωνυμία. Όπως και η γενίκευση, η απόκρυψη μπορεί να εφαρμοστεί τόσο σε κελιά όσο σε ολόκληρα γνωρίσματα. Ο συνδυασμός της Γενίκευσης και της Απόκρυψης χρησιμοποιήθηκε για την δημιουργία διαφορετικών αλγόριθμων με σκοπό την ικανοποίηση της k -ανωνυμίας.

Το συμβατικό πλαίσιο – η λογική, ενός τέτοιου αλγόριθμου πάντα ξεκινά με την απόκρυψη πολλών ευαίσθητων γνωρισμάτων και εν συνέχεια χωρίζει tuples των εναπομείνοντων γνωρισμάτων σε ομάδες καθώς υποκαθιστά τις ακριβείς τιμές των QID με γενικευμένες μορφές τους για κάθε ομάδα, όπου ονομάζονται κλάσεις ισοδυναμίας. Αυτή η γενίκευση είναι ομοιογενής γενίκευση και έχει χρησιμοποιηθεί για την εφαρμογή της k -ανωνυμίας στις αναφορές των Iwuchukwu and Naughton (2007) [4] Ghinita et al. (2007), και LeFevre et al. (2008). [5]

Η ιδιότητα μιας ομοιογενής γενίκευσης ορίζει ότι, αν μια αρχική εγγραφή t_i ταιριάζει με μια δημοσιευμένη εγγραφή t_j' της οποίας η αντίστοιχη εγγραφή είναι t_j , τότε η t_j θα ταιριάζει με την t_i' . Αυτή η ιδιότητα ονομάζεται αμοιβαιότητα.

Το πιο σημαντικό στοιχείο της ομοιογενούς γενίκευσης είναι ο τρόπος διαίρεσης των κλάσεων ισοδυναμίας. Η “στρατηγική” γενίκευσης θα επηρεάζει άμεσα την χρησιμότητα των δημοσιευμένων δεδομένων.

Υπάρχουν δύο τύποι μοντέλων γενίκευσης: Ολική και Τοπική Γενίκευση. Ακολουθεί η ανάλυση τους.

3.2.1. Ολική Γενίκευση - Global Recoding

(πλήρη ανωνυμοποίηση τομέα – domain) [7][8]

Στην Ολική Γενίκευση, σε μία στήλη εφαρμόζεται η ίδια στατηγική γενίκευσης στην ίδια τιμή. Έτσι αν δύο tuples στις αρχικές καταγραφές έχουν ίδιες QID τιμές, τότε θα πρέπει να έχουν την ίδια δημοσιευμένη τιμή.

Εν προκειμένω, τα σύνολα των δεδομένων που ανωνυμοποιήθηκαν με την εν λόγω μέθοδο έχουν όλα τα γνωρίσματα ισάξια γενικευμένα ή κρυμμένα σε όλες τις καταχωρίσεις. Εάν διαφορετικά tuples είχαν την ίδια τιμή γνωρίσματος, αυτή η τιμή θα αντιστοιχούσε στην ίδια γενικευμένη τιμή. (Πίνακας 2a)

Η Ολική Γενίκευση μπορεί να χαρακτηριστεί ως μονοδιάστατη γενίκευση ή πολυδιάστατη γενίκευση. Η μονοδιάστατη ορίζεται από την συνάρτηση $f_i : D_{x_i} \rightarrow D$ για κάθε γνώρισμα x_i του ψευδοαναγνωριστικού. Η τιμή της ανωνυμοποίησης μπορεί να ληφθεί εφαρμόζοντας την συνάρτηση f_i στις τιμές των ψευδοαναγνωριστικών x_i σε κάθε tuple του αρχικού συνόλου δεδομένων. Ενώ, η πολυδιάστατη γενίκευση ορίζεται από μια μόνο συνάρτηση

$f : D_{x_1} * \dots * D_{x_n} \rightarrow D$, η τιμή της ανωνυμοποίησης λαμβάνεται εφαρμόζοντας την συνάρτηση f στο διάνυσμα κάθε tuple του αρχικού συνόλου δεδομένων.

3.2.2 Τοπική Γενίκευση - Local Recoding

Στην περίπτωση της Τοπικής Γενίκευσης 2 tuples με ίδιες QID τιμές μπορεί να έχουν διαφορετικές γενικευμένες τιμές.

Σε αντίθεση με την Ολική Γενίκευση, η τοπική αποκρύπτει τα γνωρίσματα ανα στοιχείο κατάχωρησης. Αποτέλεσμα αυτού ένα σύνολο δεδομένων με λιγότερα κρυφά δεδομένα σε σχέση με της Ολικής Γενίκευσης. Ο χώρος των δεδομένων χωρίζεται σε διαφορετικές περιοχές και όλες οι εγγραφές της ίδιας περιοχής αντιστοιχίζονται στην ίδια γενικευμένη ομάδα. Από την στιγμή που τα όρια ποιότητας είναι αρκετά μεγάλα δεν είναι εύκολη η απόδειξη της βέλτιστης λύσης που δόθηκε στην γενίκευση. [9]

Η ολική γενίκευση επιτυγχάνει την ανωνυμία χαρτογραφώντας/αντιστοιχίζοντας τους τομείς των γνωρισμάτων των ψευδοαναγνωριστικών σε γενικές ή τροποποιημένες τιμές. Ενώ τα μοντέλα τοπικής γενίκευσης αποτυπώνουν (μη διακριτά) μεμονωμένα στοιχεία δεδομένων σε γενικευμένες τιμές (Πίνακας 2b). Το πλεονέκτημα της ολικής γενίκευσης είναι ότι ο ανωνυμοποιημένος πίνακας θα έχει ένα ομοιογενές σύνολο τιμών ενώ το μειονέκτημα του είναι ότι υπερ-γενικεύει τον αρχικό πίνακα και με αποτέλεσμα να χάνεται μεγάλη πληροφορία κατά την διαδικασία. Σε αντίθεση η τοπική γενίκευση δίνει πολύ καλύτερα αποτελέσματα σε πραγματικό σενάριο και έχει πολύ μεγαλύτερη χρησιμότητα σε σχέση με την ολική γενίκευση

Πίνακας 2: Αρχικός Πίνακας στοιχείων δεδομένων

Ηλικία	Φύλο	Τ.Κ.	Αρρώστια
23	Άντρας	34756	Βρογχίτιδα
29	Άντρας	34201	Καρκίνος
35	Γυναίκα	36020	Γρίπη
39	Άντρας	37013	Ιλαρά

Πίνακας 2a: Με χρήση ολικής γενίκευσης

Ηλικία	Φύλο	Τ.Κ.	Αρρώστια
[21 - 30]	Άνθρωποι	[34201 - 37013]	Βρογχίτιδα
[21 - 30]	Άνθρωποι	[34201 - 37013]	Καρκίνος
[31 - 40]	Άνθρωποι	[34201 - 37013]	Γρίπη
[31 - 40]	Άνθρωποι	[34201 - 37013]	Ιλαρά

Πίνακας 2b: Πίνακας με χρήση τοπικής γενίκευσης

Ηλικία	Φύλο	Τ.Κ.	Αρρώστια
[21 - 30]	A	[34201 - 34756]	Βρογχίτιδα
[21 - 30]	*	[34201 - 34756]	Καρκίνος
[31 - 40]	*	[36020 - 37013]	Γρίπη
[31 - 40]	A	[36020 - 37013]	Ιλαρά

*Οι αλγόριθμοι που χρησιμοποιούν Τοπική Γενίκευση μπορεί να εγγυηθούν περισσότερη ανωνυμία σε συσχετισμένες περιπτώσεις.(Ninghui Li and Su, 2011).

Κεφάλαιο 4. Αλγόριθμοι k-Ανωνυμίας

Στην ολική γενίκευση οι τιμές των αρχικών δεδομένων γενικεύονται στο επίπεδο του τομέα. Υπάρχουν αρκετές εργασίες που βασίζονται στην μέθοδο της ολικής γενίκευσης, κάποιες από αυτές απόδωσαν τους κάτωθι καρπούς – αλγόριθμους που θα εξετάσουμε.

4.1 Αλγόριθμος Incognito

Ο Αλγόριθμος **Incognito** που προτάθηκε στη [7] χρησιμοποιεί δυναμικό προγραμματισμό και φαίνεται να υπερέρχει από προηγούμενους αλγόριθμους σε εφαρμογή του πάνω σε σε σενάριο με 2 βάσεις δεδομένων βασισμένων σε πραγματικά στοιχεία. Η κύρια ιδέα του Incognito είναι ότι κάθε υποσύνολο tuple των QID με k-ανωνυμία, θα πρέπει να εφαρμόζει και αυτό k-ανωνυμία. (subset property)

Άρα, αν ένας πίνακας T είναι k-ανώνυμος ως προς ένα σύνολο γνωρισμάτων Q της βάσης δεδομένων, τότε είναι k-ανώνυμος και ως προς οποιοδήποτε υποσύνολο γνωρισμάτων $P \subseteq Q$.

Ο Incognito παράγει ελάχιστη γενίκευση πλήρους τομέα. Στην ολική γενίκευση, αν ένας τομέας χαμηλότερου επιπέδου πρέπει να γενικευθεί σε τομέα υψηλότερου επιπέδου, όλες οι τιμές του χαμηλότερου γενικεύονται σε υψηλότερου. Αυτό μπορεί προκαλέσει την υπεργενίκευση του πίνακα, που έχει σαν αποτέλεσμα μεγάλη απώλεια πληροφορίας.

Στην πολυδιάσταση και τοπική γενίκευση, η γενίκευση γίνεται πάντα στο επίπεδο των στοιχείων, γεγονός που δεν προκαλεί υπεργενίκευση, οδηγώντας σε πιο ευέλικτη γενίκευση με μικρότερη απώλεια πληροφορίας.

4.2 Αλγόριθμος Mondrian

Ο Αλγόριθμος **Mondrian** [11] χρησιμοποιεί την στατηγική της πολυδιάστατης Ολικής Γενίκευσης. Πρόκειται για έναν πολύ γρήγορο και κλιμακωτό αλγόριθμο εύρεσης, ο οποίος παράγει καλύτερα αποτελέσματα.

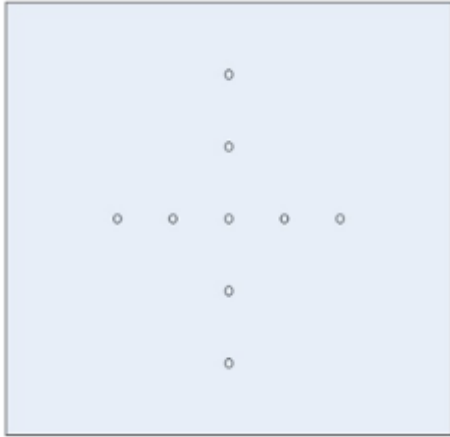
Στον Mondrian, κάθε γνώρισμα στο σύνολο των δεδομένων αντιπροσωπεύει μια διάσταση και κάθε εγγραφή αντιπροσωπεύει ένα σημείο στον χώρο[12]. Αντί, να διαιρεί τις εγγραφές, ο Mondrian διαιρεί τον χώρο σε πολλές περιοχές και σε κάθε περιοχή υπάρχουν το λιγότερα k σημεία. Χρησιμοποιεί μεθόδους αυστηρής και χαλαρής κατάτμισης, οι οποίες οδηγούν σε καλύτερα αποτελέσματα χρησιμότητας των δεδομένων.

Η τεχνική της κατάτμισης αντιστοιχίζει κάθε tuple του συνόλου δεδομένων σε ένα πολυδιάστατο χώρο. Στην συνέχεια, η γενίκευση του συνόλου δεδομένων ισούται με την κατάτμιση του αντίστοιχου πολυδιάστατου χώρου. Ένα κομμάτι του πολυδιάστατου χώρου αντιστοιχεί σε μοναδικό αποτέλεσμα ανωνυμίας. Εάν τα κομμάτια δεν διασταυρώνονται ή επικαλύπτονται το ένα με το άλλο τότε είναι γνωστό ως Αυστηρή Κατάτμιση. Αντιθέτως, αν αλληλεπικαλύπτονται μεταξύ τους τότε είναι γνωστό ως Χαλαρή Κατάτμιση.[12]

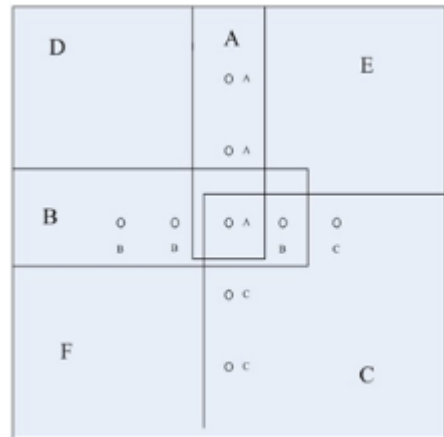
Η χαλαρή κατάτμιση θεωρείται καλύτερη τεχνική σε σχέση με την αυστηρή, και παραθέτουμε μια περίπτωση που το αποδεικνύει.

Παράδειγμα: Αν υποθέσουμε ότι θα πρέπει να εφαρμόσουμε αυστηρή κατάτμιση στον πίνακα του σχεδίου 2, σε δύο το λιγότερο περιοχές, τότε θα πρέπει να υπάρχει μια περιοχή που να περιλαμβάνει όχι περισσότερα απο δύο tuples.

Επομένως, αν απαιτηθεί 3-ανωνυμία ή περισσότερες, τότε δεν θα υπάρξει κανένας “αυστηρής κατάτμισης” αλγόριθμος για να υποστηρίξει το σκοπό μας



Σχέδ 2: Δεδομένα σε 2-διάστατο επίπεδο



Σχέδ 3: Χαλαρή Κατάτμιση

Αλλά μπορούμε να χρησιμοποιήσουμε ένα αλγόριθμο χαλαρής κατάτμισης για το διαμοιρασμό του πίνακα σε 6 περιοχές, A,B,C,D,E,F όπως υποδεικνύει και το σχέδιο 3. Οι περιοχές A,B, και C περιλαμβάνουν 3 tuples ενώ οι περιοχές D,E και F δεν περιλαμβάνουν tuples. Κάθε tuple στην διασταύρωση των περιοχών ανήκει σε μόνο μία περιοχή.[12][13]

Άρα, ο αλγόριθμος Mondrian επιστρέφει την βέλτιστη πολυδιάστατη διαίρεση σε κάθε περιοχή της οποίας ανήκουν περισσότερες από k-εγγραφές και συνεπώς ικανοποιείται η k-ανωνυμία.

4.3 Μη-ομοιογενής γενίκευση

Οι περισσότεροι αλγόριθμοι που χρησιμοποιούν την γενίκευση όπως και οι προαναφερόμενοι μοιράζονται μια κοινή δομή: στη πρώτη κατάτμιση τα tuples σε ομάδες, στην συνέχεια αντιστοιχίζονται τα ίδια γενικευμένα ψευτοαναγνωριστικά (QID) στα tuples της ίδιας ομάδας. Η ομάδα των tuples με τα ίδια QID λέγεται κατηγορία ισοδυναμίας. Μια τέτοια προσέγγιση, την οποία αποκαλούμε *ομοιογενή γενίκευση*.

Υπάρχει μια ακόμη μέθοδος που ονομάζεται μη-ομοιογενής γενίκευση. Η διαφορά της σε σχέση με την ομοιογενή είναι ότι μπορεί να επιτύχει κάποιος υψηλότερη χρηστικότητα των δεδομένων από τα δημοσιευμένα δεδομένα.

Στον Πίνακα 3, ο (b) είναι τα δημοσιευμένα δεδομένα που προκύπτουν από τον (a) χρησιμοποιώντας ομοιογενή γενίκευση, και είναι ξεκάθαρο ότι (t_1', t_2', t_5') είναι τάξεις ισοδυναμίας και (t_3', t_4') είναι άλλες. Σε κλάση ισοδυναμίας, όλες οι γενικευμένες QID τιμές είναι οι ίδιες. Ωστόσο, σε μη ομοιογενή γενικευμένο πίνακα (c), t_1', t_2' και t_5' έχουν διαφορετικές γενικευμένες QID τιμές. Ενώ και οι δύο πίνακες (b) (c) έχουν 2-ανωνυμία, ο (c) προσφέρει υψηλότερη χρηστικότητα από την στιγμή που τα γενικευμένα QID εύρη στον (c) είναι μικρότερα ή ίσοδύναμα με τα αντίστοιχα του πίνακα (b). Αυτό υποδεικνύει ότι αν κάποιος χρησιμοποιεί μη-ομοιογενή γενίκευση μπορεί να επιτύχει υψηλότερη χρηστικότητα των δεδομένων από τα δημοσιευμένα δεδομένα.

Πίνακας 3

Παράδειγμα k -ανωνυμίας ($k=2$) από ομοιογενή σε μη-ομοιογενή γενίκευση

(a) Αρχικά Δεδομένα			
Tuple ID	Γένος	Ηλικία	T.K
t_1	Θηλυκό	17	21103
t_2	Αρσενικό	29	21110
t_3	Αρσενικό	27	21210
t_4	Αρσενικό	15	21102
t_5	Θηλυκό	22	21109

(b) Ανταλλαγή δεδομένων που δημιουργούνται με ομοιογενή γενίκευση			
Tuple ID	Γένος	Ηλικία	T.K
t_1'	Θ ή A	17-29	2110*
t_2'	Θ ή A	17-29	211*
t_3'	Αρσενικό	15-27	212*
t_4'	Αρσενικό	15-27	211*
t_5'	Θ ή A	17-29	211*

(c) Κοινή χρήση δεδομένων που δημιουργούνται από μη-ομοιόμορφη γενίκευση			
Tuple ID	Γένος	Ηλικία	T.K
t_1'	Θηλυκό	17-29	2110*
t_2'	Αρσενικό	17-29	211*
t_3'	Αρσενικό	15-27	212*
t_4'	Αρσενικό	15-27	211*
t_5'	Θ ή A	17-29	211*

Κεφάλαιο 5. Επιθέσεις στην k-ανωνυμία

Η έκδοση πολλων διαφορετικών συνολων δεδομένων που βασίζεται στην ίδια ομάδα κάτοχων δεδομένων δημιουργεί περισσότερους παράγοντες επίθεσης. 3 τέτοιες επιθέσεις αναφέρονται απο κάτω. Όταν δημοσιοποιούνται αντίστοιχα των ακόλουθων δεδομένα καποιες απο τις ακόλουθες τακτικές μπορούν να αποτρέψουν τέτοιες επιθέσεις.

i) Επίθεση μη ομαδοποιημένης στοίχισης

Όταν δύο πίνακες δημοσιοποιούνται μπορεί να χρησιμοποιηθούν για να συνδέσουν σύνολα δεδομένων, εαν βασίζονται στον ίδιο αρχικό πίνακα και η θέση των tuples είναι η ίδια όπως ήταν στον προηγούμενο. Όπως το μοντέλο της k-ανωνυμίας χρησιμοποιεί το ίδιο “σχεσιακό” μοντέλο, θεωρητικά δεν θα υπάρχει καποια προδιεγεγραμμένη σειρά των tuples. Σχέσεις είναι μια ομάδα απο tuples και στις ομάδες δεν υπάρχει καμία σειρά. Όταν δημοσιοποιούνται αληθινα δεδομένα, υπάρχει μεγάλη πιθανότητα οι tuples να ταξινομούνται βάση κάποιου γνωρίσματος.

Ο σκοπός της γενίκευσης είναι να ταξινομή τα δεδομένα κατα αύξουσα ή φθίνουσα σειρά βάση ενός ή περισσοτέρων γνωρισμάτων των ψευδοαναγνωριστικών.

	Quasi – Identifier - Ψευδοαναγνωριστικά			Ευαίσθητα Δεδομένα
	Ηλικία	Γένος	Τ.Κ	Έσοδα
1	27	*	*****	40,000
2	40	*	*****	40,000
3	35	*	*****	300,000
4	40	*	*****	1,000,000
5	60	*	*****	20,000
6	46	*	*****	25,000
7	60	*	*****	100,000
8	48	*	*****	1,500,000
9	27	*	*****	25,000
10	34	*	*****	30,000
11	48	*	*****	30,000
12	27	*	*****	60,000
13	46	*	*****	70,000

Πίνακας 4:Γενικευμένος Πίνακας G2

Ο G2 στον 4ο πίνακα βασίζεται στον P(Πίνακας 1a). Η σειρά είναι η ίδια με τον G1(Πίνακας 1b). Ο G1 ικανοποιεί την k-ανωνυμία με k=4, ο G2 με k=2. Εάν και οι δύο δημοσιοποιηθούν, ο επιτιθέμενος θα μπορέσει να συνδέσει τα tuples βασισμένα στην θέση τους. Έτσι θα μπορέσει να μάθει περισσότερα για την ηλικία και το φύλο της κάθε tuple, αρκετή γνώση για αναγνώριση ενός ανθρώπου.

Αυτή η επίθεση μπορεί να αντιμετωπιστεί εύκολα με την τυχαιοποίηση της σειράς των tuples, όταν απελευθερώσουμε τα δεδομένα.

ii) Επίθεση συμπληρωματικής δημοσιοποίησης

Έαν ένας πίνακας που περιλαμβάνει ένα υποσύνολο προηγούμενου δημοσιοποιημένου πίνακα δημοσιοποιηθεί εκ νέου, μια τέτοιου είδους επίθεση μπορεί να είναι εφικτή. Γνωρίσματα τα οποία δεν είναι μέρη των ψευδοαναγνωριστικών του πρώτου πίνακα μπορούν να χρησιμοποιηθούν για να διασυνδεθούν με τους νέους πίνακες.

	Quasi – Identifier - Ψευδοαναγνωριστικά			Ευαίσθητα Δεδομένα
	Ηλικία	Γένος	Τ.Κ	Έσοδα
1	46	*	*****	25,000
2	27	*	*****	40,000
3	48	*	*****	30,000
4	40	*	*****	40,000
5	40	*	*****	1,000,000
6	27	*	*****	60,000
7	46	*	*****	70,000
8	60	*	*****	20,000
9	60	*	*****	100,000
10	27	*	*****	25,000
11	34	*	*****	30,000
12	35	*	*****	300,000
13	48	*	*****	1,500,000

Πίνακας 5: Γενικευμένος Πίνακας G3

Ο G3 (Πίνακας 5) είναι μια ανώνυμη μορφή του P. Έαν δημοσιοποιηθεί μετά τον G1, κάποια tuples μπορούν να συνδεθούν, ακόμα και αν οι θέσεις τους είναι τυχαίες. Π.χ υπάρχει μόνο ένα άτομο με το εισόδημα του 1.500.000. Αυτό είναι αρκετό για να αντιστοιχηθεί με τα αντιστοιχα tuples και να αποκτήσουμε πληροφορία για τον Qr. Η συνδέουσα τιμή για αυτό το tuple είναι [48, Αρσενικό, 7****, 1,500,000]. Για όλες τις tuples με μοναδικό συνδυασμό των τιμών για {Q G1 ∪ Έσοδα}

Για να αποτρέψουμε μια τέτοια επίθεση, όλα τα γνωρίσματα της πρώτης έκδοσης/δημοσιοποίησης θα πρέπει να χειριστούν σαν ψευδοαναγνωριστικά στον δεύτερο πίνακα. Σε αυτή την περίπτωση ο πίνακας δεν πρέπει να δημοσιοποιηθεί με αυτό τον τρόπο καθότι δεν καλύπτει k-ανωνυμία.

Επίσης αν χρησιμοποιηθεί ο G1 σαν βάση επίσης δεν θα έχουμε συντελεστή επίθεσης, καθότι αν έχουμε το ίδιο ανώνυμο σχεδιάγραμμα, δεν μπορεί να ληφθεί επιπρόσθετη πληροφορία.

iv) Χρονική Επίθεση:

Καθώς τα δεδομένα αυξάνονται και μαζεύονται τακτικά, τα σύνολα τους ενδέχεται να αυξηθούν χρονικά. Αλλαγές ή αφαιρέσεις των tuples είναι εξίσου πιθανές. Γι αυτό ακριβώς τον λόγο, δημοσιοποιούνται συχνά επιπρόσθετα δεδομένα. Αυτές οι δημοσιοποιήσεις μπορεί να είναι ευπαθής στην σύνδεση των πινάκων.

π.χ. Ας υποθέσουμε ότι μια έρευνα μάζεψε δεδομένα όπως ο P, και δημοσιοποίησε G1. Αργότερα επιπρόσθετοι tuples μαζεύτηκαν και ενημέρωσα τον P, οπότε ο P_{t1} έγινε P ∪ {51, Θηλυκό, 83581, 28,000}, {51, Αρσενικό, 81019, 32,000}}. Βάση του P_{t1} ένας πίνακας γενίκευσης G_{t1} ως G3 {{51, *, *****, 28,000}, {51, *, *****, 32,000}} εκδίδεται. Όπως φάνηκε και στο παράδειγμα 5 οι πίνακες G1 & G3 μπορούν να διασυνδεθούν.

Για να αποφύγουμε τέτοιες επιθέσεις, μπορούμε να προβούμε ως εξής:

$G1 \cup (P_{t1} - P)$. Δηλαδή $G1 \cup \{[* , \Theta\eta\lambda\upsilon\kappa\acute{o}, 83^{***}, 28,000], [\geq 45, \text{Αρσενικό}, 81^{***}, 32,000]\}$

Τα προαναφερόμενα ήδη επιθέσεων, όπως παρατηρήσαμε, αντιμετωπίζονται χρησιμοποιώντας αλγορίθμους είτε/και επιπρόσθετες προσωπικές ενέργειες - μεθόδους προσέγγισης χωρίς να αποτελούν πανάκεια - καθότι υπάρχουν νέες επιθέσεις και νέες αντιμετώπισεις που μπορούν να προσφέρουν ισχυρότερη ανωνυμία λειτουργώντας εναλλακτικά είτε συμπληρωματικά. Παρακάτω θα εξεταστούν αναλυτικότερα οι δύο βασικότερες: η I- ποικιλομορφία και η t – εγγύτητα.

Κεφάλαιο 6. Σύνθετες επιθέσεις στην k-ανωνυμία και τρόποι επίλυσης

6.1 Εισαγωγή

Η προστασία που η k-ανωνυμία προσφέρει είναι απλή και εύκολη στην κατανόηση. Εάν ένας πίνακας ικανοποιεί k-ανωνυμία για κάποια τιμή k, τότε ο οποιοσδήποτε γνωρίζει μόνο τις τιμές του ψευδο-αναγνωριστικού ενός συγκεκριμένου ατόμου δεν μπορεί να προσδιορίσει – αναγνωρίσει την καταγραφή που αρμόζει σε αυτό το συγκεκριμένο άτομο με σιγουριά μεγαλύτερη του $1/k$.

Ενώ η k-ανωνυμία προστατεύει την αποκάλυψη ταυτότητας, δεν προσφέρει επαρκής προστασία έναντι της αποκάλυψης του γνωρίσματος. Αυτό αναγνωρίστηκε από αρκετούς συγγραφείς. Δυο είναι οι επιθέσεις που βρέθηκαν: Η επίθεση της ομοιογένειας και η επίθεση με πρότερη γνώση. Παρακάτω θα εξετάσουμε 2 παραδείγματα για την καλύτερη αντίληψη των ευπαθειών.

6.2. Επίθεση της Ομοιογένειας

Η επίθεση της Ομοιογένειας αξιοποιεί την περίπτωση όπου όλες οι τιμές που μπορεί να πάρει μια ευαίσθητη τιμή μέσα σε ένα σύνολο k εγγραφών, είναι ίδιες.

Για παράδειγμα, ας υποθέσουμε ότι έχουμε μια ομάδα k διαφορετικών αρχείων και όλα μοιράζονται ένα συγκεκριμένο ψευδοαναγνωριστικό. Ένας επιτιθέμενος δεν μπορεί να εντοπίσει το άτομο βάσει των ψευδοαναγνωριστικών. Αλλά αν ενδιαφέρεται για τα ευαίσθητα γνωρίσματα και όλες οι ομάδες έχουν την ίδια τιμή τότε τα δεδομένα έχουν διαρρεύσει. Παραθέτουμε και ένα παράδειγμα για την καλύτερη κατανόηση:

Το σενάριο μας περιλαμβάνει την *Αλίκη* και τον *Πέτρο* ως “θύτης” και “θύμα”. Η Αλίκη - βάση του ρόλου της - , παρακολουθεί την προσωπική ζωή του Πέτρου με απώτερο στόχο να γνωρίσει τις αδυναμίες του. Ανακαλύπτει ότι ο Πέτρος επισκέπτεται τακτικά το νοσοκομείο. Πρέπει να μάθει αν υποφέρει από κάποια αρρώστια. Εισέρχεται στην προσοχή της ένας 4ης – Ανωνυμίας πίνακας ενός δημοσιευμένου εγγράφου των ασθενών του νοσοκομείου που επισκέπτεται ο Πέτρος. (Πίνακας 7). Έτσι, γνωρίζει ότι σε κάποια από τις εγγραφές αντιστοιχούν τα δεδομένα του Πέτρου. Η Αλίκη γνωρίζει ότι ο Πέτρος είναι κάπου στα 30 με 35, Αμερικανός και τον T.K (13053) του. Έτσι λοιπόν γνωρίζει ότι οι τελευταίες εγγραφές (9-12) ανταποκρίνονται στα γνωστά της δεδομένα. Βγαίνει λοιπόν το πόρισμα ότι ο Πέτρος έχει Καρκίνο.[16]

1η Παρατήρηση Η k-ανωνυμία μπορεί να δημιουργήσει ομάδες που μπορούν να διαρρεύσουν πληροφορίες χάριν της ελλείψεως διαφορετικότητας στα ευαίσθητα γνωρίσματα.

Τέτοιες περιπτώσεις δεν είναι ασυνήθιστες, και πόσο μάλλον όταν εξετάζονται δεδομένα 60.000 διαφορετικών tuples, όπου τα ευαίσθητα γνωρίσματα θα μπορούν να πάρουν μόνο 3 διακριτές τιμές. Οπότε, προτείνεται επιπρόσθετα να εφαρμόζεται “διαφορετικότητα” σε όλα τα tuples που μοιράζονται τις ίδιες τιμές των ψευδο-αναγνωριστικών τους, να έχουν διαφορετικές τιμές στα ευαίσθητα γνωρίσματα.

	Μη - Ευαίσθητα			Ευαίσθητα
	Τ.Κ.	Ηλικία	Εθνικότητα	Κατάσταση
1	13053	28	Ρωσική	Καρδιακή Ασθένεια
2	13068	29	Αμερικάνικη	Καρδιακή Ασθένεια
3	13068	21	Γαλική	Ύωση
4	13053	23	Αμερικάνικη	Ύωση
5	14853	50	Ινδική	Καρκίνος
6	14853	55	Ρωσική	Καρδιακή Ασθένεια
7	14850	47	Αμερικάνικη	Ύωση
8	14850	49	Αμερικάνικη	Ύωση
9	13053	31	Αμερικάνικη	Καρκίνος
10	13053	37	Ινδική	Καρκίνος
11	13068	36	Ιαπωνική	Καρκίνος
12	13068	35	Αμερικάνικη	Καρκίνος

Πίνακας 6. Μικροδεδομένα Ασθενών

	Μη - Ευαίσθητα			Ευαίσθητα
	Τ.Κ.	Ηλικία	Εθνικότητα	Κατάσταση
1	130**	<30	*	Καρδιακή Ασθένεια
2	130**	<30	*	Καρδιακή Ασθένεια
3	130**	<30	*	Ύωση
4	130**	<30	*	Ύωση
5	1485*	≥40	*	Καρκίνος
6	1485*	≥40	*	Καρδιακή Ασθένεια
7	1485*	≥40	*	Ύωση
8	1485*	≥40	*	Ύωση
9	130**	3*	*	Καρκίνος
10	130**	3*	*	Καρκίνος
11	130**	3*	*	Καρκίνος
12	130**	3*	*	Καρκίνος

Πίνακας 7. 4-ανώνυμιας μικροδεδομένα ασθενών

6.3. Επίθεση με πρότερη γνώση

Αυτή η επίθεση αξιοποιεί μια συσχέτιση μεταξύ ενός ή περισσότερων γνωρισμάτων ψευδοαναγνωριστικών με το ευαίσθητο γνώρισμα, με σκοπό την μείωση των πιθανών τιμών του εύρους του ευαίσθητο γνώρισματος [15]

Σε συνέχεια του προηγούμενου σεναρίου μας, η Αλίκη μετά την επιτυχή επίθεση ομοιογένειας που εφάρμοσε στον (Πίνακα 7), παίρνει θάρρος και δοκιμάζει τις δυνατότητες της στην εύρεση νέου στόχου. Το θύμα της τώρα η Μαρία. Τα γνωστά γι αυτή στοιχεία είναι ότι η Μαρία είναι γειτόνισσα της, άρα έχει το ίδιο T.K. 13053 είναι 28 χρονών και νοσηλεύεται τελευταία στο ίδιο νοσοκομείο. Η Μαρία είναι Γαλλίδα. Παράλληλα, η Αλίκη γνωρίζει ότι οι Γάλλοι έχουν τις χαμηλότερες εμφανίσεις εμφραγμάτων παγκόσμια και κατ' επέκταση τις λιγότερες πιθανότητες να εμφανίσουν καρδιαγγειακά νοσήματα. Άρα συμπεραίνει ότι η Μαρία έχει μια απλή ίωση.

2η Παρατήρηση Η k-ανωνυμία δεν μπορεί να προστατεύσει από επιθέσεις που βασίζονται στην πρότερη γνώση.

6.4 I – Ποικιλομορφία (I-diversity)

Ένας νέος ορισμός προστασίας της ιδιωτικότητας που ήρθε να επιλύσει το πρόβλημα της αποκάλυψης ευαίσθητων γνωρισμάτων που συνάδει με την k-ανωνυμία είναι η I-ποικιλομορφία, και όπως προτάσει και το όνομα, διασφαλίζει την διαφορετικότητα – ποικιλομορφία των τιμών των ευαίσθητων γνωρισμάτων σε κάθε κλάση ισοδυναμίας, διατηρώντας το ελάχιστο μέγεθος του k-συνόλου.

Η βασικός τρόπος λειτουργίας της απαιτεί κάθε ομάδα ψευδοαναγνωριστικών, να έχει τουλάχιστον I - “καλά εκπροσωπούμενες” διαφορετικές τιμές, οι οποίες θα χρησιμοποιηθούν για να κάνουν μια tuple, ασαφή. Για να προσδιορίσουμε του πόσο καλά οι τιμές των γνωρισμάτων θα εκπροσωπευτούν, μπορούμε να εφαρμόσουμε τα ακόλουθα πιθανά μοντέλα.

i) **Διακριτή I - ποικιλομορφία.** Η απλούστερη κατανόηση του όρου “καλά εκπροσωπούμενοι” θα ήταν η εξασφάλιση ότι υπάρχουν το λιγότερο I διακριτές τιμές για τα ευαίσθητα γνωρίσματα σε κάθε μια κλάση ισοδυναμίας. Η διακριτή I-ποικιλομορφία δεν εμποδίζει πιθανολογικές επιθέσεις συμπερασμάτων. Μια κλάση ισοδυναμίας μπορεί να εμφανίζει μια τιμή πιο συχνά σε σχέση με άλλες, επιτρέποντας στον επιτιθέμενο να καταλήξει στο συμπέρασμα ότι μια οντότητα ,σε αυτή την κλάση ισοδυναμίας, είναι πολύ πιθανό να έχει ίδια τιμή. Αυτό οδήγησε στην ανάπτυξη των 2 ακόλουθων ισχυρότερων εννοιών της I - ποικιλομορφίας.

ii) **Εντροπία I - ποικιλομορφίας.** Η εντροπία E μιας κλάσης ισοδυναμίας ορίζεται ως:

$$Entropy(E) = - \sum_{s \in S} p(E, s) \log p(E, s)$$

Όπου S είναι το πεδίο ορισμού του ευαίσθητου γνωρίσματος, και $p(E, s)$ είναι το κλάσμα των εγγραφών της εντροπίας E με s την είσοδο των γνωρισμάτων.

Ένας πίνακας λέμε ότι έχει εντροπία I -ποικιλομορφίας εάν για κάθε κλάση ισοδυναμίας E , $Entropy(E) \geq \log I$. Δηλαδή, η εντροπία ολόκληρου του πίνακα πρέπει να είναι το λιγότερο $\log(I)$. Μερικές φορές αυτό μπορεί να είναι πολύ περιοριστικό, καθώς η εντροπία ολόκληρου του πίνακα μπορεί να είναι μικρή εάν λίγες τιμές είναι πολύ συνηθισμένες. Αυτό οδηγεί στη παρακάτω, λιγότερο περιοριστική έννοια της I -ποικιλομορφίας.

iii) **Επαναλαμβανόμενη (c, I) – ποικιλομορφία.** Η επαναλαμβανόμενη ποικιλομορφία εξασφαλίζει ότι η πιο συχνή τιμή δεν θα εμφανίζεται πολύ συχνά, και οι λιγότερο συχνές τιμές δεν θα εμφανίζονται πολύ σπάνια. Εάν m είναι ο αριθμός των τιμών σε μια κλάση ισοδυναμίας, και r_i με $1 \leq i \leq m$ είναι ο αριθμός των φορών όπου ο $i^{\text{ο}}$ πιο συχνή τιμή εμφανίζεται σε μια κλάση ισοδυναμίας E . Τότε E λέγεται ότι έχει επαναλαμβανόμενη (c, I) - ποικιλομορφία εάν $r_1 < c(r_1 + r_{1+1} + \dots + r_m)$.

Ο πίνακας αναφέρεται ότι έχει επαναλαμβανόμενη (c, I) - ποικιλομορφία εάν όλες οι ισοδύναμες τάξεις της έχουν επαναλαμβανόμενη(c, I) – ποικιλομορφία.

6.4.1 Περιορισμοί I - ποικιλομορφίας

Ενώ η αρχή - κανόνας της I - ποικιλομορφίας αντιπροσωπεύει ένα σημαντικό βήμα στην βελτίωση της k -ανωνυμίας, ως προς την προστασία της αποκάλυψης γνωρισμάτων, έχει αρκετές ελλείψεις που και θα αναφέρουμε.

Η εφαρμογή της I - ποικιλομορφίας μπορεί να είναι δύσκολη και αχρείαστη.

Παράδειγμα 1

Εάν υποθέσουμε ότι τα αρχικά δεδομένα έχουν μόνο ένα ευαίσθητο γνώρισμα: π.χ το αποτέλεσμα ενός συγκεκριμένου ιού. Θα παίρνει τότε μόνο 2 τιμές: θετικό και αρνητικό. Ακόμη, αν υποθέσουμε ότι υπάρχουν 10000 εγγραφές, με το 99% να είναι αρνητικές, και μόνο το 1% θετικές. Τότε οι 2 τιμές έχουν διαφορετικό βαθμό ευαισθησίας. Ο ένας δεν θα πείραζε να είναι γνωστό ότι βγήκε αρνητικός, καθότι είναι ίδιος με το υπόλοιπο 99% του πληθυσμού, αλλά ο άλλος δεν θα ήθελε να είναι γνωστό ότι βγήκε θετικός. Σε αυτή την περίπτωση η 2- διαφορετικότητα είναι αχρείαστη σε μια κλάση ισοδυναμίας που περιλαμβάνει εγγραφές που είναι μόνο αρνητικές. Για να μπορέσουμε να έχουμε πίνακα διακριτής 2 – διαφορετικότητας πρέπει να είναι το πολύ $10000 \times 1\% = 100$ ισοδύναμες τάξεις και η απώλεια της πληροφορίας να είναι μεγάλη. Επίσης, άξιο παρατήρησης είναι ότι εξαιτίας της πολύ μικρής εντροπίας των ευαίσθητων γνωρισμάτων του πίνακα, εάν κάποιος χρησιμοποιήσει εντροπία I - ποικιλομορφίας, το I θα πρέπει να έχει μικρή τιμή.

Η I - ποικιλομορφία είναι ανεπαρκής στην αποτροπή της αποκάλυψης των γνωρισμάτων.

Ακολουθούν 2 επιθέσεις στην 1 - ποικιλομορφία.

i) Επίθεση αλλοίωσης (Skewness Attack): Όταν η συνολική κατανομή είναι στρεβλή, η εφαρμογή της 1 – ποικιλομορφίας δεν πρόκειται να αποτρέψει την αποκάλυψη των γνωρισμάτων.

Κρίνοντας από το παραπάνω Παράδειγμα (1) αν υποθέσουμε ότι μια κλάση ισοδυναμίας έχει ίδιο αριθμό αρνητικών όσο και θετικών εγγραφών. Ικανοποιεί διακριτή 2- ποικιλομορφία, εντροπία 2-ποικιλομορφίας, και οποιαδήποτε επαναλαμβανόμενη (c,2) – ποικιλομορφία μπορεί να επιβληθεί. Ωστόσο, αυτό αποτελεί σοβαρή απειλή για την ιδιωτικότητα, γιατί ο οποιοσδήποτε μέσα στην τάξη μπορεί να θεωρηθεί ότι έχει 50% πιθανότητα να είναι θετικός, συγκρινόμενος με το 1% του συνολικού πληθυσμού.

ii) Επίθεση ομοιότητας (Similarity Attack): Όταν οι τιμές των ευαίσθητων γνωρισμάτων στην κλάση ισοδυναμίας είναι διακριτές αλλά σημασιολογικά παρόμοια, ο επιτιθέμενος μπορεί να μάθει σημαντικές πληροφορίες. Εξατάζουμε το παρακάτω παράδειγμα:

Παράδειγμα 2 Ο πίνακας 8.1 είναι ο αρχικός πίνακας, και ο Πίνακας 8.2 δείχνει την ανώνυμη εκδοχή του (με εφαρμογή διακριτής και εντροπία 3 – ποικιλομορφίας). Υπάρχουν 2 ευαίσθητα γνωρίσματα:

Ο Μισθός και η Αρρώστια. Εάν υποθέσουμε ότι η καταγραφή του Πέτρου βρίσκεται σε μια από τις 3 πρώτες εγγραφές, τότε σίγουρα γνωρίζουμε ότι ο μισθός του είναι μεταξύ 3K – 5K και μπορούμε να υποθέσουμε ότι είναι σχετικά χαμηλός. Αυτή η επίθεση ισχύει όχι μόνο για αριθμητικά γνωρίσματα αλλά και γνωρίσματα κατηγορίας όπως η “Αρρώστια”. Γνωρίζοντας λοιπόν ότι οι εγγραφές του Πέτρου βρίσκονται στην πρώτη τάξη ισοδυναμίας συμπεραίνουμε ότι ο Πέτρος έχει προβλήματα με το στομάχι του, γιατί όλες αυτές οι καταγραφές αφορούν αρρώστιες του στομάχου.

Αυτή η ροή ευαίσθητων πληροφοριών συμβαίνει παρόλο που εφαρμόζεται 1-ποικιλομορφία, η οποία εγγυάται ποικιλομορφία στις ευαίσθητες τιμές της κάθε ομάδας, καθότι δεν λαμβάνει υπόψη τη σημασιολογική εγγύτητα αυτών των τιμών.

	T.K	Ηλικία	Μισθός	Αρρώστια
1	47677	29	3.000	Γαστρικό έλκος
2	47602	22	4.000	Γαστρίτιδα
3	47678	27	5.000	Καρκίνος στομάχου
4	47905	43	6.000	Γαστρίτιδα
5	47909	52	11.000	Γρίπη
6	47906	47	8.000	Βρογχίτιδα
7	47605	30	7.000	Βρογχίτιδα
8	47673	36	9.000	Πνευμονία
9	47607	32	10.000	Καρκίνος στομάχου

Πίνακας 8.1: Πίνακας αρχικών στοιχείων μισθού/αρρώστιας

	T.K	Ηλικία	Μισθός	Αρρώστια
1	476**	2*	3.000	Γαστρικό έλκος
2	476**	2*	4.000	Γαστρίτιδα
3	476**	2*	5.000	Καρκίνος στομάχου
4	4790*	≥40	6.000	Γαστρίτιδα
5	4790*	≥40	11.000	Γρίπη
6	4790*	≥40	8.000	Βρογχίτιδα
7	476**	3*	7.000	Βρογχίτιδα
8	476**	3*	9.000	Πνευμονία
9	476**	3*	10.000	Καρκίνος στομάχου

Πίνακας 8.2 Πίνακας 3 – ποικιλομορφίας

Για την αντιμετώπιση τέτοιων επιθέσεων, ορίστηκε η t - εγγύτητα, μια συμπληρωματική έννοια ιδιωτικότητας που αναπαριστά το γενικότερο γνωστικό πεδίο του επιτιθέμενου πάνω στην κατανομή των τιμών του ευαίσθητου γνωρίσματος.

6.5 t - εγγύτητα

“Μια τάξη ισοδυναμίας λέγεται ότι έχει t - εγγύτητα εάν, η απόσταση μεταξύ της κατανομής ενός ευαίσθητου γνωρίσματος σε αυτή την τάξη και της κατανομής του γνωρίσματος σε ολόκληρο τον πίνακα δεν υπερβαίνει ένα t όριο. Τότε λέγεται ότι ένας πίνακας έχει t -εγγύτητα, όταν όλες οι ισοδύναμες τάξεις έχουν t -εγγύτητα”. [17]

Η βασική ιδέα της t -εγγύτητας απαιτεί ότι, η κατανομή κάθε ευαίσθητου γνωρίσματος σε κάθε ομάδα θα πρέπει να είναι παρόμοια με εκείνη του συνολικού πίνακα. Επιπλέον, για την μέτρηση της προαναφερόμενης μέτρησης της απόστασης, οι συγγραφείς εισήγαγαν έναν νέο όρο – μέτρο μέτρησης της απόστασης, την EMD (Earth Mover Distance). Η σταθερά t χρησιμοποιείται ως όριο για την ικανοποίηση της t - εγγύτητας.

Η ιδιωτικότητα μετράται από το κέρδος των πληροφοριών που θα αποκομίσει ένας παρατηρητής. Πριν δει το δημοσιοποιημένο πίνακα, ο παρατηρητής έχει μόνο κάποια σχετική ιδέα σχετικά με τα ευαίσθητα γνωρίσματα – πληροφορίες ενός ατόμου. Από την στιγμή που δει τον δημοσιοποιημένο πίνακα, αποκτάει μια διαφορετική εικόνα - άποψη. Το κέρδος της πληροφορίας μπορεί να αναπαρασταθεί ως η διαφορά μεταξύ της προηγούμενης ιδέας και της μεταγενέστερης ιδέας.

Για την καλύτερη κατανόηση του συγκεκριμένου κανόνα αλλά και την έκφραση του με τύπο, ακολουθεί παράδειγμα.

Παράδειγμα:

Υποθέτουμε πρότερη εικόνα - ιδέα παρατηρητή B_0 για τα ευαίσθητα γνωρίσματα καποιου ιδιώτη. Ύστερα από κάποια υποθετική ενέργεια, δίνεται στον παρατηρητή μια πλήρη γενικευμένη έκδοση ενός πίνακα δεδομένων όπου όλα τα γνωρίσματα στα ψευδο-αναγνωριστικά έχουν αφαιρεθεί (ή έχουν γενικευθεί στις πιο γενικευμένες τιμές τους). Η ιδέα του παρατηρητή επηρεάζεται από την κατανομή Q της τιμής του ευαίσθητου γνωρίσματος σε όλο τον πίνακα και αλλάζει σε B_1 . Τέλος, δίνεται στον παρατηρητή ο δημοσιευμένος πίνακας. Γνωρίζοντας τις τιμές των ψευδο-αναγνωριστικών, ο παρατηρητής δύναται να αναγνωρίσει την τάξη ισοδυναμίας, όπου βρίσκεται η καταγραφή του ιδιώτη, και μαθαίνει την κατανομή P των τιμών των ευαίσθητων γνωρισμάτων στην τάξη. Η άποψη του παρατηρητή αλλάζει σε B_2 .

Η σκοπός της 1 – ποικιλομορφίας είναι ο περιορισμός της μείωσης της διαφοράς μεταξύ B_0 και B_2 . Μειώνοντας τη διαφορά μεταξύ B_1 και B_2 καταλήγουμε στο συμπέρασμα ότι Q η κατανομή των ευαίσθητων γνωρισμάτων στον συνολικό πληθυσμό του πίνακα, είναι η δημόσια πληροφορία. Δεν περιορίζεται το κέρδος της πληροφορίας των παρατηρητών ως σύνολο, αλλά περιορίζεται ο βαθμός στον οποίο ο παρατηρητής μπορεί να μάθει περισσότερα πράγματα για συγκεκριμένα άτομα.[17]

Για να δικαιολογηθεί η παραδοχή ότι η Q θα πρέπει να διαχειριστεί σαν δημόσια πληροφορία, παρατηρείται ότι με τις γενικεύσεις το μέγιστο που κάποιος μπορεί να καταφέρει είναι να γενικεύσει όλες τις τιμές των ψευδο-αναγνωριστικών στην πιο γενική τιμή τους. Έτσι με την δημοσίευση μιας έκδοσης των δεδομένων θα δημοσιευθεί παράλληλα και μια κατανομή Q . Η μεγάλη διαφορά του B_0 στο B_1 , σημαίνει ότι ο πίνακας περιέχει περισσότερα στοιχεία και νέα πληροφορία ως σύνολο. Όσο μεγαλύτερη είναι η διαφορά του B_0 στο B_1 τόσο πιο πολύτιμα είναι τα δεδομένα. Δεδομένου ότι η αύξηση της πληροφορίας μεταξύ B_0 και B_1 αφορά ολόκληρο το πληθυσμό, δεν περιορίζουμε αυτό το κέρδος.

Περιορίζουμε το κέρδος μεταξύ του B_1 B_2 , περιορίζοντας την απόσταση μεταξύ του P και Q . Έτσι αν $P = Q$ τότε $B_1 = B_2$. Αν P και Q είναι κοντά, τότε ανάλογα B_1 και B_2 θα είναι το ίδιο κοντά, ακόμα και αν B_0 μπορεί να είναι διαφορετικό από B_1 και B_2 . [17]

Απαιτώντας όμως τα P και Q να είναι κοντά θα περιορίζε το σύνολο της οφέλιμης δημοσιευμένης πληροφορίας, καθώς περιορίζει την συσχέτιση μεταξύ των ψευδο-αναγνωριστικών και των ευαίσθητων γνωρισμάτων. Όμως αυτό ακριβώς είναι το ζητούμενο. Εάν λάβει καθαρή εικόνα σχετικά με αυτή την συσχέτιση, τότε επέρχεται και η αποκάλυψη των γνωρισμάτων. Η παράμετρος t στην t - εγγύτητα επιτρέπει την ανταλλαγή μεταξύ χρηστικότητας και ιδιωτικότητας.

Το επόμενο βήμα είναι η μέτρηση της απόστασης μεταξύ δυο πιθανολογικών κατανομών. Υπάρχουν αρκετοί τρόποι για να οριστεί η απόσταση μεταξύ τους. Βάση των δύο κατανομών $P = (p_1, p_2, \dots, p_m)$, $Q = (q_1, q_2, \dots, q_m)$ δύο γνωστά μέτρα απόστασης είναι τα ακόλουθα:

i) Η μεταβολική απόσταση (variational distance) ορίζεται ως:

$$D[P, Q] = \sum_{i=1}^m \frac{1}{2} |p_i - q_i|.$$

και

ii) Η Kullback-Leibler (KL) απόσταση ορίζεται ως:

$$D[P, Q] = \sum_{i=1}^m p_i \log \frac{p_i}{q_i} = H(P) - H(P, Q)$$

όπου

$H(P) = \sum_{i=1}^m p_i \log p_i$ είναι η εντροπία του P και

$H(P, Q) = \sum_{i=1}^m p_i \log q_i$ είναι η cross-εντροπία του P και Q .

Αυτά τα μέτρα απόστασης δεν αντικατοπτρίζουν τη σημασιολογική απόσταση μεταξύ των τιμών

Αν και ορίζονται, η μεταβολική απόσταση και η απόσταση Kullback-Leibler, ο σκοπός της μέτρησης της απόστασης έχει σκοπό την αναπαράσταση της διαφοροποίησης των τιμών του ευαίσθητου γνωρίσματος με τέτοιο τρόπο ώστε μικρότερη απόσταση να ερμηνεύεται ως: λιγότερο διαφορετικές τιμές και συνεπώς μεγαλύτερος κόπος εξόρυξης της προσωπικής πληροφορίας. Την καλύτερη αναπαράσταση της απόστασης, επιτυγχάνει η μετρική Earth Mover's Distance (που αρχικώς προαναφέραμε), η οποία ορίζεται από το πρόβλημα μετακίνησης. Βασίζεται στο ελάχιστο απαιτούμενο έργο για την μετατροπή της μίας κατανομής στην άλλη με την μαζική ανακατανομή ανάμεσά τους.

Έτσι αν $P = (p_1, p_2, \dots, p_m)$, $Q = (q_1, q_2, \dots, q_m)$ με d_{ij} τη σταθερή απόσταση μεταξύ του στοιχείου i του P και στοιχείου j του Q . Θέλουμε να βρούμε μια ροή $F = [f_{ij}]$ όπου f_{ij} η ροή της μάζας από το p_i στο q_j έτσι ώστε να ελαχιστοποιηθεί το συνολικό έργο το οποίο είναι ίσο με την μετρική EMD και δίνεται από τη σχέση:

$$WORK(P, Q, F) = \sum_{i=1}^m \sum_{j=1}^m d_{ij} f_{ij}$$

με τους ακόλουθους περιορισμούς:

$$f_{ij} \geq 0 \quad 1 \leq i \leq m, 1 \leq j \leq m \quad (c1)$$

$$p_i - \sum_{j=1}^m f_{ij} + \sum_{j=1}^m f_{ji} = q_i \quad 1 \leq i \leq m \quad (c2)$$

$$\sum_{i=1}^m \sum_{j=1}^m f_{ij} = \sum_{i=1}^m p_i = \sum_{i=1}^m q_i = 1 \quad (c3)$$

Αυτοι οι περιορισμοι εγγυώνται ότι ο P μετατρέπεται σε Q απο την ροή της μάζας F . Απο την στιγμή που το πρόβλημα μεταφοράς λυθεί ο EMD ορίζεται ως συνολικό έργο

$$D[P, Q] = WORK(P, Q, F) = \sum_{i=1}^m \sum_{j=1}^m d_{ij} f_{ij}$$

Βάση τη t - εγγύτητα και τους πίνακες που την ικανοποιούν όπως και τη χρήση της EMD προκύπτουν οι δύο κάτωθι ιδιότητες:

Ιδιότητα γενίκευσης: Έστω T ένας πίνακας δεδομένων και A, B δύο πίνακες-γενικεύσεις του T τέτοιες ώστε ο A να είναι πιο γενικευμένος από τον B . Αν ο T ικανοποιεί την t -εγγύτητα με χρήση του B , τότε ο T ικανοποιεί την t -εγγύτητα και με χρήση του πίνακα A .

Ιδιότητα υποσύνολου: Έστω T ένας πίνακας δεδομένων και C το σύνολο των γνωρισμάτων που συμμετέχουν στα δεδομένα του T . Αν ο T ικανοποιεί την t -εγγύτητα ως προς το σύνολο C τότε ικανοποιεί την t -εγγύτητα ως προς κάθε υποσύνολο γνωρισμάτων $D \subset C$.

Οι δύο ιδιότητες εγγυώνται τη δυνατότητα ένταξης της έννοιας της t -εγγύτητας με χρήση της μετρικής EMD στο γενικό πλαίσιο του αλγορίθμου Incognito.

Παρόλο το σημαντικό πλεονέκτημα χρήσης του εν λόγω αλγορίθμου - αρχής, δημιουργούνται κάποια ζητήματα:

i) Δεν είναι εύκολη η προστασία της ιδιωτικότητας σύμφωνα με πολλά διαφορετικά επίπεδα ασφάλειας

ii) Η εισαχθείσα, μέτρηση απόστασης (EMD) εκλείπει της ευελιξίας της διαχείρισης αριθμητικών γνωρισμάτων

iii) Ενδέχεται πιθανότητα “θυσίας” της χρηστικότητας των δημοσιευμένων δεδομένων διότι είναι πολύ δύσκολο ένας κανόνας να επιτρέπει όλες τις διανομές γνωρισμάτων να είναι μεταξύ τους παρόμοιες.

6.6 δ – Παρουσία (δ -Presence)

Η δ -Παρουσία αποσκοπεί στην προστασία του συνόλου των δεδομένων από την αποκάλυψη της ιδιότητας του μέλους. Απαιτεί την ρητή διαμόρφωση του συνόλου δεδομένων που έχει αποκαλυφθεί σε υποσύνολο μεγαλύτερου συνόλου δεδομένων που αντιπροσωπεύει τις γνώσεις του επιτιθέμενου. Το κριτήριο επιτάσσει, περιορισμούς στις πιθανότητες με τις οποίες μπορεί να προσδιοριστεί, εάν ή όχι ένα άτομο από το ολικό σύνολο δεδομένων περιλαμβάνεται στο υποσύνολο της έρευνας. Τα ανώτερα όρια για αυτές τις πιθανότητες υπολογίζονται με βάση τα μεγέθη ισοδύναμων κλάσεων.

6.7 Ανωνυμοποίηση βάση κινδύνου

Η ανωνυμοποίηση βάση του κινδύνου μπορεί να χρησιμοποιηθεί για την αυτόματη μετατροπή των συνόλων δεδομένων για την διασφάλιση ότι οι εκτιμώμενοι κίνδυνοι επαναπροσφορισμού θα υποχωρήσουν – πέσουν κάτω ενός συγκεκριμένου ορίου

Κεφάλαιο 7. Εργαλεία Ανωνυμοποίησης

Για την εφαρμογή ανωνυμοποίησης σε ένα σύνολο δεδομένων υπάρχουν διαθέσιμα πολλά εργαλεία. Χωρίζονται ανάλογα με το πως διατίθενται – δωρεάν, επι πληρωμή, ανοιχτού κώδικα, κλειστού κώδικα - και κατα πόσο ολοκληρωμένα είναι όσον αφορά την ποικιλία των αλγόριθμων ανωνυμοποίησης που υποστηρίζουν όπως των μεθόδων που χρησιμοποιούν.

Θα ασχοληθούμε αποκλειστικά με το ARX, το οποίο θεωρείται ένα από τα πληρέστερα όσον αφορά τα προαναφερθέντα και πρόκειται για ανοιχτού κώδικα λογισμικό. Ακολουθεί μια συνοπτική παρουσίαση αντίστοιχων εργαλείων.

7.1 UTD

Το [UTD Anonymization Toolbox](#) είναι γραμμένο σε JAVA και είναι ανοιχτού κώδικα. Υποστηρίζει 3 διαφορετικά μοντέλα αυθεντικοποίησης (k -anonymity, l -diversity και t -closeness). Το εργαλείο εφαρμόζει τον Incognito και τον [Datafly](#) (αλγόριθμος ευρετικής αναζήτησης). Όταν χρησιμοποιείται ο l -ποικιλομορφία ή ο t -εγγύτητα, ο μετασχηματισμός δεδομένων περιορίζεται σε γενίκευση πλήρους τομέα, διότι το εργαλείο δεν είναι δυνατό να διαχειριστεί μη μονοτονικά προβλήματα ιδιωτικότητας. Αυτό μπορεί να οδηγήσει σε δεδομένα κακής ποιότητας. Επίσης υποστηρίζει ακόμα δύο μοντέλα γενίκευσης: πολυδιάσταση ολική γενίκευση με χρήση του Mondrian αλγόριθμου και του ανατομισμού (μορφή “τεμαχισμού”...δεν αναφέρεται στην παρούσα εργασία). Το εργαλείο χρησιμοποιεί SQLite συστήμα υποστήριξης βάσης δεδομένων. Δεν έχει γραφικό περιβάλλον και σε σχέση με το ARX δεν περιλαμβάνει μεθόδους για ανάλυση κινδύνου ή ανωνυμοποίηση βάσης του ρίσκου.

7.2 CAT

Το εργαλείο ανωνυμοποίησης του πανεπιστημίου του Cornell ([CAT](#)) είναι γραμμένο σε C++. Εφαρμόζει l -ποικιλομορφία και t -εγγύτητα. Σαν μοντέλο γενίκευσης χρησιμοποιεί Ολική Γενίκευση (γενίκευση πλήρους τομέα). Υποστηρίζει μόνο μη αυτόματη απόκρυψη tuple. Για την είσοδο των δεδομένων απαιτεί αυστηρά καθορισμένη μορφή και εφαρμόζει τον αλγόριθμο Incognito. Τέλος παρέχει λίγες μεθόδους αυτόματης αξιολόγησης δεδομένων και των κινδύνων γνωστοποίησης.

7.3 TIAMAT

Το TIAMAT είναι κλειστού κώδικα λογισμικό και είναι γραμμένο σε JAVA. Υποστηρίζει μόνο το μοντέλο της k -ανωνυμίας. Εφαρμόζει τον Mondrian και τον k -Member, (αλγόριθμος ομαδοποίησης). Το μοντέλο γενίκευσης είναι πολυδιάστατης ολικής επαναγενίκευσης με γενίκευση. Το εργαλείο περιλαμβάνει απλό γραφικό πρόγραμμα επεξεργασίας των γενικευμένων ιεραρχιών. Για την αυτόματη μέτρηση της χρησιμότητας των δεδομένων το TIAMAT χρησιμοποιεί την μέθοδο [GCP](#) και Classification Metric[19] Ο στόχος αυτού του εργαλείου είναι να συγκρίνει τον αλγόριθμο του Mondrian με την προσέγγιση του k -Member και εξ αυτού εφαρμόζει διάφορες οπτικοποιήσεις που ομαδοποιούν τους χρόνους εκτέλεσης τους, για καλύτερη χρηστικότητα των δεδομένων

7.4 SECRETA

Το [SECRETA](#) (Ελληνικής ομάδας ερευνητών) είναι κλειστού κώδικα λογισμικό και είναι γραμμένο σε C++. Αναπτύχθηκε για να την σύγκριση των μεθόδων ανωνυμοποίησης σχεσιακών και συναλλακτικών δεδομένων. Για τα δεδομένα συσχέτισης εφαρμόζει αλγόριθμο ομαδοποίησης και Incognito. Τα μοντέλα γενίκευσης για τα δεδομένα συσχέτισης περιλαμβάνουν πλήρη γενίκευση, γενίκευση υποδέντρου και τοπική γενίκευση. Ενσωματώνει εννέα γνωστούς αλγόριθμους ανωνυμοποίησης καθώς και 3 μεθόδους δέσμευσης για τον συνδυασμό των προαναφερόμενων αλγορίθμων κάτω από κοινό πλαίσιο. Επίσης έχει γραφικό περιβάλλον οπτικής αναπαράστασης των χαρακτηριστικών των δεδομένων (όπως το ARX & CAT) με μεθόδους εμπειρικών στατιστικών. Το εργαλείο εφαρμόζει μια προχωρημένη μέθοδο μέτρησης της χρησιμότητας των δεδομένων, ορίζοντας τον φόρτο των αναζητήσεων και προσδιορίζοντας το μέσο σχετικό σφάλμα των αποτελεσμάτων της αναζήτησης το βοηθητικό πρόγραμμα δεδομένων, καθορίζοντας φόρτου εργασίας ερωτήματος και προσδιορίζοντας το μέσο σχετικό σφάλμα (average relative error) των αποτελεσμάτων του ερωτήματος.

Δεδομένου ότι ο στόχος του εργαλείου είναι να συγκρίνει διαφορετικές στρατηγικές ανωνυμοποίησης, χρησιμοποιεί διάφορες οπτικοποιήσεις χρόνων εκτέλεσης και χρησιμότητας δεδομένων.

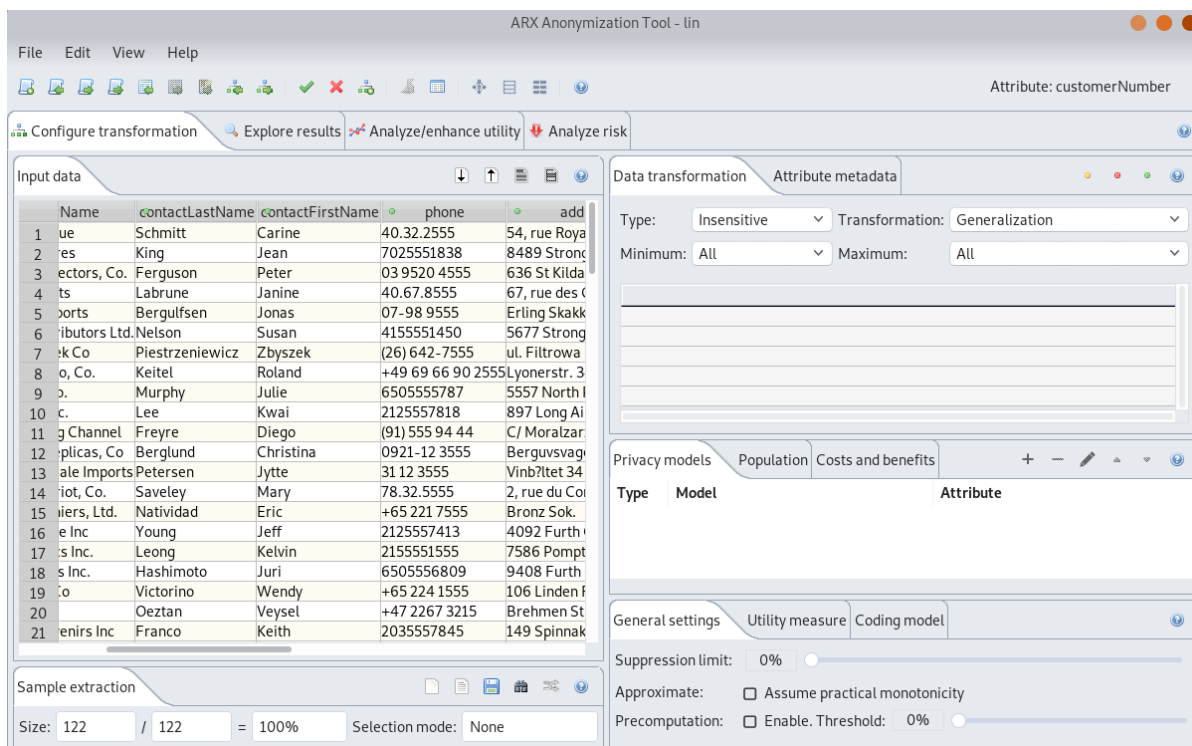
7.5 ARX

Το ARX είναι ένα ολοκληρωμένο λογισμικό ανοιχτού κώδικα για την ανωνυμοποίηση των ευαίσθητων προσωπικών δεδομένων. Υποστηρίζει μια ευρεία ποικιλία (1) μοντέλων ιδιωτικότητας και κινδύνου, (2) μεθόδους μετασχηματισμού δεδομένων και (3) μεθόδους για την ανάλυση της χρησιμότητας των δεδομένων εξόδου.

Τα βασικά πλεονεκτήματα του ARX είναι[21] :

- **Η ανάλυση του κινδύνου:** Το ARX εφαρμόζει πολλαπλές μεθόδους εκτίμησης του κινδύνου του επαναπροσδιορισμού, συμπεριλαμβάνοντας 2 διαφορετικά μοντέλα υπερπληθυσμού.
- **Ανωνυμοποίηση βάση κινδύνου:** Το ARX παρέχει κριτήρια απορρήτου που χρησιμοποιούν τα ενσωματωμένα μοντέλα κινδύνου. Μπορεί αυτόματα να μετασχηματίσει σύνολα δεδομένων προς την εξασφάλιση ότι οι κίνδυνοι είναι κάτω του ορίου που ορίζει ο χρήστης.
- **Συντακτικά κριτήρια απορρήτου:** Το APX υλοποιεί μια ποικιλία από πολλαπλά συντακτικά κριτήρια απορρήτου , συμπεριλαμβανομένης της k-ανωνυμίας, 3 παραλλαγών της l-diversity, 2 παραλλαγές του t-closeness και δ-presence (δεν αναφέρεται στην παρούσα). Επιπλέον το ARX υποστηρίζει αυθαίρετους συνδυασμούς κριτηρίων προστασίας της ιδιωτικότητας (συμπεριλαμβανομένων των μεθόδων που βασίζονται στον κίνδυνο)
- **Αξιολόγηση χρησιμότητας:** Το ARX εφαρμόζει διάφορες μεθόδους αξιολόγησης της χρησιμότητας των δεδομένων, τόσο αυτόματα όσο και κατά προσωπική παραμετροποίηση. Όσον αφορά τον πρώτο (αυτόματα) οι μέθοδοι είναι πλήρως παραμετροποιήσιμοι.

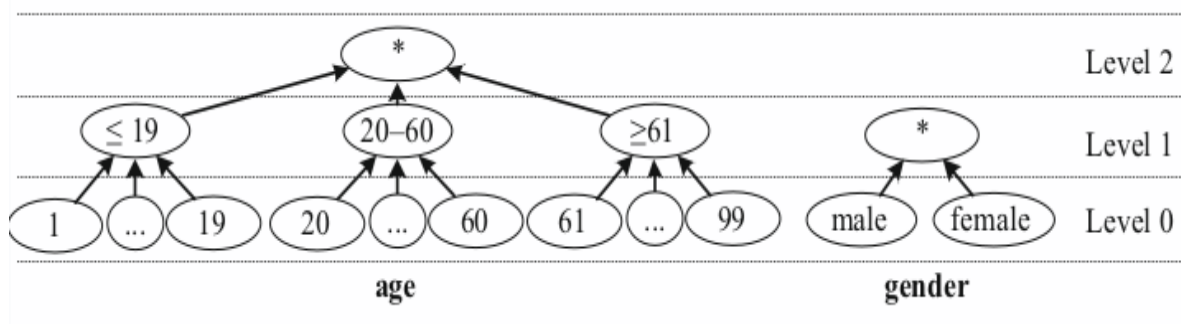
- **Μοντέλο “διαισθητικού” μετασχηματισμού:** Το ARX χρησιμοποιεί την ολική γενίκευση μέσω της πλήρους γενίκευσης τομέα – domain σε συνδυασμό με την τοπική γενίκευση μέσω της απόκρυψης tuple. Αυτό το σχήμα γενίκευσης έχει προταθεί και στον τομέα της βιοϊατρικής. Επιπλέον, υποστηρίζει την top – και bottom γενίκευση καθώς και της μικροσυσσωμάτωσης.
- **Βέλτιστη λειτουργία:** Το ARX υλοποιεί μια ολική βέλτιστη στρατηγική αναζήτησης όπου αυτόματα θα ταξινομεί το χώρο λύσης και θα καθορίζει τον βέλτιστο μετασχηματισμό σύμφωνα με το καθορισμένο μέτρο χρησιμότητας
- **Επεκτασιμότητα:** Το ARX έχει σχεδιαστεί από την αρχή στο να είναι ιδιαίτερα επεκτάσιμο, και είναι ικανό να διαχειριστεί πολύ μεγάλα σύνολα δεδομένων (εκατομμυρίων εισαγωγών) σε οικιακά μηχανήματα.
- **Συμβατότητα:** Το ARX υποστηρίζει δεδομένα εισόδου από πηγές format, όπως .csv αρχεία, excel φύλλα, και σχεσιακών βάσεων δεδομένων συστήματα (MS SQL Server, MySQL, PostgreSQL)
- **Πλήρες γραφικό περιβάλλον:** Το ARX παρέχει ένα ολοκληρωμένο γραφικό περιβάλλον με χαρακτηριστικά όπως οδηγούς δημιουργίας ιεραρχιών γενίκευσης και διάφορες οπτικοποιήσεις των εκτιμήσεων κινδύνου, λύσεων βασιζόμενων στον χώρο και της χρησιμότητας των δεδομένων
- **Προσεκτικά σχεδιασμένο API:** Το ARX προσφέρει επίσης ένα προγραμματιστικό περιβάλλον (API) που επιτρέπει πρόσβαση σε όλες τις λειτουργίες του. Τόσο οι μέθοδοι όσο και το γραφικό περιβάλλον είναι προσεκτικά σχεδιασμένα.



Εικόνα 1: Το γραφικό περιβάλλον του ARX σε Linux

- **Cross-platform:** Το ARX έχει υλοποιηθεί σε Java και είναι διαθέσιμο σε όλα τα κοινά λειτουργικά συστήματα.

- **Ανοιχτού Κώδικα:** Το ARX είναι ανοιχτού κώδικα λογισμικό, παροτρύνοντας τις κριτικές από την κοινότητα και επιτρέποντας στους χρήστες να το προσαρμόζουν στις δικές του απαιτήσεις



Σχήμα 2: Ιεραρχίες γενίκευσης για τα χαρακτηριστικά ηλικία και γένος

Στην ενότητα 1.1 αναλύσαμε τις τρεις απειλές της ιδιωτικότητας

Για την αντιμετώπιση των απειλών κατά της ιδιωτικότητας με το ARX, τα δεδομένα μετασχηματίζονται με ιεραρχίες γενίκευσης. Το εργαλείο υποστηρίζει ιεραρχίες τόσο για κατηγορηματικά όσο και για συνεχή γνωρίσματα.

Παραδείγματα φαίνονται στο σχήμα 2. Εδώ, οι τιμές του χαρακτηριστικού ηλικία μεταβάλλονται σε “ηλικιακές ομάδες”, ενώ οι τιμές της ιδιότητας του φύλου μπορούν μόνο να κρυφτούν.

Οι ιεραρχίες γενίκευσης είναι κατάλληλες για κατηγορηματικά γνωρίσματα. Μπορούν επίσης να χρησιμοποιηθούν για συνεχή χαρακτηριστικά, εκτελώντας επί τόπου κατηγοριοποίηση. Στο ARX, αυτή η κατηγοριοποίηση υλοποιείται με την απεικόνιση στρατηγικών γενίκευσης με ένα λειτουργικό τρόπο, π.χ., ως ένα σύνολο διαστημάτων.

Για την αύξηση της χρηστικότητας των δεδομένων, η γενίκευση των γνωρισμάτων συνδυάζεται με την απόκρυψη των εγγραφών των αρχείων. Αυτό σημαίνει ότι οι γραμμές που “παραβιάζουν” τα κριτήρια ιδιωτικότητας διαγράφονται αυτόματα από το σύνολο των δεδομένων. Το σύνολο των διεγγραμμένων καταγραφών φυλάσσεται κάτω από ένα καθορισμένο όριο, το οποίο ονομάζεται όριο απόκρυψης (suppresion limit). Ως αποτέλεσμα αυτού απαιτείται λιγότερη γενίκευση για την εξασφάλιση της ικανοποίησης του μοντέλου ιδιωτικότητας από τις υπόλοιπες εγγραφές.[21]

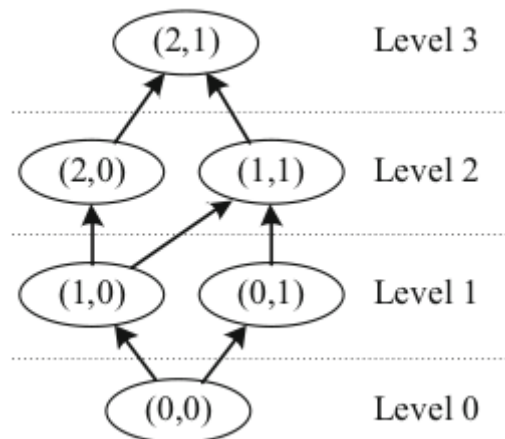
Το ARX εφαρμόζει πολυ-διάσταση ολική γενίκευση με πλήρους τομέα γενίκευση και τοπική γενίκευση με απόκρυψη tuple.

Αυτός ο συνδυασμός των μεθόδων έχει βρεθεί να ταιριάζει τέλεια στο τομέα της βιοϊατρικής, καθότι είναι: i) διαισθητική, ii) Παράγει σύνολα δεδομένων που είναι κατάλληλα αναλύσεων από επιδημιολόγους και iii) μπορούν να διαμορφώνονται από μη ειδικούς. Η παραμετροποίηση (και ως εκ τούτου η ισορροπία της ιδιωτικότητας και της χρηστικότητας) μπορεί να πραγματοποιηθεί μεταβάλλοντας τις ιεραρχίες γενίκευσης ή απλά επιλέγοντας ένα κατάλληλο μετασχηματισμό από το σύνολο των λύσεων (solution space). Για συνεχείς μεταβλητές, το ARX υποστηρίζει την μικροσυσσωμάτωση. Αυτό σημαίνει ότι οι τιμές του quasi-identifier μέσα σε μια ομάδα μετασχηματίζονται εφαρμόζοντας αθροιστικές λειτουργίες, όπως το αριθμητικό ή το γεωμετρικό μέσο. Σε έναν ανώνυμο πίνακα, όλη η ομάδα των τιμών αντικαθιστάται από το

αποτέλεσμα της συναθροιστικής συνάρτησης. Σε αντίθεση με την γενίκευση ή την απόκρυψη, αυτή η μέθοδος μπορεί να χρησιμοποιηθεί για τη διατήρηση των τύπων των δεδομένων και κλίμακα μέτρησης μιας αριθμητικής μεταβλητής.[21]

Όταν χρησιμοποιούμε ολική γενίκευση με γενίκευση πλήρους τομέα, ο χώρος αναζήτησης μπορεί να διαμορφωθεί ως επαναλαμβανόμενη δομή γενίκευσης, η οποία είναι ένα σύνολο διατεταγμένων πιθανών συνδυασμών των επιπέδων γενίκευσης για κάθε γνώρισμα. Οι επαναλαμβανόμενες δομές μπορούν να απεικονιστούν ως διαγράμματα Hasse, που στην περίπτωση του πλαισίου μας σημαίνει, την σχεδίαση της μεταβατικής μείωσης ενός συνόλογου από μετασχηματισμούς σε ένα άκυκλο γράφημα, όπου ο κάθε κόμβος θα είναι συνδεδεμένος με όλους τους άμεσους διαδόχους και προκάτοχους του. Ο κάθε κόμβος αντιπροσωπεύει ένα μετασχηματισμό και ορίζει τα επίπεδα γενίκευσης όλων των quasi-identifiers. Ένα βέλος υποδηλώνει ότι ένας μετασχηματισμός είναι μια άμεση γενίκευση ενός πιο εξειδικευμένου μετασχηματισμού, που μπορεί να δημιουργεί αυξάνοντας τα επίπεδα γενίκευσης που ορίζει ο προκάτοχος του.

Ένα παράδειγμα που χρησιμοποιεί τις ιεραρχίες από το σχήμα 2, απεικονίζεται στο σχήμα 3.



Σχήμα 3: Παράδειγμα χώρου αναζήτησης

7.5.1 Η Ανάλυση κινδύνου και η Ανωνομοποίηση βάση κινδύνου με την χρήση του ARX

Στο ARX, μπορούν να χρησιμοποιηθούν μοντέλα υπερπληθυσμού για την εκτίμηση της μοναδικότητας του πληθυσμού, δηλαδή το κλάσμα των καταχωρήσεων στο σύνολο των δεδομένων που είναι μοναδικά στο σύνολο του πληθυσμού.

Αυτές οι στατιστικές μέθοδοι υπολογίζουν χαρακτηριστικά του συνολικού πληθυσμού με κατανομές πιθανοτήτων που παραμετροποιούνται με τα χαρακτηριστικά του δείγματος.

Επίσης, οι εκτιμήσεις κινδύνου μπορούν να χρησιμοποιηθούν για την ανωνυμοποίηση βάση κινδύνου. Το στατικό μοντέλο k-ανωνυμίας ορίζει ένα ανώτερο όριο στον επαναπροσδιορισμό του κινδύνου, που υπολογίζεται βάση δείγματος συχνοτήτων. Το ARX υποστηρίζει πολλαπλές παραλλαγές αυτού του κριτηρίου που χρησιμοποιούν τα προαναφερόμενα μοντέλα κινδύνου επιβεβαιώνοντας ότι οι εκτιμήσεις κινδύνου υπολείπονται ενός συγκεκριμένου ορίου.

7.5.2 Αξιολόγηση χρησιμότητας

Για την αυτόματη μέτρηση της χρησιμότητας των δεδομένων, το ARX υποστηρίζει μεθόδους που βασίζονται σε κλάσεις ισοδυναμίας, που ονομάζονται μετρήσεις μονοδιάστατης χρησιμότητας μέσα στο εργαλείο. Παραδείγματα περιλαμβάνουν την *Διακριτότητα* και το *Μέγεθος μέσης ισοδύναμης κλάσης*. Και τα 2 μέτρα υπολογίζουν το περιοχόμενο της πληροφορίας βάσει του μεγέθους των τάξεων ισοδυναμίας, ενώ η ευκρίνεια συνεπάγεται επίσης της ποινής για τις κρυμμένες tuples.

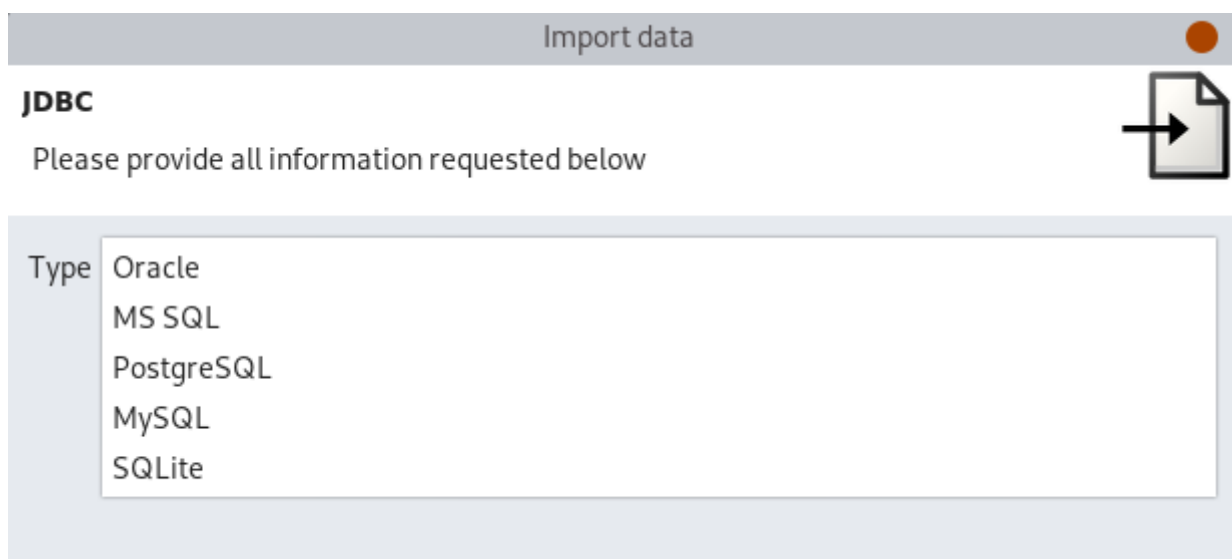
Επιπρόσθετα, το ARX υποστηρίζει μεθόδους που βασίζονται σε τιμές πραγματικών δεδομένων.

Καθώς αυτές οι μέθοδοι υπολογίζουν ανεξάρτητες τιμές για κάθε γνώρισμα, οι οποίες στην συνέχεια μεταγλωττίζονται σε μια ολική τιμή και αποκαλούνται πολυ-διάστατες μετρήσεις χρησιμότητας. Επίσης το ARX υποστηρίζει αρκετές συναθροιστικές λειτουργίες για την σύνταξη μέτρων ολικής χρησιμότητας, συμπεριλαμβανομένου του αθροίσματος, του αριθμητικού και του γεωμετρικού μέσου. [21]

7.5.3 Επιπρόσθετα Χαρακτηριστικά

Το ARX είναι συμβατό με ένα ευρύ φάσμα εργαλείων επεξεργασίας δεδομένων. Στην παρούσα υποστηρίζει διεπαφές για την εισαγωγή δεδομένων για αρχεία CSV, MS Excel φύλλα, συστήματα διαχείρισης σχεσιακών βάσεων δεδομένων (RDBMS), όπως MS SQLServer, PostgreSQL and MySQL.

Το συντακτικό των CSV αρχείων αναγνωρίζεται αυτόματα. Σε σημασιολογικό επίπεδο, το ARX υποστηρίζει διαφορετικούς τύπους δεδομένων και κλίμακες μέτρησης. Τα format των τύπων δεδομένων εντοπίζονται εξίσου αυτόματα. Επίσης υποστηρίζει εκκαθάριση δεδομένων κατά την διάρκεια της εισόδου, εννοώντας ότι τα μη-έγκυρα δεδομένα αντικαθίστονται με μια τιμή που αντιπροσωπεύει τα στοιχεία που λείπουν. Παράλληλα χειρίζεται με ασφάλεια τις ελλείπουσες τιμές, διασφαλίζοντας ότι οι διαφορετικές τιμές που λείπουν, αντιστοιχούν μεταξύ τους.



Εικόνα 2: Επιλογές βάσεων δεδομένων

7.5.4 Χρήση του ARX

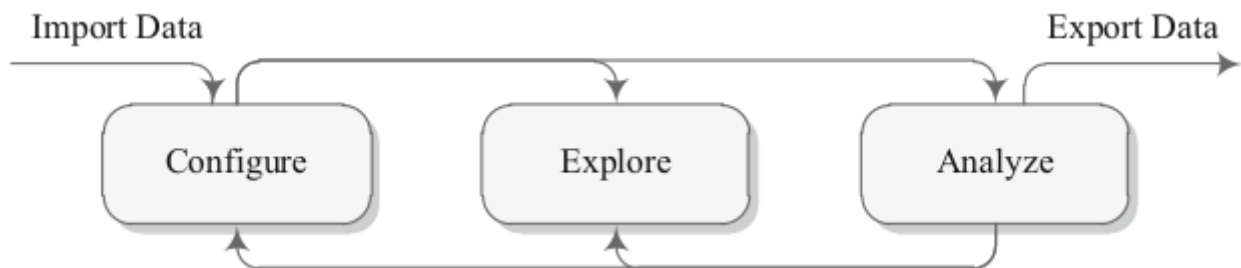
Στην κάτωθι ενότητα θα περιγραφεί αρχικώς η ροή εργασιών που υποστηρίζεται από το γραφικό περιβάλλον χρήστη του ARX. Εν συνεχεία θα παρουσιαστεί αναλυτικότερα το γραφικό περιβάλλον και θα περιγραφούν αναλυτικότερα τα προαναφερόμενα

7.5.4.1 Η διαδικασία της ανωνυμοποίησης

Η κεντρική πρόκληση στην ανωνυμοποίηση δεδομένων είναι η επίτευξη ισορροπίας μεταξύ της χρηστικότητας των δεδομένων και της ιδιωτικότητας. Στο ARX, οι μέθοδοι που μοντελοποιούν τις διαφορετικές πτυχές αυτής της διαδικασίας συνδυάζονται σε μια ροή εργασίας πολλαπλών σταδίων που επιτρέπει στους χρήστες να προσαρμόζουν με παραμέτρους τις παραμέτρους, μέχρις ότου το αποτέλεσμα να ανταποκρίνεται στις προσδοκίες τους. Όπως απεικονίζεται στο Σχήμα 3, τα βασικά βήματα αποτελούνται από

- i) Τη παραμετροποίηση του μοντέλου ιδιωτικότητας και μετασχηματισμού
- ii) Εξερεύνηση του χώρου της λύσης
- iii) Ανάλυση εισερχόμενων και εξερχόμενων δεδομένων. Στη φάση διαμόρφωσης, τα δεδομένα εισόδου φορτώνονται και σχολιάζονται, δημιουργούνται ή εισάγονται ιεραρχίες γενίκευσης και καθορίζονται όλες οι υπόλοιπες παράμετροι, όπως τα κριτήρια ιδιωτικότητας.

Όταν ο χώρος λύσης χαρακτηρίζεται από την εκτέλεση του ARX αλγόριθμου ανωνυμοποίησης, η φάση εξερεύνησης υποστηρίζει την αναζήτηση μετασχηματισμών δεδομένων που προστατεύουν την ιδιωτικότητα και πληρούν τις απαιτήσεις του χρήστη. Για να εγκριθεί η καταλληλότητα, η φάση της ανάλυσης επιτρέπει τη σύγκριση μετασχηματισμένων συνόλων δεδομένων στο αρχικό σύνολο δεδομένων εισόδου με μεθόδους περιγραφικών στατιστικών. Επιπλέον, μπορούν να πραγματοποιηθούν αναλύσεις κινδύνου για δεδομένα εισόδου όπως και μετασχηματισμένες παραστάσεις τους. Βάση αυτών των αναλύσεων, μπορούν να εξεταστούν και περαιτέρω πιθανές λύσεις και να αξιολογηθούν ή η διαμόρφωση της διαδικασίας ανωνυμίας μπορεί να μεταβληθεί.



- | | | |
|--|---|---|
| - Προσδιορισμός κανόνων παραμετροποίησης | - Καθαρισμός και ανάλυση του χώρου των λύσεων | - Σύγκριση & ανάλυση εισαγωγής και εξαγωγής |
| - Προσδιορισμός μοντέλου ιδιωτικότητας | - Οργάνωση παραμετροποιήσεων | - Κίνδυνοι και χρηστικότητα |
| - Προσδιορισμός μοντέλου γενίκευσης | | |

Διάγραμμα 1. Η διαδικασία της ανωνυμοποίησης, εφαρμοσμένη στο γραφικό περιβάλλον του ARX

Τα 3 βήματα της διαδικασίας της ανωνυμοποίησης αντιστοιχίζονται σε τέσσερις προοπτικές στο γραφικό περιβάλλον του ARX:

Παραμετροποίηση:

Σε αυτή την προοπτική, *πρώτον*, ένα σύνολο δεδομένων μπορεί να εισαχθεί στο εργαλείο και να επισημανθεί, π.χ., χαρακτηρισμός γνωρίσματος.

Δεύτερον, με απευθείας εισαγωγή ή χρήση του οδηγού, οι ιεραρχίες γενίκευσης για τα ψευδοαναγνωριστικά ή ευαίσθητα γνωρίσματα μπορούν να δημιουργηθούν ημι-αυτόματα

Τρίτον: Μπορούν να προσδιοριστούν τα κριτήρια ιδιωτικότητας, η μέθοδος μέτρησης χρησιμότητας δεδομένων και άλλες παράμετροι, όπως οι ιδιότητες του μοντέλου μετασχηματισμού.

Εξερεύνηση:

Ως αποτέλεσμα της διαδικασίας της ανωνυμοποίησης, δημιουργείται ένας χώρος λύσεων και χαρακτηρίζεται βάση των δεδομένων παραμέτρων. Αυτή η προοπτική επιτρέπει στους χρήστες την περιήγηση των διαθέσιμων μετασχηματισμένων δεδομένων, να οργανώνουν και να τα φιλτράρουν ανάλογα με τις ανάγκες τους, και να επιλέγουν επιπλέον μετασχηματισμούς για ανάλυση.

Αξιολόγηση χρησιμότητας:

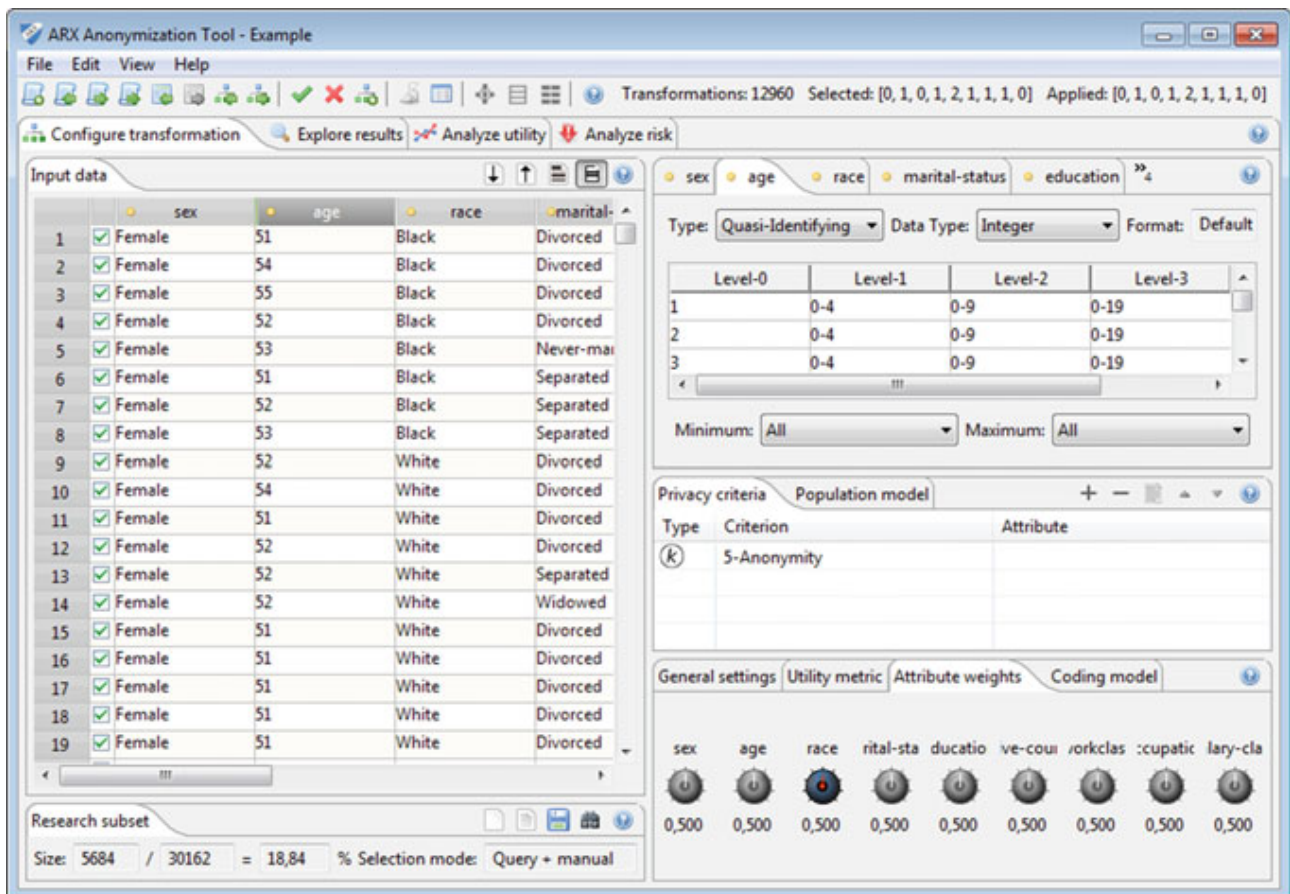
Η συγκεκριμένη προοπτική επιτρέπει τη σύγκριση μετασχηματισμών του συνόλου δεδομένων εισόδου με τα αρχικά δεδομένα, για να αξιολογήσει της καταλληλότητας ενός συγκεκριμένου μετασχηματισμού για το δεδομένο σενάριο χρήσης. Για το σκοπό αυτό, ενσωματώνει διάφορες γραφικές παραστάσεις των αποτελεσμάτων των στατιστικών αναλύσεων ενώ επιτρέπει την σύγκριση των πεδίων 1-1.

Ανάλυση Κινδύνου:

Σε αυτή την προοπτική, μπορεί να αναλυθεί η κατανομή των μεγεθών τάξης, οι κίνδυνοι που σχετίζονται με τα μεμονομένα ψευδοαναγνωριστικά, όπως επίσης εκτιμήσεις κινδύνου βάση δειγμάτων και πληθυσμού. Εμφανίζονται επίσης λεπτομέρειες σχετικά με εκτιμώμενους κινδύνους επαναπροσδιορισμού που προέρχονται από διαφορετικά μοντέλα.

7.5.4.2 Παραμετροποιώντας την διαδικασία ανωνυμοποίησης

Σε αυτήν την προοπτική, που φαίνεται και στην εικόνα 3, μπορεί να “σχολιαστεί” (χαρακτηριστεί) ένα σύνολο δεδομένων και να παραμετροποιηθεί η διαδικασία της ανωνυμοποίησης. Με την έννοια του σχόλιου, εννοούμε τον χαρακτηρισμό ενός γνώρισματος ως ψευδογνώρισμα ή ευαίσθητο, ο ορισμός ιεραρχιών γενίκευσης, ο ορισμός υποσυνόλου έρευνας και η εισαγωγή πληροφοριών του συνολικού πληθυσμού. Η παραμετροποίηση σημαίνει, τον προσδιορισμό του μοντέλου ιδιωτικότητας και διαμόρφωση της διαδικασίας μετασχηματισμού. Στη αριστερή πλευρά της εικόνας



Εικόνα 3: Η παραμετροποίηση του ARX

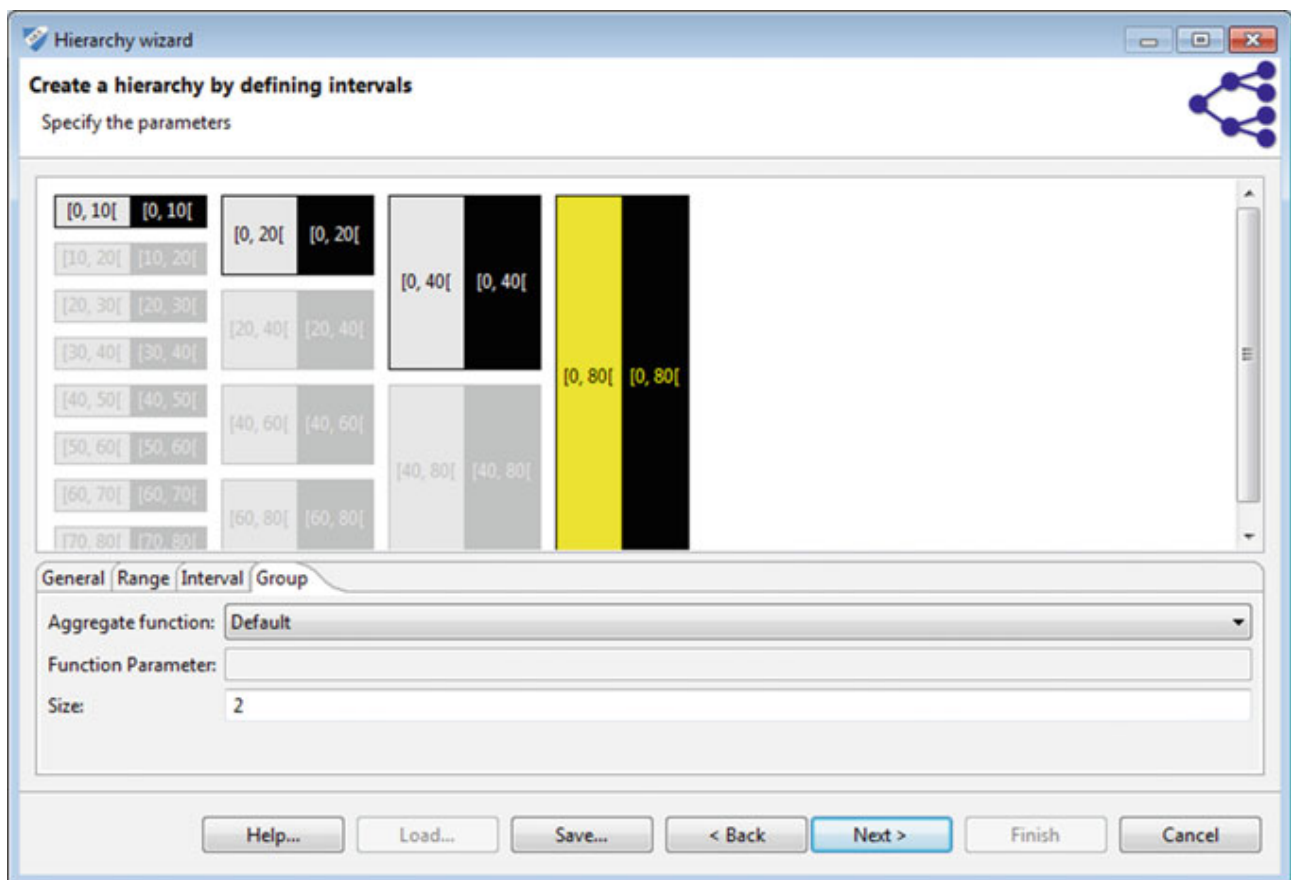
Μια διεπαφή αναζήτησης, η οποία βρίσκεται στην κάτω αριστερή γωνία μπορεί να χρησιμοποιηθεί για την επιλογή ενός υποσύνολου αναζήτησης, το οποίο είναι ουσιαστικά ένα υποσύνολο αρχείων δεδομένων τα οποία θα συμπεριληφθούν στο τελικό ανώνυμο σύνολο δεδομένων. Αυτό ο μηχανισμός μπορεί να χρησιμοποιηθεί για διάφορους σκοπούς.

Πρώτον, μπορεί να χρησιμοποιηθεί για την επιβολή δ-παρουσίας στο υποσύνολο, για να προστατέψει το υποσύνολο από την αποκάλυψη της ιδιότητας του μέλους, εμποδίζοντας τους επιτιθέμενους που γνωρίζουν τη συνολική ομάδα δεδομένων να καθορίσουν εάν ένα συγκεκριμένο περιέχεται ή όχι στο υποσύνολο.

Δεύτερον, μπορεί να χρησιμοποιηθεί για τη χειροκίνητη κατάργηση των καταχωρήσεων των δεδομένων. Η περιοχή στην επάνω δεξιά γωνία περιγράφει όλα τα χαρακτηριστικά και δίνει την επιλογή των “χαρακτηρισμών”, καθορίζοντας τύπους δεδομένων και τύπους γνωρισμάτων. Στην ίδια περιοχή εμφανίζονται επίσης οι ιεραρχίες γενίκευσης που σχετίζονται με τα γνωρίσματα.

Το ARX διακρίνει τέσσερις διαφορετικούς τύπους γνωρίσματος:

1. Χαρακτηριστικά γνωρίσματα, θα αφαιρεθούν από το σύνολο δεδομένων. Τέτοιου τύπου είναι π.χ. τα ονόματα ή οι αριθμοί κοινωνικής ασφάλισης.
2. Τα ψευδοαναγνωρίσματα θα μετασχηματιστούν μέσω της γενίκευσης, απόκρυψης ή μικροσυσσωμάτωσης. Τυπικά παραδείγματα είναι το φύλο, η ημερομηνία γέννησης και οι ταχυδρομικοί κώδικες.
3. Τα ευαίσθητα γνωρίσματα θα διατηρηθούν χωρίς τροποποίηση, αλλά ενδέχεται να υπόκεινται σε περαιτέρω περιορισμούς, όπως t-εγγύτητα t ή l- ποικιλομορφία. Χαρακτηριστικά παραδείγματα είναι οι διαγνώσεις.
4. Τα μη ευαίσθητα γνωρίσματα θα παραμείνουν αμετάβλητα.



Εικόνα 4: Οδηγός δημιουργίας της ιεραρχίας γενίκευσης με διαστήματα

Για τον καθορισμό των ιεραρχιών γενίκευσης, το εργαλείο προσφέρει έναν οδηγό, ο οποίος υποστηρίζει τη δημιουργία ιεραρχιών για γνωρίσματα με διαφορετικές κλίμακες μέτρησης. Ο οδηγός για τη δημιουργία ιεραρχιών με διαστήματα εμφανίζεται στην εικόνα 4.

Αρχικά, μπορεί να οριστεί μια ακολουθία διαστημάτων. Στο επόμενο βήμα, μπορούν να οριστούν τα επόμενα επίπεδα της ιεραρχίας, τα οποία αποτελούνται από ομάδες διαστημάτων από το προηγούμενο επίπεδο. Όπως μπορεί να φανεί, κάθε ακολουθία διαστημάτων ή ομάδων επαναλαμβάνεται αυτόματα για να καλύψει το πλήρες εύρος του γνωρίσματος. Κάθε στοιχείο συσχετίζεται με μια συνθετική συνάρτηση. Αυτές οι συναρτήσεις ενσωματώνουν μεθόδους για την δημιουργία ετικετών για τα διαστήματα, ομάδες και τιμές που πρέπει να μεταφραστούν σε μια κοινή γενικευμένη τιμή.

Από τη προοπτική της παραμετροποίησης, οι ιεραρχίες οπτικοποιούνται με την μορφή πινάκων, όπου κάθε σειρά τους περιέχει τον κανόνα γενίκευσης για μια μοναδική τιμή γνωρίσματος.

Αυτή η “διαισθητική” αναπαράσταση επιτρέπει τη συμβατότητα με τρίτες εφαρμογές όπως τα προγράμματα υπολογιστικών φύλλων (excel). Οι ιεραρχίες μπορούν επίσης να ρυθμιστούν με έναν ενσωματωμένο επεξεργαστή, ο οποίος είναι διαθέσιμος κάνοντας δεξί κλικ στην πινακοποιημένη αναπαράστασή τους.

Η άποψη στο κέντρο της δεξιάς πλευράς της προοπτικής υποστηρίζει την προδιαγραφή του μοντέλου ιδιωτικότητας. Για το σκοπό αυτό, εμφανίζει μια λίστα με τα κριτήρια ιδιωτικότητας που είναι ενεργοποιημένα και υποστηρίζονται σαν επιλογές η προσθήκη, η αφαίρεση και διαμόρφωση των κριτηρίων μέσω των αντίστοιχων πλαισίων. Στη δεύτερη καρτέλα, μπορούν να οριστούν τα χαρακτηριστικά του πληθυσμού που απαιτούνται για αναλύσεις κινδύνου και ανωνυμοποίηση βάσει κινδύνου.

Επιλογές για τη παρεμετροποίηση της διαδικασίας μετασχηματισμού ενσωματώνονται στις προβολές στο κάτω μέρος της δεξιάς πλευράς της προοπτικής. Στις γενικές ρυθμίσεις, μπορεί να καθοριστεί το όριο κατάργησης και να προσαρμοστούν περαιτέρω παράμετροι που σχετίζονται με την απόδοση. Στη δεύτερη καρτέλα, μπορούν να επιλεγούν και να διαμορφωθούν μέτρα χρησιμότητας. Οι άλλες δύο καρτέλες επιτρέπουν στους χρήστες να παραμετροποιήσουν το μοντέλο μετασχηματισμού με τη σταθμίζοντας γνωρίσματα και με την ιεράρχηση διαφορετικών τύπων γενίκευσης δεδομένων.

7.5.4.3 Διαχείριση του χώρου των αποτελεσμάτων

Όταν ο χώρος αναζήτησης έχει χαρακτηριστεί με βάση τις συγκεκριμένες παραμέτρους, η προοπτική εξερεύνησης επιτρέπει στους χρήστες να αναζητούν πιθανές λύσεις, να τις οργανώνουν και να τις φιλτράρουν σύμφωνα με τις απαιτήσεις τους. Σκοπός της προοπτικής (εικόνα 5) είναι η επιλογή ενός συνόλου ενδιαφερόντων μετασχηματισμών για περαιτέρω αξιολόγηση.

The screenshot displays the ARX Anonymization Tool interface. The main window shows a list of transformations with columns for Transformation, Anonymity, Min. Info. Loss, and Max. Info. Loss. A tooltip is visible over one of the rows, showing detailed information loss for various attributes. Below the list, there are three panels: Filter, Clipboard, and Properties.

Transformation	Anonymity	Min. Info. Loss	Max. Info. Loss
[0, 2, 0, 1, 2, 1, 1, 1, 0]	ANONYMOUS	0.28307655989412295 [0,000%]	0.28307655989412295 [0,000%]
[0, 2, 0, 1, 1, 1, 1, 1, 0]	ANONYMOUS	0.30521601292778766 [3,088%]	0.30521601292778766 [3,088%]
[0, 2, 0, 1, 2, 1, 0, 1, 0]	ANONYMOUS	0.3311292560270289 [6,703%]	0.3311292560270289 [6,703%]
[0, 1, 0, 1, 2, 1, 1, 1, 0]	ANONYMOUS	0.3315901556540435 [6,767%]	0.3315901556540435 [6,767%]
[0, 2, 0, 1, 2, 1, 1, 0, 0]	ANONYMOUS	0.3385456412213499 [7,737%]	0.3385456412213499 [7,737%]
[0, 1, 0, 1, 1, 1, 1, 1, 0]	ANONYMOUS	0.3623878009327741 [11,063%]	0.3623878009327741 [11,063%]
[0, 2, 0, 1, 1, 1, 0, 1, 0]	ANONYMOUS	0.3678197028861945 [11,820%]	0.3678197028861945 [11,820%]
[0, 2, 0, 1, 1, 1, 1, 0, 0]	ANONYMOUS	0.3759756805657366 [12,958%]	0.3759756805657366 [12,958%]
[0, 1, 0, 1, 2, 1, 0, 1, 0]	ANONYMOUS	0.3910095752566891 [15,055%]	0.3910095752566891 [15,055%]
[0, 2, 0, 1, 2, 1, 0, 0, 0]	ANONYMOUS	0.4161224233049243 [18,558%]	0.4161224233049243 [18,558%]
[0, 1, 0, 1, 2, 1, 1, 0, 0]	ANONYMOUS	0.4288125783288368 [20,328%]	0.4288125783288368 [20,328%]

Filter Panel:

Attribute	0	1	2	3	4
sex	✓	✗			
age	✗	✓	✓	✗	✗
race	✓	✗			
marital-status	✗	✓	✗		

Clipboard Panel:

Node	Comment
[0, 2, 0, 1, 2, 1, 1, 1, 0]	Minimal information loss
[0, 1, 0, 1, 2, 1, 1, 1, 0]	Age is less generalized

Properties Panel:

Property	Value
Anonymous	ANONYMOUS
Min. info. loss	0.3315901556540435 [6,767%]
Max. info. loss	0.3315901556540435 [6,767%]
Successors	9
Predecessors	6
Transformation	[0, 1, 0, 1, 2, 1, 1, 1, 0]
Checked	true

Στο κέντρο αυτής της οπτικής εμφανίζεται ένα υποσύνολο του χώρου λύσης. Στο παράδειγμα μας, οι μετασχηματισμοί παρουσιάζονται ως λίστα, ο οποίος ταξινομείται βάση χρηστικότητας δεδομένων. Κάθε μετασχηματισμός εκπροσωπείται από τα επίπεδα γενίκευσης που καθορίζει για τα ψευδοαναγνωριστικά στην είσοδο του σύνολο δεδομένων. Οι μετασχηματισμοί χαρακτηρίζονται από τέσσερα διαφορετικά χρώματα:

το πράσινο υποδεικνύει ότι ένας μετασχηματισμός οδηγεί σε ένα ανώνυμο σύνολο δεδομένων.

Το κόκκινο υποδεικνύει ότι ο μετασχηματισμός δεν έχει ως αποτέλεσμα ένα ανώνυμο σύνολο δεδομένων.

Το κίτρινο υποδεικνύει ότι ο μετασχηματισμός είναι ο βέλτιστος σε πλήρες επίπεδο, και το γκρι δείχνει ότι η ιδιότητα ανωνυμίας ενός μετασχηματισμού είναι άγνωστη ως αποτέλεσμα του κλάδου.

Αν ένας τέτοιος μετασχηματισμός εφαρμοστεί σε σύνολο των δεδομένων, η πραγματική ιδιότητά της ανωνυμίας θα υπολογιστεί στο παρασκήνιο και η κατάσταση του χώρου αναζήτησης θα ενημερωθεί.

Το κάτω αριστερό σημείο επιτρέπει το φιλτράρισμα του εμφανιζόμενου υποσυνόλου του χώρου λύσης, π.χ., περιορίζοντας τους μετασχηματισμούς σε ορισμένα επίπεδα γενίκευσης ή καθορίζοντας κατώτατα όρια στην χρηστικότητα των δεδομένων. Οι μετασχηματισμοί μπορούν επίσης να προστεθούν στο πρόχειρο, το οποίο βρίσκεται στο κάτω-κεντρικό τμήμα της οπτικής μας. Εδώ, οι μετασχηματισμοί μπορούν να σχολιαστούν και να οργανωθούν. Η κάτω δεξιά περιοχή εμφανίζει βασικές πληροφορίες για τον τρέχων επιλεγμένο μετασχηματισμό. Κάνοντας δεξιά κλικ σε ένα μετασχηματισμό εμφανίζεται ένα μενού περιεχομένου, το οποίο επιτρέπει την εφαρμογή του στο σύνολο δεδομένων. Το αποτέλεσμα μπορεί στη συνέχεια να εξαχθεί ή να αναλυθεί.

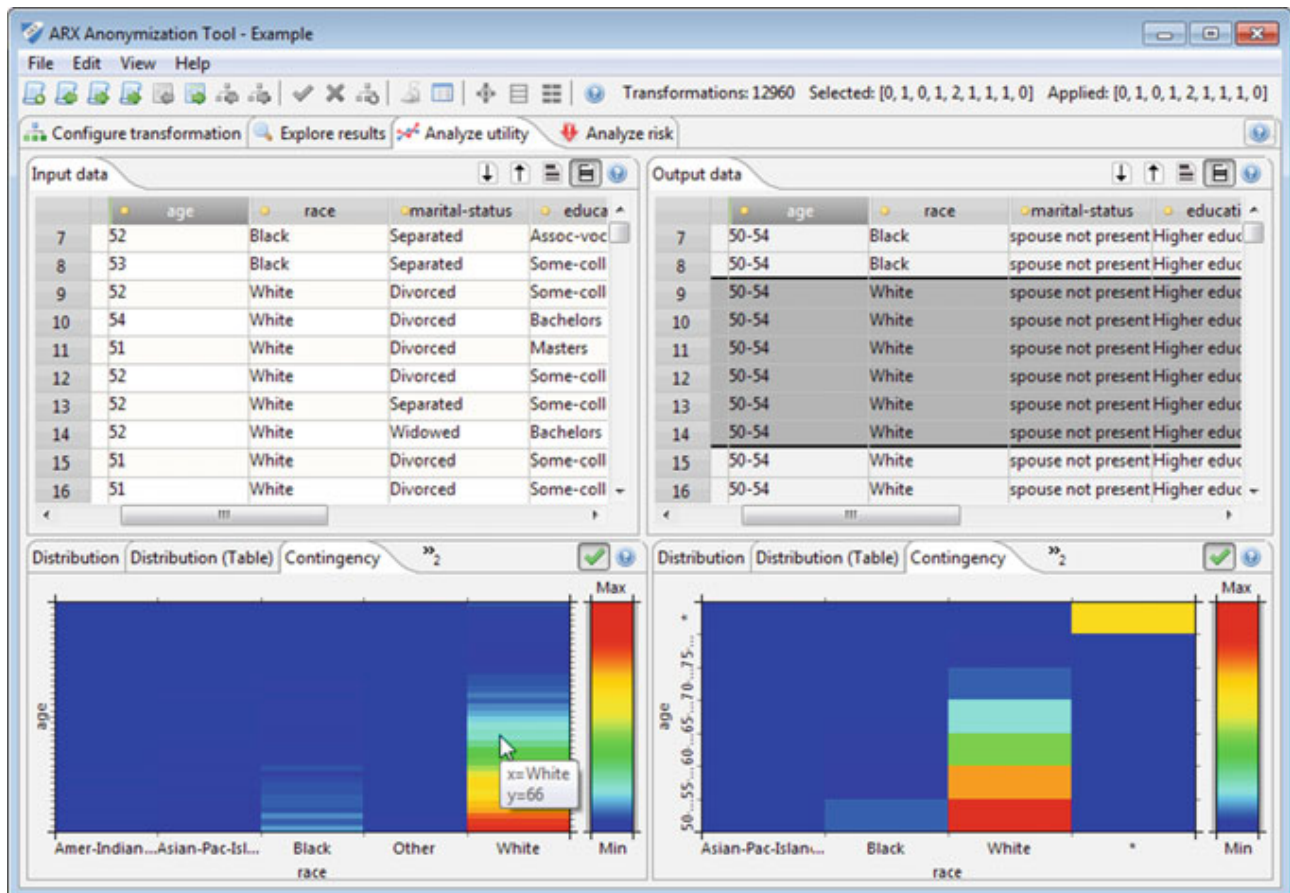
Δύο επιπλέον οπτικοποιήσεις του χώρου λύσεων παρέχονται. Πρώτον, μπορεί να απεικονιστεί ως διάγραμμα Hasse. Εδώ, χρησιμοποιούνται τα ίδια χρώματα όπως στην προβολή λίστας για την κωδικοποίηση της ιδιότητας ανωνυμίας. Η προβολή είναι διαδραστική και υποστηρίζει ζουμ. Δεύτερον, ο χώρος των αποτελεσμάτων μπορεί επίσης να απεικονιστεί ως σύνολο “πλακιδίων”, τα οποία ταξινομούνται και χρωματίζονται με βάση την χρηστικότητα των δεδομένων. Σε αντίθεση με τις άλλες δύο παραστάσεις, αυτή η οπτική είναι ικανή να εμφανίζει ταυτόχρονα μεγάλο αριθμό μετασχηματισμών.

7.5.4.4 Αξιολόγηση Χρηστικότητας Δεδομένων

Αυτή η προοπτική επιτρέπει στους χρήστες να συγκρίνουν τις παραμετροποιημένες απεικονίσεις δεδομένων με το αρχικό σύνολο δεδομένων. Ένα παράδειγμα φαίνεται στη εικόνα 6. Εμφανίζει το σύνολο δεδομένων εισόδου στην αριστερή πλευρά και το σύνολο δεδομένων εξόδου στη δεξιά πλευρά. Και οι δύο πίνακες είναι συγχρονίζονται κατά την κύλιση, υποστηρίζοντας συγκρίσεις μεταξύ των πεδίων. Τα δεδομένα μπορούν να ταξινομηθούν κατά ένα μόνο χαρακτηριστικό ή κατά όλα τα ψευδοαναγνωριστικά, τα οποία επίσης θα επισημάνουν τις οι προκύπτουσες κλάσεις ισοδυναμίας, εάν εφαρμοστούν στην έξοδο του σύνολο δεδομένων.

Για τη σύγκριση δύο απεικονίσεων των δεδομένων, η προβολή εμφανίζει τυπικές περιγραφικές στατιστικές, τόσο σε μορφή πίνακα όσο και σε μορφή γραφικών. Περιλαμβάνουν μεθόδους μοναδικής και ποικίλης ποικιλίας, όπως εμπειρικές κατανομές, κεντρική τάση και

διασπορά, καθώς και πίνακες διασταύρωσης. Μαζί με τις στατιστικά μετρήσεις, αυτή η οπτική εμφανίζει επίσης πληροφορίες αντανακλώντας τη διαδικασία ανωνυμοποίησης.



Εικόνα 6 Αξιολόγηση Χρηστικότητα Δεδομένων στο ARX

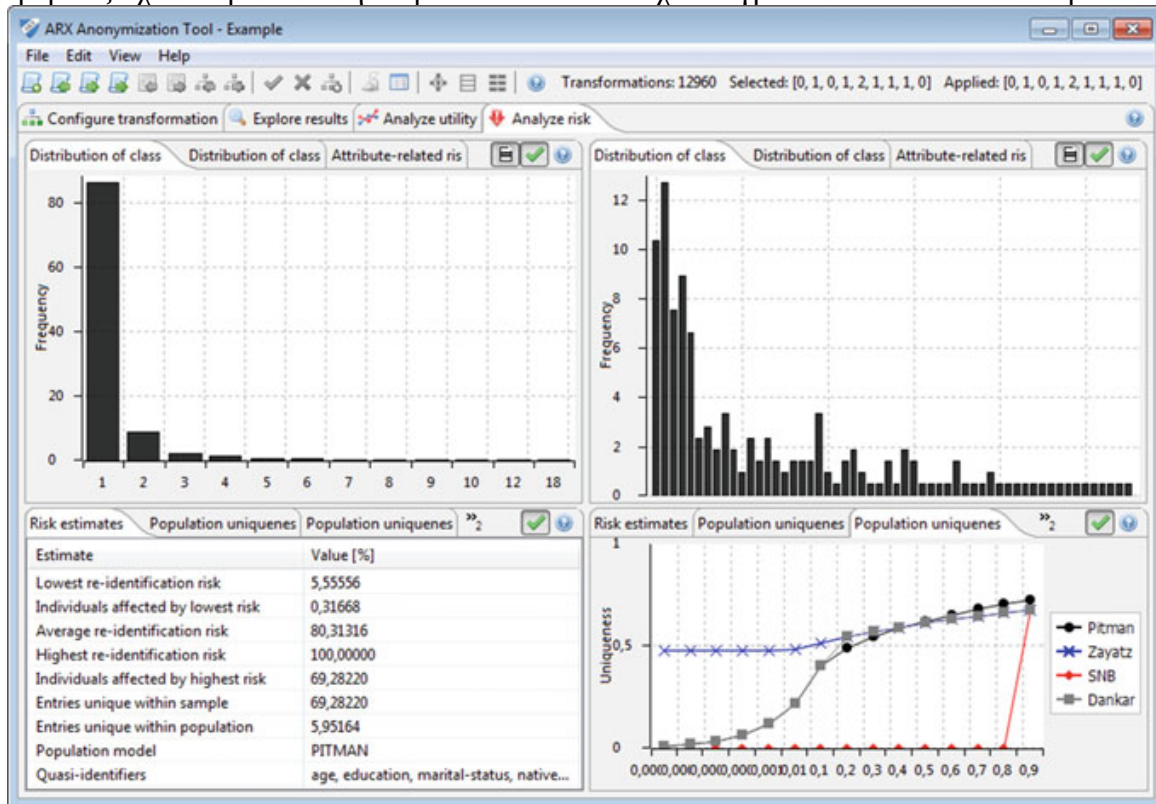
7.5.4.5 Ανάλυση κινδύνων επαναπροσδιορισμού

Με αυτή την προοπτική, οι κίνδυνοι επαναπροσδιορισμού ενός συνόλου δεδομένων εισόδου ή εξόδου, μπορούν να αναλυθούν με διάφορα μέτρα βασισμένα στο δείγμα ή στο πληθυσμό. Επιπλέον, αυτή η προοπτική παρέχει μια σύνοψη της κατανομής των μεγεθών των κλάσεων ισοδυναμίας στο σύνολο δεδομένων και υποστηρίζει την ανάλυση των κινδύνων που σχετίζονται με ομάδες γνωρισμάτων για την εύρεση ψευδοαναγνωριστικών.

Στο εικόνα 7, εμφανίζεται η κατανομή των μεγεθών τάξης τόσο εισόδου όσο και εξόδου. Η άποψη στην κάτω δεξιά γωνία δείχνει μια σύγκριση της εκτιμώμενης μοναδικότητας του πληθυσμού από τρία μοντέλα κινδύνου για διαφορετικά κλάσματα δειγμάτων. Παρουσιάζει επίσης τα αποτελέσματα του κανόνα απόφασης που προτείνεται και επικυρώνεται για τα κλινικά σύνολα δεδομένων από τους Dankar et al., τα οποία επιλέγουν ένα διαφορετικό μοντέλο κινδύνου για τα διαφορετικά κλάσματα του δείγματος [20]

Κατα τον υπολογισμό εκτιμήσεων κινδύνου, το ARX πρέπει να επιλύει αριθμητική συστήματα εξισώσεων. Ο υπολογιστής που χρησιμοποιείται από το ARX μπορεί να ρυθμιστεί στο πλαίσιο των ρυθμίσεων. Εδώ μπορείτε να ορίσετε επιλογές όπως ο συνολικός αριθμός επαναλήψεων, ο μέγιστος αριθμός επαναλήψεων ανά δοκιμή και η απαιτούμενη ακρίβεια. Αυτό μπορεί να επηρεάσει την ακρίβεια των αποτελεσμάτων και των χρόνων εκτέλεσης. Επιπλέον, στην κάτω αριστερή γωνία, η προβολή εμφανίζει βασικές εκτιμήσεις κινδύνου που υπολογίζονται είτε από το ίδιο το δείγμα είτε από ένα μοντέλο κινδύνου.

Αυτή η ανάλυση χρησιμοποιεί εκτιμήσεις της μοναδικότητας είτε βασισμένη σε ένα δείγμα είτε σε ένα στατιστικό μοντέλο. Όταν ένα σύνολο γνωρισμάτων έχει επιλεγεί για ανάλυση, το ARX θα καθορίσει τον μέσο όρο κινδύνου επαναπροσδιορισμού που σχετίζεται με το σύνολο ισχύος αυτών των γνωρισμάτων. Αυτή η πληροφορία μπορεί να φανεί χρήσιμη στην απόφαση του ποια γνωρίσματα πρέπει να γενικευτούν. Επιπλέον, μπορούν να ρυθμιστούν σε ξεχωριστό πλαίσιο οι λεπτομέρειες σχετικά με τον πληθυσμό από τον οποίο έχει δειγματοληφτεί το σύνολο δεδομένων



Εικόνα 7: Η ανάλυση κινδύνου του ARX

7.6 Τελικές παρατηρήσεις

Μέσα από αυτή την συνοπτική παρουσίαση των εργαλείων ανωνυμοποίησης, κρίνοντας όμως κυρίως από το ARX που δοκιμάσαμε, παρατηρούμε ότι η πρόοδος που υπήρξε τα τελευταία χρόνια πάνω σε ένα όχι τόσο γνωστό θέμα στο ευρύ κοινό - αυτό της ανωνυμοποίησης - είναι αξιοσημείωτη. Τόσο από πλευράς πληρότητας όσο ποικιλίας, δεδομένου ότι η ανωνυμοποίηση των δεδομένων παραμένει ένα πολύπλοκο ζήτημα που θα πρέπει να εξετάζεται αλλά και εκτελείται από ειδικούς. Δυστυχώς δεν υπάρχει ενιαίο μέτρο που να είναι ικανό να προστατέψει τις βάσεις δεδομένων από όλους τους πιθανούς κινδύνους, ειδικά αν είναι αρκετά ευέλικτο για να υποστηρίξει πολλά σενάρια χρήσης.

Σαν εργαλείο το ARX θεωρείται πλήρες και με συνεχόμενη εξέλιξη, καθότι έχει ενεργή ομάδα ανάπτυξης και είναι ανοιχτού κώδικα. Με κάθε νέα έκδοση γίνονται νέες προσθήκες για να υποστηρίχουν νέες λειτουργίες.

Κεφάλαιο 8. Επίλογος

Με την παρούσα διπλωματική εργασία έγινε μια προσπάθεια παρουσίασης του προβλήματος της εφαρμογής ανωνυμίας – απο την απλή ομαδοποίηση μέχρι τις πιο σύνθετες γενικεύσεις - σαν μια οργανική συνολική κατάσταση εξέλιξης και βελτίωσης του τρόπου αντιμετώπισης του μέσα απο την σκοπιά των νέων κανόνων, αλγόριθμων και εφαρμογών που την απαρτίζουν.

Δώσαμε μεγαλύτερη βάση στην ανάλυση της k-ανωνυμίας, τις αδυναμίες της, τις επιθέσεις που τις εκμεταλλεύονται και τους τρόπους επίλυσης μέσω νέων αλγορίθμων ανωνυμοποίησης. Βάση αυτού αναλύουμε περισσότερο τους δύο πιο δυνατούς και καταλήγουμε στην παρουσίαση έτοιμων εφαρμογών, ακαδημαϊκού χαρακτήρα, που μπορούν να τους εφαρμόσουν. Επικεντρωνόμαστε στο ARX, και αναλύουμε τον τρόπο ανάλυσης του προβλήματος της ανωνυμίας μέσω των ανάλυσης που γίνεται απο την φύση του κώδικα του.

Τέλος, να αναφέρουμε, ότι το συνεχόμενο ενδιαφέρον των ερευνητικών ομάδων γι αυτή – ανωνυμία- συντελεί στην εξέλιξη της και στους νέους τρόπους εφαρμογής της. Εκμεταλλεζόμενη την εποχή της συνειδητοποίησης του τρόπου λειτουργίας του διαδικτύου, χρησιμοποιείται περισσότερο και απο περισσότερους.

Δείγμα εφαρμογής της μπορούμε να δούμε στις υλοποιήσεις που γίνονται, εκτός των εργαλείων που προαναφέραμε. Τέτοιο παράδειγμα αποτελεί το <https://haveibeenpwned.com/> όπου ο δημιουργός του συγκεκριμένου εφαρμόζει k-ανωνυμία, των queries προς την βάση σε σύγκριση με της εισόδου του επισκεπτόμενου.

Κύριο κίνητρο για την συγγραφή της ήταν η παρουσίαση αυτής και του τρόπου που μπορεί να γίνει επιτεύξιμη. Ενώ παραμένει η ελπίδα χρήσης της όλο περισσότερο απο τις Ελληνικές Αρχές, Υπηρεσίες, Οργανισμούς και Εταιρείες στο πλαίσιο ικανοποίησης του GDPR.

Βιβλιογραφία

- 1] Raymond Choo, Man Ho Au. Mobile Security and Privacy
- 2] Aris Gkoulalas-Divanis, Grigorios Loukides. Medical Data Privacy Handbook && Li, T., Li, N., Zhang, J., Molloy, I.: Slicing: a new approach for privacy preserving data publishing.
- 3] Sweeney L. k-anonymity: a model for protecting privacy. *Int. J. Uncertain. Fuzz. Knowl. Based Syst.* 2002;10(5):557–570.
- 4] Iwuchukwu T., Naughton J. K-anonymization as spatial indexing: toward scalable and incremental anonymization. In: Koch C., ed. *Proceedings of the 33rd International Conference on Very Large Data Bases*, September 23–27, 2007, University of Vienna, Austria; ACM; 2007:746–757.
- 5] Ghinita G., Karras P., Kalnis P., Mamoulis N. *Fast data anonymization with low information loss*. In: Koch C., ed. *Proceedings of the 33rd International Conference on Very Large Data Bases*, 23-27 September 2007, University of Vienna, Austria; ACM; 2007:758–769.
- 6] LeFevre K., DeWitt D., Ramakrishnan R. Workload-aware anonymization techniques for large-scale datasets. *ACM Trans. Database Syst.* 2008;33(3):1–47.
- 7] LeFevre K., DeWitt D.J., Ramakrishnan R. *Incognito: efficient full-domain k-anonymity*. In: Özcan F., ed. *Proceedings of the ACM SIGMOD International Conference on Management of Data*, June 14–16, 2005, Baltimore, MD; ACM; 2005:49–60.
- 8] El Emam K., Dankar F.K., Issa R., Jonker E., Amyot D., Cogo E., Corriveau J.P., Walker M., Chowdhury S., Vaillancourt R., Roffey T., Bottomley J. Research paper: a globally optimal k-anonymity method for the de-identification of health data. *JAMIA*. 2009;16(5):670–682.
- 9] Jian Xu, Wei Wang, Jian Pei, Xiaoyuan Wang, Baile Shi, Ada Wai-Chee Fu. Utility-based anonymization using local recoding
- 10] A. Gionis and T. Tassa, “k-Anonymization with minimal loss of information,” *IEEE Trans. Knowl. Data Eng.*, vol. 21, no. 2, 2009, pp. 206–219.
- 11] LeFevre K., DeWitt D., Ramakrishnan. Mondrian Multidimensional K-Anonymity
- 12] Tang, Qingming, et al. "Improving Strict Partition for Privacy Preserving Data Publishing." *Networking and Distributed Computing (ICNDC)*, 2010 First International Conference on. IEEE, 2010.
- 13] Divya Sadhwani, Dr. Sanjay Silakari, Mr. Uday Chourasia. Preserving Privacy during Big Data Publishing using K-Anonymity Model – A Survey
- 14] Janosch Maier . Anonymity: Formalization of Privacy – k-anonymity
Wong W.K., Mamoulis N., Cheung D.W.L. *Non-homogeneous generalization in privacy preserving data publishing*. In: Elmagarmid A., Agrawal D., eds. *Proceedings of the ACM SIGMOD International Conference on Management of Data (SIGMOD 2010)*, June 6–10, Indianapolis, IN; ACM; 2010:747–758.

- 15]Charu C. Aggarwal, Philip S. Yu . Privacy-Preserving Data Mining
- 16]Ashwin Machanavajjhala, Johannes Gehrke, Daniel Kifer, Muthuramakrishnan Venkatasubramanian . ℓ -Diversity: Privacy Beyond k -Anonymity
- 17]Ninghui Li, Tiancheng Li, and Suresh Venkatasubramanian (2007). "t-Closeness: Privacy beyond k -anonymity and l -diversity
- 18]Yu Zong, Guandong Xu . Applied Data Mining by Zhenglu Yang
- 19] Iyengar, V.: Transforming data to satisfy privacy constraints. In: Proceedings of International Conference on Knowledge Discovery and Data Mining, pp. 279–288 (2002)
- 20] Dankar, F., Emam, K.E., Neisa, A., Roffey, T.: Estimating the re-identification risk of clinical data sets. *BMC Med. Inform. Decis. Mak.* 12(1), 66 (2012)
- 21] Fabian Prasser, Florian Kohlmayer, Putting Statistical Disclosure Control into Practice: The ARX Data Anonymization Tool
- Aggarwal G., Panigrahy R., Feder T., Thomas D., Kenthapadi K., Khuller S., Zhu A. Achieving anonymity via clustering. *ACM Trans. Algor.* 2010;6(3):1–19.
- Li N., Qardaji W.H., Su D. Provably private data anonymization: or, k -anonymity meets differential privacy. *CoRR*. 2011 Abs/1101.2604.
- A Framework for Efficient Data Anonymization under Privacy and Accuracy Constraints
GABRIEL GHINITA, PANAGIOTIS KARRAS, and PANOS KALNIS National University of Singapore and NIKOS MAMOULIS University of Hong Kong, p9 (2009)
- Ciriani, V., De Capitani di Vimercati, S., Foresti, S., Samarati, P.: Microdata protection. In: Yu, T., Jajodia, S. (eds.) *Secure Data Management in Decentralized Systems. Advances in Information Security*, vol. 33, pp. 291–321. Springer, Berlin (2007) – μικροσυσσωμάτωση
- Nergiz, M.E., Atzori, M., Clifton, C.: Hiding the presence of individuals from shared databases. In: Proceedings of International Conference on Management of Data,(2007)