

ΠΑΝΕΠΙΣΤΗΜΙΟ ΚΡΗΤΗΣ
ΣΧΟΛΗ ΘΕΤΙΚΩΝ ΕΠΙΣΤΗΜΩΝ
ΤΜΗΜΑ ΕΠΙΣΤΗΜΗΣ ΥΠΟΛΟΓΙΣΤΩΝ

**Χρήση Γνώσης του Πεδίου Εφαρμογής για
την Ενίσχυση των Επαγωγικών
Μηχανισμών Μάθησης μέσω
Παραδειγμάτων**

Εργασία που υποβλήθηκε από την
Ελένη Γκάγκα
ως μερική απαίτηση για την απόκτηση του
ΜΕΤΑΠΤΥΧΙΑΚΟΥ ΔΙΠΛΩΜΑΤΟΣ ΕΞΕΙΔΙΚΕΥΣΗΣ

Χρήση Γνώσης του Πεδίου Εφαρμογής για την Ενίσχυση των Επαγωγικών Μηχανισμών Μάθησης μέσω Παραδειγμάτων

Εργασία που υποβλήθηκε από την
Ελένη Γκάγκα

ως μερική απαίτηση για την απόκτηση του
ΜΕΤΑΠΤΥΧΙΑΚΟΥ ΔΙΠΛΩΜΑΤΟΣ ΕΞΕΙΔΙΚΕΥΣΗΣ

18 Ιουνίου, 1996

Συγγραφέας:

Εισηγητική Επιτροπή:

Στέλιος Ορφανουδάκης, Καθηγητής, Επόπτης

Πάνος Κωνσταντόπουλος, Αναπληρωτής Καθηγητής, Μέλος

Γεώργιος Τζιρίτας, Αναπληρωτής Καθηγητής, Μέλος

Βασίλης Μουστάκης, Επίκουρος Καθηγητής, Εξωτερικό Μέλος

Δεκτή:

Πάνος Κωνσταντόπουλος, Αναπληρωτής Καθηγητής
Πρόεδρος Επιτροπής Μεταπτυχιακών Σπουδών

1 Ευχαριστίες

Η εργασία αυτή πραγματοποιήθηκε με την καθοδήγηση του καθηγητή του Τμήματος Επιστήμης Υπολογιστών Πανεπιστημίου Κρήτης, Στέλιου Ορφανουδάκη τον οποίο ευχαριστώ θερμά για την υποστήριξη και τις πολύτιμες συμβουλές του.

Ευχαριστώ επίσης τους Βασίλη Μουστάκη και Γιώργο Ποταριά για τις πολύτιμες συζητήσεις που είχαμε στο χώρο της Μηχανικής Μάθησης.

Ιδιαίτερος ευχαριστώ τον παιδοχειρουργό Γιώργο Χαρίση γιατί χωρίς τις επιστημονικές του γνώσεις αλλά και τη συνεργασία του θα ήταν αδύνατο να ολοκληρωθεί η εργασία αυτή. Παράλληλα, ευχαριστώ τα μέλη της ομάδας του κ. Χαρίση στην Παιδιατρική Κλινική του Πανεπιστημιακού Νοσοκομείου. Την ομάδα αποτελούν οι ειδικευόμενοι Γιαννόπουλος και Βλαχάκης και η κ. Παπαγεωργίου. Χωρίς την συνεργασία της ομάδας αυτής δεν θα ήταν δυνατή η αξιολόγηση της παρούσας εργασίας.

Τέλος, ευχαριστώ τον Αντώνη Αργυρό για τις πολύτιμες συμβουλές, την υποστήριξη αλλά και τον χρόνο που διέθεσε για την διαμόρφωση του τελικού κειμένου.

Περιεχόμενα

1 Ευχαριστίες	1
2 Εισαγωγή	1
3 Μάθηση μέσω Παραδειγμάτων	12
3.1 Ο αλγόριθμος ID3	12
3.2 Μάθηση μέσω Παραδειγμάτων και Εμπειρα Συστήματα	19
3.3 Ο αλγόριθμος NewId	20
3.3.1 Τύποι Ιδιοτήτων-Τιμών του NewId	20
3.3.2 Σχέση Διάταξης στον NewId	22
4 Ο αλγόριθμος IDDD	24
4.1 Απλή Εξάρτηση	32
4.2 Εξάρτηση Αποκλεισμού	34
4.3 Διαδοχικές Εξαρτήσεις	39
4.4 Πολυπλοκότητα του αλγορίθμου IDDD	40
5 Πειραματικά Αποτελέσματα	41
5.1 Περιγραφή Πεδίου Εφαρμογής	41
5.2 Προεπεξεργασία των Παραδειγμάτων Εκπαίδευσης	44
5.3 Σχεδιασμός Πειραμάτων	47
5.4 Αξιολόγηση Αποτελεσμάτων	52
5.4.1 Ποσοτική Αξιολόγηση	52

5.4.2	Ποιοτική Αξιολόγηση	54
6	Συμπεράσματα και Δυνατότητες Μελλοντικής Έρευνας	59
6.1	Συμπεράσματα	59
6.2	Δυνατότητες Μελλοντικής Έρευνας	62

Κατάλογος Σχημάτων

1	Ένα απλό δέντρο απόφασης	17
2	Ένα πολύπλοκο δέντρο απόφασης	17
3	B εξαρτάται από την A : $> a_i$ και $a_j \leq a_i$	37
4	B εξαρτάται από την A : $> a_i$ και $a_j > a_i$	37
5	(a) Το δέντρο του NewId, (b) Το δέντρο του IDDD	61

Κατάλογος Πινάκων

1	Ένα μικρό σύνολο εκπαιδευτικών παραδειγμάτων	16
2	Μέσοι Οροι Αριθμού Κανόνων και Ακρίβειας Κατάταξης.	53
3	Σύγκριση μεταξύ Αρχικού και Διατεταγμένου Συνόλου Κανόνων.	56
4	Σύγκριση μεταξύ Αρχικού και Διάταξης-Εξάρτησης Συνόλου Κανόνων.	57
5	Σύγκριση μεταξύ Αρχικού και Εξάρτησης Συνόλου Κανόνων.	57
6	Σύγκριση μεταξύ Διατεταγμένου και Διάταξης-Εξάρτησης Συνόλου Κανόνων.	58

Abstract

Artificial intelligence (AI) is now experiencing extraordinary growth, and applications of its ideas and methods are appearing in many fields. Among its most visible and important successes are the development of expert systems. In this context, it is important to ask what the limitations of the current methods are and what new directions research in this field should take. One of the obvious limitations relates to *machine learning*.

Current AI systems have very limited learning abilities or none at all. All of their knowledge must be programmed into them. When they contain an error, they cannot correct it on their own. They will repeat it endlessly, no matter how many times the procedure is executed. Generally speaking, these systems lack the ability to draw inductive inferences from information given to them.

The ability to classify objects or events as members of known classes is a very common task for learning systems. A well known approach to heuristic classification is decision tree induction. An alternative to decision tree induction is given by the well known algorithm ID3. ID3 uses a heuristic search process to find a set of discriminant descriptions between classes, given: (1) A set of observational statements each of which is assigned to a certain class and (2) a universe of classes.

Working in a recursive manner, the algorithm selects the most discriminant attribute by maximizing an information gain function at each step. The result is a tree in which nodes represent tests on attributes, while branches are possible values of the corresponding attributes.

ID3 does not take into account any background information resulting to a set of rules that are far from the expert's model. In this work, we present IDDD, which extends ID3 by using dependency relationships, between attributes and/or attribute value sets, as domain knowledge agents. IDDD is based on **NewId**, an enhanced implementation of ID3, developed by the Turing Institute. The basic premise of IDDD lies in the deployment of domain knowledge in the decision tree induction process. We introduce dependency relationships between attributes that are provide some structure over the rather 'flat' data representation used by NewId itself. The attribute that depends on another attribute is called Daughter while the latter attribute is called Mother.

A **simple** dependency relationship states that information represented by attribute **Daughter** should be useful only when it is combined with information represented by attribute **Mother**. Another form of the relation defined above is the exclusion of an attribute when a specific value has already been assigned to another attribute. We call this dependency **exclusive**. We should note here that any **Daughter** attribute may have only one **mother**, while this restriction does not apply to a **mother** attribute.

More often than not, **NewId** is able to ‘handle’ relationships such as the ones defined above, yet this is achieved in an implicit manner and it is not based on explicit modeling of domain knowledge. However, when the value of a **mother** attribute is unknown, **NewId** may drift into irrelevant attribute selection and splitting. Furthermore, when two attributes **A** and **B** score equally in information gain, **NewId** may select **B** instead of **A**. The outcome is a set of rules which may be accurate, from the point of view of classification accuracy but meaningless from the expert’s point of view.

The effectiveness of **IDDD** is demonstrated through carefully designed experiments involving a medical application. The effectiveness of dependency relationships defined in **IDDD** indicates that further work is necessary towards the investigation, formal definition, and handling of domain knowledge in the inductive process.

Περίληψη

Η ανάπτυξη εμπειρών συστημάτων αποτελεί μία από τις μεγαλύτερες επιτυχίες του χώρου της Τεχνητής Νοημοσύνης αλλά γεννά και πολλά ερωτηματικά για τα όρια των σημερινών τεχνικών. Τα σημερινά συστήματα τεχνητής νοημοσύνης έχουν ελάχιστες ή και καθόλου ικανότητες μάθησης. Λειτουργούν με βάση την γνώση που έχουν από την αρχή χωρίς δυνατότητες αναπροσαρμογής. Ο κλάδος της Μηχανικής Μάθησης αναπτύχθηκε για την αντιμετώπιση αυτού του προβλήματος. Ο στόχος είναι η ανάπτυξη αποδοτικών συστημάτων μάθησης τα οποία θα αποτελέσουν κομμάτια των εμπειρών συστημάτων και θα επιτρέψουν την αναπροσαρμογή της γνώσης που αυτά έχουν όταν αντιμετωπίζονται άγνωστες ή ιδιόμορφες καταστάσεις.

Ένα πολύ συνηθισμένο πρόβλημα για τα συστήματα μάθησης είναι η εύρεση κανόνων κατάταξης παραδειγμάτων σε κλάσεις. Μία πολύ γνωστή προσέγγιση σε αυτό το πρόβλημα είναι η επαγωγική εύρεση ενός δέντρου απόφασης. Αυτή την προσέγγιση ακολουθεί ο γνωστός, από την βιβλιογραφία, αλγόριθμος ID3.

Πρόκειται για έναν αναδρομικό αλγόριθμο ο οποίος χρησιμοποιεί μία ευριστική διαδικασία εύρεσης των πιο διαχωριστικών περιγραφών μεταξύ κλάσεων. Τα δεδομένα είναι: (1) Ένα σύνολο παραδειγμάτων με την πληροφορία της κλάσης στην οποία ανήκουν και (2) ένα σύνολο κλάσεων.

Σε κάθε βήμα επιλέγεται η πιο πληροφοριακή ιδιότητα μεγιστοποιώντας μία συνάρτηση πληροφορίας. Το αποτέλεσμα είναι ένα δέντρο απόφασης στο οποίο κάθε εσωτερικός κόμβος αναπαριστά έλεγχο της τιμής μιας ιδιότητας και τα κλαδιά είναι οι δυνατές τιμές των αντιστοιχών ιδιοτήτων. Κάθε φύλλο αντιστοιχεί σε μία κλάση και κάθε μονοπάτι από τη ρίζα μέχρι κάποιο φύλλο αποτελεί έναν διαφορετικό κανόνα κατάταξης παραδειγμάτων στην αντιστοιχη κλάση.

Ο αλγόριθμος δεν λαμβάνει υπόψην του ιδιομορφίες του πεδίου εφαρμογής με αποτέλεσμα να κατασκευάζονται κανόνες που αγνοούν στοιχεία που θεωρούνται σημαντικά από τους ειδικούς. Στα πλαίσια της εργασίας αυτής υλοποιήθηκε η ενίσχυση του μηχανισμού κατασκευής του δέντρου με την εισαγωγή γνώσης που αφορά την συγκεκριμένη εφαρμογή. Αποτέλεσμα της προσπάθειας αυτής ήταν η ανάπτυξη του αλγορίθμου IDDD.

Το σύστημα IDDD βασίστηκε στο σύστημα NewId που είναι μία βελτιωμένη

έκδοση του ID3. Η γνώση που χρησιμοποιεί ο IDDD παίρνει την ειδική μορφή των σχέσεων εξάρτησης μεταξύ των ιδιοτήτων που περιγράφουν την εφαρμογή. Ο στόχος είναι η εξαγωγή κανόνων που υπόκεινται στους ίδιους περιορισμούς που λαμβάνουν υπόψη οι ειδικοί ακόμα και στις περιπτώσεις που οι περιορισμοί αυτοί δεν αντανakλώνται στα δεδομένα. Στη σχέση εξάρτησης εισάγουμε τον όρο Κόρη για την ιδιότητα που εξαρτάται από κάποια άλλη και τον όρο Μάνα γι'αυτή την δεύτερη ιδιότητα.

Συγκεκριμένα ορίζεται η σχέση της απλής εξάρτησης ιδιοτήτων σύμφωνα με την οποία η επιλογή μιας ιδιότητας στις συνθήκες ενός κανόνα εξαρτάται από την επιλογή κάποιας άλλης ιδιότητας σε προηγούμενο βήμα. Επιπλέον, ορίζεται η σχέση της αποκλειστικής εξάρτησης η οποία καθορίζει μία συγκεκριμένη τιμή της Μάνας για την οποία απαγορεύεται η επιλογή της Κόρης. Για τον IDDD απαγορεύεται η εξάρτηση μιας ιδιότητας από περισσότερες από μία ιδιότητες ενώ ο περιορισμός αυτός δεν ισχύει για τις ιδιότητες Μάνες.

Μερικές φορές, ο NewId χειρίζεται ικανοποιητικά τέτοιες σχέσεις. Η διαφορά είναι ότι ο μηχανισμός που χρησιμοποιεί δεν σχετίζεται με τον τυπικό ορισμό αυτών των σχέσεων αλλά απαιτεί από τον χρήστη πολύ καλή γνώση του αλγορίθμου και μεταβολές στην αναπαράσταση των δεδομένων. Ακόμα και με τις προϋποθέσεις αυτές, είναι αρκετά συχνές οι περιπτώσεις αποτυχίας σωστού χειρισμού τέτοιων σχέσεων. Όταν η τιμή της Μάνας είναι άγνωστη, ο αλγόριθμος μπορεί να οδηγηθεί στη λανθασμένη επιλογή της Κόρης. Επιπλέον όταν η Μάνα και η Κόρη έχουν την ίδια πληροφοριακή ισχύ ο NewId μπορεί να επιλέξει την Κόρη αντί της Μάνας. Το αποτέλεσμα είναι ένα σύνολο κανόνων οι οποίοι μπορεί να επιτυχάνουν μεγάλη ακρίβεια στην κατάταξη παραδειγμάτων αλλά δεν έχουν νόημα για τους ειδικούς της εφαρμογής.

Στον IDDD οι σχέσεις αυτές χρησιμοποιούνται σαν επιπλέον κριτήρια (περιορισμοί) για την επιλογή μιας ιδιότητας. Οι ιδιότητες που δεν μετέχουν σε καμία σχέση εξάρτησης δεν επηρεάζονται αλλά περιορίζεται η επιλογή ιδιοτήτων που είναι Κόρες. Ο ορισμός των σχέσεων εξάρτησης εξαρτάται και από τον τύπο των ιδιοτήτων. Διακρίνουμε μεταξύ κατηγορηματικών και αριθμητικών ιδιοτήτων. Οι διαφορετικοί τύποι υπαγορεύουν ιδιαιτερότητες στην υλοποίηση των σχέσεων οι οποίες ελήφθησαν υπόψη κατά την ανάπτυξη του IDDD.

Για την αξιολόγηση της απόδοσης του IDDD χρησιμοποιήθηκε μία συγκεκριμένη ιατρική εφαρμογή, με βάση την οποία έγινε η συγκριτική μελέτη των IDDD και NewId. Τα πειραματικά αποτελέσματα δείχνουν τη χρησιμότητα της προσέγγισης, οδηγούν σε γενικότερα συμπεράσματα για την χρήση γνώσης του πεδίου εφαρμογής σαν μέσο ενίσχυσης των επαγωγικών μηχανισμών μάθησης μέσω παραδειγμάτων και θέτουν σαφείς ερευνητικές κατευθύνσεις.

2 Εισαγωγή

Η Τεχνητή Νοημοσύνη είναι ένας από τους κλάδους της επιστήμης των υπολογιστών που έχει γνωρίσει, ιδιαίτερα τα τελευταία χρόνια, αλματώδη ανάπτυξη και προσελκύει το ενδιαφέρον ολοένα και περισσότερων ερευνητών. Από τα μεγαλύτερα επιτεύγματα του χώρου είναι η ανάπτυξη έμπειρων συστημάτων, η σημαντική πρόοδος στον τομέα της Μηχανικής Ορασης και στην αναγνώριση φωνής.

Η ανάπτυξη αυτή δημιουργεί βάσιμες ελπίδες για παραπέρα επιτυχίες αλλά γεννά και πολλά ερωτηματικά για τα όρια του χώρου. Τα σημερινά συστήματα τεχνητής νοημοσύνης έχουν ελάχιστες, ή και καθόλου, ικανότητες μάθησης. Ο προγραμματιστής ενός τέτοιου συστήματος πρέπει να το τροφοδοτήσει εκ των προτέρων με την αναγκαία γνώση. Αν η γνώση αυτή περιέχει κάποιο λάθος, το σύστημα αδυνατεί να το διορθώσει και επαναλαμβάνει επ'άπειρο την λανθασμένη λειτουργία του. Γενικά, θα μπορούσαμε να πούμε πως αυτά τα συστήματα είναι ικανά να εξάγουν μόνο συμπεράσματα με βάση την γνώση που έχουν ενώ τους λείπουν μηχανισμοί εξαγωγής επαγωγικών συμπερασμάτων. Σε αντίθεση, ο άνθρωπος είναι ικανός να διευρύνει τις γνώσεις του και να αποκτά καινούργιες ικανότητες τις οποίες μάλιστα βελτιώνει με την πρακτική εξάσκηση.

Ο μεγάλος ρόλος της μάθησης στην ανθρώπινη νοημοσύνη μας οδηγεί στο λογικό συμπέρασμα πως ανάλογοι μηχανισμοί πρέπει να αναπτυχθούν και για τα έμπειρα συστήματα [1, 2, 3]. Αυτός είναι και ο λόγος της δημιουργίας του κλάδου της **Μηχανικής Μάθησης**, ο οποίος μελετά τη φύση της ανθρώπινης μάθησης και προσπαθεί να υλοποιήσει τους μηχανισμούς της στα υπολογιστικά συστήματα. Σαν επιστήμη η Μηχανική Μάθηση έχει τρεις βασικές ερευνητικές κατευθύνσεις [4]:

1. **Προσανατολισμός ως προς το στόχο** (Task-Oriented Studies) που αφορά στην ανάπτυξη συστημάτων για την επίλυση συγκεκριμένων προβλημάτων.
2. **Προσομοίωση Γνωστικών Λειτουργιών** (Cognitive Simulation), που αφορά στην έρευνα και προσομοίωση των μηχανισμών του ανθρώπινου τρόπου μάθησης.
3. **Θεωρητική ανάλυση** (Theoretical Analysis), που αφορά στην έρευνα των δυ-

νατών μηχανισμών μάθησης ανεξαρτήτως πεδίου εφαρμογής.

Φυσικά η εξέλιξη σε οποιοδήποτε από τους παραπάνω τομείς επιδρά στην εξέλιξη και των υπολοίπων.

Μελετώντας τον μηχανισμό της ανθρώπινης μάθησης μπορεί κανείς να διακρίνει δύο βασικές μορφές του: την **απόκτηση γνώσης** (knowledge acquisition) και την **αναπροσαρμογή χειρισμών** σε συγκεκριμένες καταστάσεις (skill refinement). Η απόκτηση γνώσης μπορεί να οριστεί σαν απόκτηση καινούργιας πληροφορίας που συνδυάζεται με ικανότητα αποτελεσματικής εφαρμογής της στην αντιμετώπιση προβλημάτων. Σε αντίθεση με την απόκτηση γνώσης, η οποία μπορεί να είναι υποσυνείδητη, η αναπροσαρμογή χειρισμών είναι η συνειδητή προσπάθεια του ανθρώπου να βελτιώσει τις ικανότητές του, στο χειρισμό κάποιων καταστάσεων, με βάση την εξάσκηση. Η ανθρώπινη νοημοσύνη αποτελεί ένα συνδυασμό και των δύο αυτών μορφών μάθησης. Η γνώση υποδεικνύει στον άνθρωπο κάποιους χειρισμούς οι οποίοι αναπροσαρμόζονται αν αποδειχτούν αναποτελεσματικοί και έτσι οδηγείται σε απόκτηση νέας γνώσης.

Ο σημαντικότερος παράγοντας για το σχεδιασμό συστημάτων Μηχανικής Μάθησης αφορά στο είδος της πληροφορίας που παρέχεται από το περιβάλλον και ειδικότερα στο **επίπεδο** και στην **ποιότητα** αυτής της πληροφορίας. Το επίπεδο της πληροφορίας αναφέρεται στο βαθμό γενικότητάς της. Πληροφορία που σχετίζεται με μια γενικότερη κλάση προβλημάτων είναι πληροφορία υψηλού επιπέδου ενώ η χαμηλού επιπέδου πληροφορία αφορά στο συγκεκριμένο προς επίλυση πρόβλημα. Το σύστημα πρέπει είτε να εξειδικεύσει την πληροφορία υψηλού επιπέδου είτε να γενικεύσει την πληροφορία χαμηλού επιπέδου. Για το σκοπό αυτό πρέπει να διαθέτει αντίστοιχα τους κατάλληλους μηχανισμούς **εξειδίκευσης** (specialization) ή **γενίκευσης** (generalization). Η ελλιπής πληροφορία απαιτεί **δημιουργία υποθέσεων** (generation of hypotheses). Οι υποθέσεις αυτές πρέπει να αξιολογηθούν με βάση την αποτελεσματικότητά τους και να αναπροσαρμοστούν, αν αυτό είναι απαραίτητο.

Με βάση τις παραπάνω παρατηρήσεις, ο χώρος της Μηχανικής Μάθησης διαφεύγει σε δύο βασικούς τομείς με γνώμονα τον τρόπο δημιουργίας των υποθέσεων. Έτσι δια-

κρίνεται ο χώρος της **Συμπερασματικής** (Deductive) και ο χώρος της **Επαγωγικής** (Inductive) Μάθησης [5, 6, 7].

Με τον όρο Συμπερασματική Μάθηση εννοούμε τη διαδικασία δημιουργίας υποθέσεων με την εφαρμογή συμπερασματικής λογικής στην πληροφορία που έχει το σύστημα.

Περισσότερο ενδιαφέρον παρουσιάζει ο χώρος της Επαγωγικής Μάθησης δηλαδή της διαδικασίας που επιτρέπει σε ένα σύστημα να δημιουργεί, με βάση τη γνώση που έχει, καινούργιες υποθέσεις για έννοιες (concepts), εξαιρέσεις, επεξηγήσεις ή υψηλού επιπέδου χαρακτηρισμούς της πληροφορίας εισόδου.

Το επίπεδο της πληροφορίας καθορίζει το είδος των υποθέσεων που το σύστημα καλείται να δημιουργήσει. Διακρίνουμε τέσσερα τέτοια είδη υποθέσεων, καθένα από τα οποία ορίζει ένα συγκεκριμένο κλάδο του γενικότερου χώρου της επαγωγικής μάθησης:

1. **Μάθηση χωρίς υποθέσεις** (Rote Learning). Το σύστημα διαθέτει ακριβώς την πληροφορία που του χρειάζεται και δεν είναι αναγκαία η δημιουργία υποθέσεων.
2. **Μάθηση μέσω οδηγιών** (Learning by being told). Η πληροφορία είναι γενική και το σύστημα πρέπει να δημιουργήσει υποθέσεις για την πληροφορία που λείπει.
3. **Μάθηση μέσω παραδειγμάτων** (Learning from examples). Η πληροφορία είναι πολύ εξειδικευμένη και το σύστημα πρέπει να δημιουργήσει υποθέσεις για πιο γενικούς κανόνες.
4. **Μάθηση μέσω αναλογιών** (Learning by analogy). Η πληροφορία σχετίζεται με ανάλογα προβλήματα οπότε το σύστημα πρέπει να εντοπίσει τις αναλογίες και να δημιουργήσει ανάλογες υποθέσεις για το τρέχον πρόβλημα.

Όλες οι παραπάνω κατηγορίες μάθησης αφορούν στην απόκτηση γνώσης και όχι στην αναπροσαρμογή χειρισμών. Παρατίθενται μάλιστα με αυξανόμενο βαθμό πολυπλοκότητας των μηχανισμών μάθησης.

Η ποιότητα της πληροφορίας επιδρά σημαντικά στο βαθμό δυσκολίας του προβλήματος μάθησης. Για παράδειγμα όταν η μάθηση γίνεται μέσω παραδειγμάτων ο βαθμός

δυοκοιλίας μειώνεται όταν τα παραδείγματα είναι προσεκτικά επιλεγμένα, ώστε να αντιπροσωπεύουν το πεδίο εφαρμογής, και δεν περιέχουν λάθη.

Τα συστήματα μάθησης που έχουν δημιουργηθεί ως τώρα κατατάσσονται αυστηρά σε μία μόνο από τις κατηγορίες επαγωγικής μάθησης. Τελευταία, η ερευνητική προσπάθεια έχει στραφεί στον συνδυασμό περισσότερων τεχνικών για την ανάπτυξη συστημάτων που θα επεξεργάζονται πληροφορίες διαφορετικών επιπέδων αλλά και θα προσαρμόζονται στις ανάγκες και ιδιαιτερότητες του εκάστοτε πεδίου εφαρμογής [8, 9, 10]. Οι προσπάθειες αυτές αν και βρίσκονται ακόμα σε εμβρυική μορφή έχουν ελπιδοφόρα αποτελέσματα.

Στα πλαίσια της εργασίας αυτής, θα περιοριστούμε στην ειδική κατηγορία της επαγωγικής μάθησης μέσω παραδειγμάτων. Στο χώρο της Μηχανικής Μάθησης ο κλάδος αυτός προσέλκυσε πλήθωρα ερευνητών τόσο για την απλότητά του αλλά και για την ποιότητα των αποτελεσμάτων του. Είναι γνωστό, από το χώρο της Τεχνητής Νοημοσύνης, πως το μεγαλύτερο πρόβλημα στην ανάπτυξη έμπειρων συστημάτων αποτελεί η εξαγωγή και η αναπαράσταση της γνώσης του ειδικού στον τομέα εφαρμογής. Ενώ άλλες μέθοδοι Μηχανικής Μάθησης κάνουν χρήση της γνώσης αυτής και αναζητούν μεθόδους βέλτιστης αναπαράστασης και χρησιμοποίησή της, η μάθηση μέσω παραδειγμάτων στοχεύει να εξαλείψει αυτή την ανάγκη. Έτσι τα συστήματα της κατηγορίας αυτής είναι εύκολα στη χρήση και δεν απαιτούν γνώση των μηχανισμών μάθησης. Ωστόσο υπολείπονται, όσον αφορά το τελικό αποτέλεσμα, των συστημάτων που κάνουν χρήση γνώσης για το πεδίο εφαρμογής. Κι αυτό γιατί αναπτύχθηκαν για να χρησιμοποιηθούν σε οποιοδήποτε πρόβλημα χωρίς να εκμεταλλεύονται τυχόν ιδιαιτερότητές του. Η εργασία αυτή επικεντρώθηκε στην εισαγωγή γνώσης σε συστήματα μάθησης μέσω παραδειγμάτων προσπαθώντας να διατηρήσει την απλότητα του συστήματος. Στόχος ήταν η εύρεση ενός απλού τρόπου εισαγωγής γνώσης έτσι ώστε να μην επιβαρύνεται ο χρήστης, να λαμβάνεται υπόψη η φύση της εφαρμογής αλλά και να διατηρηθεί η γενικότητα του συστήματος ώστε να μην περιοριστούν οι δυνατότητές του.

Στο σύστημα Μηχανικής Μάθησης παρέχεται χαμηλού επιπέδου πληροφορία με την μορφή παραδειγμάτων που περιγράφουν συγκεκριμένες καταστάσεις του πεδίου εφαρ-

μογής. Στα περισσότερα ίσως συστήματα τα παραδείγματα παρέχονται με την μορφή διανυσμάτων ιδιοτήτων (feature vectors). Ειδικότερα, για κάθε πεδίο εφαρμογής επιλέγονται κάποιες ιδιότητες (χαρακτηριστικά) που το περιγράφουν. Για παράδειγμα στην περιοχή της διαγνωστικής ιατρικής οι ιδιότητες αυτές είναι το σύνολο των κλινικών και εργαστηριακών εξετάσεων που απαιτούνται προκειμένου να καταλήξει ο γιατρός στη διάγνωση. Κάθε συγκεκριμένο παράδειγμα έχει μια συγκεκριμένη τιμή για κάθε ιδιότητα. Έτσι αναπαριστάται σαν ένα διάνυσμα τιμών ιδιοτήτων (attribute-value vectors). Το σύστημα καλείται είτε να **ομαδοποιήσει σε κλάσεις** τα παραδείγματα (clustering) δημιουργώντας και τις αντίστοιχες περιγραφές των κλάσεων είτε, αν τα παραδείγματα είναι ήδη ταξινομημένα σε κλάσεις, να εξάγει κανόνες **κατάταξης παραδειγμάτων** (classification). Στη δεύτερη περίπτωση, στο σύστημα παρέχονται όχι μόνο παραδείγματα αλλά η απόφαση που παίρνεται στη συγκεκριμένη κατάσταση που το παράδειγμα αναπαριστά. Κάθε τέτοια απόφαση ονομάζεται **κλάση** και το σύνολο των δυνατών αποφάσεων είναι οι κλάσεις του πεδίου εφαρμογής. Η πληροφορία που παρέχεται από τα παραδείγματα και τις αντίστοιχες κλάσεις τους γενικεύεται με αποτέλεσμα τη δημιουργία ενός συνόλου γενικών κανόνων συμπεριφοράς/απόφασης.

Ένας πολύ γνωστός αλγόριθμος κατάταξης που περιγράφεται στη βιβλιογραφία και έχει γνωρίσει ευρεία αποδοχή για την απλότητα και την αποτελεσματικότητά του είναι ο ID3 [11, 12]. Ο ID3 χρησιμοποιεί έναν ευριστικό μηχανισμό αναζήτησης ενός συνόλου κανόνων που διαχωρίζουν μεταξύ των κλάσεων. Είσοδος του συστήματος είναι η **περιγραφή του πεδίου εφαρμογής** και ένα σύνολο παραδειγμάτων. Με τον όρο περιγραφή του πεδίου εννοούμε μία λεπτομερή καταγραφή των ιδιοτήτων που χαρακτηρίζουν/περιγράφουν το πεδίο καθώς και αναλυτική παράθεση των τιμών που παίρνει καθεμιά από αυτές. Το σύνολο των παραδειγμάτων εξυπηρετεί εκπαιδευτικούς σκοπούς και χρησιμοποιείται, μαζί με την περιγραφή του πεδίου, για την εξαγωγή των κανόνων κατάταξης των παραδειγμάτων. Κάθε παράδειγμα είναι ένα διάνυσμα τιμών των ιδιοτήτων του πεδίου. Το διάνυσμα αυτό πρέπει να περιλαμβάνει και την πληροφορία για την συγκεκριμένη κλάση στην οποία αυτό ανήκει. Αν για ένα παράδειγμα η τιμή μιας ιδιότητας δεν έχει καταγραφεί, ο αλγόριθμος παρέχει ειδικό σύμβολο για την αναπαράσταση της άγνωστης τιμής. Το σύνολο των κανόνων που παράγεται παίρνει την μορφή

ενός δέντρου απόφασης. Οι κόμβοι του δέντρου αντιστοιχούν σε ελέγχους των τιμών κάποιων ιδιοτήτων ενώ τα κλαδιά είναι οι δυνατές τιμές των αντιστοιχών ιδιοτήτων. Κάθε φύλλο αναπαριστά μία συγκεκριμένη κλάση. Κάθε μονοπάτι του δέντρου από τη ρίζα του μέχρι κάποιο συγκεκριμένο φύλλο αποτελεί έναν διαφορετικό κανόνα κατάταξης παραδειγμάτων σε μία από τις κλάσεις.

Ο ID3 είναι ένας αναδρομικός αλγόριθμος, σε κάθε βήμα του οποίου επιλέγεται η ιδιότητα που διαχωρίζει καλύτερα μεταξύ των παραδειγμάτων διαφορετικών κλάσεων. Ο τρόπος επιλογής της ιδιότητας αυτής σχετίζεται με τη μεγιστοποίηση μιας συνάρτησης υπολογισμού της πληροφοριακής ισχύος της ιδιότητας. Πρακτικά, μία ιδιότητα έχει μεγάλη πληροφοριακή ισχύ αν οι τιμές της στο σύνολο των παραδειγμάτων είναι τέτοιες ώστε να επιτυγχάνεται όσο το δυνατόν μεγαλύτερος διαχωρισμός μεταξύ των παραδειγμάτων που ανήκουν σε διαφορετική κλάση. Η ιδιότητα με τη μεγαλύτερη πληροφοριακή ισχύ γίνεται κόμβος του δέντρου και δημιουργούνται τόσα κλαδιά όσες και οι τιμές της. Το σύνολο των παραδειγμάτων χωρίζεται σε υποσύνολα έτσι ώστε κάθε υποσύνολο να περιέχει εκείνα τα παραδείγματα που έχουν μία συγκεκριμένη τιμή στην ιδιότητα που επιλέχτηκε. Κάθε τέτοιο υποσύνολο αποτελεί το σύνολο παραδειγμάτων για το αντίστοιχο κλαδί. Ο αλγόριθμος προσπαθεί σε κάθε βήμα να καλύψει τα παραδείγματα, δηλαδή να δημιουργήσει υποσύνολα του αρχικού συνόλου παραδειγμάτων έτσι ώστε τα παραδείγματα κάθε υποσυνόλου να ανήκουν σε μία και μόνο κλάση.

Ο ID3 χρειάζεται σαν είσοδο την περιγραφή του πεδίου και μερικά εκπαιδευτικά παραδείγματα. Δεν δέχεται οποιαδήποτε μορφή γνώσης που αφορά στο πεδίο εφαρμογής (Domain Knowledge). Το χαρακτηριστικό αυτό αποτελεί πλεονέκτημα στις περιπτώσεις όπου δεν υπάρχει κάποιος ειδικός να παρέχει τη γνώση αυτή. Όταν όμως η γνώση αυτή είναι διαθέσιμη, η αδυναμία χειρισμού και αξιοποίησής της από τον ID3 αποτελεί τον σημαντικότερο παράγοντα περιορισμού της αξίας των κανόνων.

Βασική αιτία γι' αυτό είναι ότι, όσον αφορά τον αλγόριθμο, οι ιδιότητες που περιγράφουν το πεδίο εφαρμογής περικλείουν πληροφορία αλλά δεν σχετίζονται μεταξύ τους με οποιοδήποτε τρόπο. Υποθέτει δηλαδή ανεξαρτησία των ιδιοτήτων. Κατά την δημιουργία του δέντρου εξετάζεται κάθε ιδιότητα και αποφασίζεται αν αυτή θα αποτελέσει κόμβο με

βάση μόνο τη πληροφοριακή της ισχύ. Η συνάρτηση υπολογισμού της πληροφοριακής ισχύος της ιδιότητας υπολογίζεται με βάση τα παραδείγματα. Τα παραδείγματα αυτά περιέχουν πολλές φορές ‘θόρυβο’ που προέρχεται είτε από λάθος μέτρηση της τιμής κάποιας ιδιότητας είτε από την ύπαρξη παραγόντων που επηρεάζουν το πεδίο εφαρμογής αλλά δεν καταγράφονται [13, 14]. Ο ID3 αντιμετωπίζει περιπτώσεις ‘θορύβου’ λόγω λανθασμένης μέτρησης αλλά δεν προβλέπει τη περίπτωση του δεύτερου είδους ‘θορύβου’. Στους προαναφερθέντες παράγοντες περιλαμβάνονται και περιπτώσεις όπου οι ιδιότητες δεν είναι ανεξάρτητες μεταξύ τους.

Δυστυχώς στις περισσότερες πραγματικές εφαρμογές η υπόθεση για ανεξαρτησία μεταξύ ιδιοτήτων δεν ευσταθεί. Συνήθως, η τιμή μιας ιδιότητας όχι μόνο επηρεάζει την τιμή κάποιας άλλης αλλά μπορεί να υποδεικνύει τη μη χρησιμότητα της άλλης. Όλα αυτά τα χαρακτηριστικά της εφαρμογής, που αγνοούνται από τον ID3, αποτρέπουν την παραγωγή καλών κανόνων από το σύστημα μάθησης. Αν το σύστημα διέθετε χειρισμούς αξιοποίησης της γνώσης αυτής θα μπορούσε να απλοποιηθεί το πρόβλημα μάθησης αλλά και να βελτιωθεί το τελικό αποτέλεσμα.

Στη βιβλιογραφία υπάρχουν αρκετές προτάσεις αντιμετώπισης του συγκεκριμένου προβλήματος της μελέτης αλλά και του χειρισμού αυτών των εξαρτήσεων [15, 16, 17, 18, 19, 20, 21]. Οι προτάσεις αυτές είτε είναι αναποτελεσματικές όπως η [15], της οποίας η γενική ιδέα μελετήθηκε και σχολιάζεται στο κεφάλαιο 4, είτε σχετίζονται με κάποιους άλλους αλγόριθμους και δεν είναι προφανής ο τρόπος προσαρμογής τους στην λειτουργία του ID3.

Για τη μελέτη του γενικού αλγόριθμου ID3 χρησιμοποιήθηκε μία συγκεκριμένη υλοποίησή του από το Turing Institute της Γλασκώβης, με την ονομασία NewId. Στη συγκεκριμένη υλοποίηση ο αλγόριθμος έχει ενισχυθεί με κάποιο μηχανισμό αντιμετώπισης μίας ειδικής σχέσης εξάρτησης ιδιοτήτων. Ο μηχανισμός αυτός συνίσταται στον ορισμό μιας διάταξης κατά την επιλογή των ιδιοτήτων. Βοηθητικά, έχουν αναπτυχθεί μηχανισμοί για τον χειρισμό περιπτώσεων αγνώστων τιμών. Πολλές φορές, κατά τη συλλογή δεδομένων, δεν είναι γνωστή η τιμή μιας ή περισσότερων ιδιοτήτων. Το γεγονός αυτό μπορεί να οφείλεται είτε στη μη ύπαρξη της πληροφορίας είτε στην αδυναμία καταγραφής της. Σε

άλλες περιπτώσεις η τιμή αυτή δεν καταγράφεται όχι γιατί δεν μπορούμε να τη μετρήσουμε αλλά γιατί δεν μας ενδιαφέρει για το συγκεκριμένο παράδειγμα. Η τελευταία αυτή περίπτωση συνήθως εκφράζει σχέσεις εξάρτησης μεταξύ τιμών ιδιοτήτων. Ο NewId διαθέτει μηχανισμούς αντιμετώπισης τέτοιων περιπτώσεων και έτσι θα περίμενε κανείς ότι έστω και με αυτό τον έμμεσο τρόπο χειρίζεται σχέσεις μεταξύ ιδιοτήτων. Συγκεκριμένα, εκτός από την άγνωστη τιμή, εισάγει την έννοια της αδιάφορης τιμής για παραδείγματα όπου δεν ενδιαφέρει η καταγραφή κάποιας ιδιότητας, και χειρίζεται με ειδικό τρόπο αυτές τις περιπτώσεις.

Παρά την ύπαρξη αυτών των μηχανισμών, το δέντρο απόφασης δεν παίρνει πάντα υπόψη του τις σχέσεις εξάρτησης μεταξύ των ιδιοτήτων. Το αποτέλεσμα εξαρτάται άμεσα από τα παραδείγματα εκπαίδευσης τα οποία δεν είναι πάντα εύκολο να ελέγξουμε για την ορθότητά τους. Για να γίνει πιο κατανοητή η παραπάνω παρατήρηση θεωρούμε την περίπτωση όπου κάποια συγκεκριμένη τιμή μιας ιδιότητας καθιστά απαγορευτική τη χρήση κάποιας άλλης ιδιότητας στον ίδιο κανόνα. Η πρώτη ιδιότητα ονομάζεται **Μάνα** και η δεύτερη **Κόρη**. Στα παραδείγματα όπου η Μάνα έχει την συγκεκριμένη απαγορευτική τιμή δίνουμε στην Κόρη την τιμή που υποδηλώνει στον αλγόριθμο ότι δεν θέλουμε να χρησιμοποιήσει την ιδιότητα Κόρη. Σύμφωνα με τα παραπάνω δεν θα έπρεπε να δημιουργείται πρόβλημα. Όμως, σύμφωνα με τα πειραματικά αποτελέσματα, ο αλγόριθμος αποτυγχάνει στις περιπτώσεις όπου η Μάνα έχει, για μερικά παραδείγματα, άγνωστη τιμή. Για τα παραδείγματα αυτά η Κόρη δεν θα έχει αδιάφορη τιμή. Το σύνολο των παραδειγμάτων αυτών απαιτεί κάλυψη και τίποτα δεν εμποδίζει τον αλγόριθμο να οδηγηθεί σε επιλογή της Κόρης άσχετα αν ενδεχομένως η επιλογή της Μάνας σε προηγούμενο βήμα θα έπρεπε να το απαγορεύει. Η επιλογή αυτή είναι η κυριότερη αιτία παραγωγής δέντρων τα οποία είναι πολύ αποτελεσματικά όσον αφορά την **ακρίβεια** κατάταξης των κανόνων (classification accuracy) αλλά δεν επιτυγχάνουν στην εξαγωγή **ποιοτικά** καλών κανόνων. Με τον όρο ακρίβεια κατάταξης αναφερόμαστε στο ποσοστό επιτυχίας των κανόνων να κατατάξουν τα παραδείγματα στις σωστές κλάσεις. Με τον όρο **ποιότητα** κανόνων αναφερόμαστε στο βαθμό αποδοχής των κανόνων από τους ειδικούς. Ειδικότερα όσον αφορά το δέντρο του NewId, επιτυγχάνει συνήθως αρκετά μεγάλη ακρίβεια κατάταξης αλλά οι ειδικοί στο εκάστοτε πεδίο εφαρμογής συχνά κρίνουν ότι

λαμβάνει υπόψιν στοιχεία που αυτοί δεν θα χρησιμοποιούσαν ή αγνοεί κάποια άλλα που θεωρούν σημαντικά.

Στα πλαίσια της εργασίας αυτής διαπιστώθηκε πως αιτία για την μη εξαγωγή καλών κανόνων αποτελεί, πολλές φορές, η αναποτελεσματικότητα των μηχανισμών χειρισμού εξαρτήσεων. Η παραπάνω διαπίστωση οδήγησε σε προσεκτικότερη μελέτη της φύσης των εξαρτήσεων μεταξύ ιδιοτήτων. Έτσι ορίστηκαν τυπικότερα οι σχέσεις εξάρτησης Μάνας-Κόρης και διαπιστώθηκε πως εμφανίζονται με δύο διαφορετικές μορφές, τις **απλές εξαρτήσεις** και τις **εξαρτήσεις αποκλεισμού**. Με τον όρο απλή εξάρτηση περιγράφεται η περίπτωση όπου η πληροφορία μιας συγκεκριμένης ιδιότητας δεν είναι χρήσιμη παρά μόνο όταν η ιδιότητα αυτή εμφανίζεται σε ένα κανόνα σε συνδυασμό με κάποια άλλη ιδιότητα. Με τον όρο εξάρτηση αποκλεισμού περιγράφεται η περίπτωση όπου η τιμή μιας ιδιότητας, όπως αυτή εμφανίζεται σε ένα κανόνα, απαγορεύει τη χρήση μιας άλλης ιδιότητας για τον κανόνα αυτό. Ο εντοπισμός των εξαρτήσεων αυτών απαιτεί την ανάπτυξη αντίστοιχων μηχανισμών για το χειρισμό τους. Ουσιαστικά η υλοποίηση τέτοιων σχέσεων τροφοδοτεί τον αλγόριθμο με επιπλέον πληροφορία που αφορά στο συγκεκριμένο πεδίο εφαρμογής και του επιτρέπει την εξαγωγή κανόνων που ταιριάζουν καλύτερα στο μοντέλο με το οποίο δουλεύει ο ειδικός. Επιπλέον, οι σχέσεις αυτές είναι τις περισσότερες φορές γνωστές στον ειδικό ή ακόμα και αυτονόητες. Έτσι με ελάχιστο επιπλέον κόστος για τον χρήστη βελτιώνεται κατά πολύ το τελικό αποτέλεσμα.

Στην προσπάθεια υλοποίησης των σχέσεων εξάρτησης, χρησιμοποιήθηκε και τροποποιήθηκε ο αλγόριθμος NewId, γεγονός που οδήγησε στην δημιουργία του αλγόριθμου IDDD (Inductive Domain Dependent Decision) [22] που αποτελεί και το βασικό προϊόν της εργασίας αυτής. Ο IDDD λειτουργεί όπως και ο NewId όταν δεν υπάρχουν σχέσεις εξάρτησης. Αξίζει να σημειωθεί πως οι σχέσεις που μελετήθηκαν και ορίστηκαν είναι ανεξάρτητες του συγκεκριμένου αλγόριθμου μάθησης για τον οποίο υλοποιήθηκαν. Είναι εύκολη η αναπροσαρμογή της υλοποίησης τους ώστε να εφαρμοστούν σε κάποιο άλλο αλγόριθμο μάθησης μέσω παραδειγμάτων.

Για να αξιολογηθεί ο IDDD, εφαρμόστηκε σε μία πραγματική εφαρμογή. Σημειώ-

νεται πως η εφαρμογή αυτή είναι αρκετά γενική με την έννοια πως τα προβλήματα που παρουσιάζει δεν είναι μοναδικά αλλά παρουσιάζονται και σε πάρα πολλές άλλες εφαρμογές.

Οι αλγόριθμοι μάθησης μέσω παραδειγμάτων προποθέτουν ότι σε κάθε παράδειγμα έχει αντιστοιχηθεί μία και μόνο κλάση. Δυστυχώς σε εφαρμογές που σχετίζονται με αποφάσεις παρατηρείται το φαινόμενο της ανθρώπινης ασυνέπειας. Για να γίνει σαφέστερη η παραπάνω παρατήρηση σημειώνουμε πως σε ίδιες περιπτώσεις η απόφαση του ειδικού στον τομέα μπορεί να διαφέρει είτε λόγω λάθους είτε λόγω ιδιαιτέρων καταστάσεων που δεν καταγράφονται ή δεν σχετίζονται άμεσα με το πεδίο εφαρμογής. Ωστόσο κατά τη δημιουργία του συνόλου κανόνων είναι επιθυμητή η εξαίρεση τέτοιων περιπτώσεων. Επειδή πολλές φορές τα παραδείγματα παρέχονται στο σύστημα αυτόματα από ήδη υπάρχοντα αρχεία είναι πιθανή η ύπαρξη δύο ή περισσότερων παραδειγμάτων τα οποία αν και είναι απολύτως όμοια έχουν ωστόσο χαρακτηριστεί σαν μέλη διαφορετικών κλάσεων. Από τέτοια παραδείγματα είναι αδύνατη η εξαγωγή οποιουδήποτε λογικού συμπεράσματος ακόμα και για την ανθρώπινη νοημοσύνη και πολύ περισσότερο για οποιοδήποτε σύστημα μηχανικής μάθησης. Για το λόγο αυτό αναπτύχθηκε ένα διαλογικό σύστημα προεπεξεργασίας των παραδειγμάτων. Το σύστημα αυτό ενεργοποιείται από το χρήστη, εντοπίζει ασυνέπειες στο σύνολο των παραδειγμάτων και επιτρέπει τη διόρθωση ενδεχομένων λαθών.

Η εργασία αυτή οργανώνεται ως εξής. Στο κεφάλαιο 3 περιγράφεται ο γενικός αλγόριθμος ID3, ενώ στην παράγραφο 3.3 δίνεται η περιγραφή του NewId. Στο κεφάλαιο 4 παρουσιάζεται ο αλγόριθμος IDDD καθώς και η μελέτη των σχέσεων εξάρτησης η οποία παρατίθεται στις παραγράφους 4.1 και 4.2. Στο κεφάλαιο 5 παρατίθενται τα πειραματικά αποτελέσματα για την εφαρμογή που περιγράφεται στη παράγραφο 5.1. Στη παράγραφο 5.2 περιγράφεται το σύστημα προεπεξεργασίας παραδειγμάτων εκπαίδευσης. Στις παραγράφους 5.3 και 5.4 περιγράφονται τα πειράματα που έγιναν και παραθέτονται και αξιολογούνται τα αντίστοιχα αποτελέσματα. Η αξιολόγηση των αποτελεσμάτων έγινε από ειδικούς του πεδίου εφαρμογής. Η σύγκριση μεταξύ των κανόνων που παράγουν οι αλγόριθμοι IDDD και NewId αποτελεί τη βάση αξιολόγησης του αλγόριθμου IDDD

που προτείνεται σε αυτή την εργασία. Παράλληλο προϊόν της εργασίας αποτελεί και ο ορισμός μιας τυπικά ορισμένης κλίμακας της ποιότητας των κανόνων, παράγραφος 5.4.

Τέλος στο κεφάλαιο 6 παραθέτονται τα συμπεράσματα που εξάγονται από αυτή την εργασία με βάση τα πειραματικά αποτελέσματα. Μέσα από τα συμπεράσματα αυτά προκύπτουν και οι δυνατότητες που υπάρχουν για επέκταση της ερευνητικής προσπάθειας προς την κατεύθυνση του εντοπισμού νέων σχέσεων και επομένως της ανάπτυξης νέων αποδοτικότερων αλγορίθμων.

3 Μάθηση μέσω Παραδειγμάτων

3.1 Ο αλγόριθμος ID3

Ο αλγόριθμος ID3 ανήκει στην κατηγορία επαγωγικής μάθησης μέσω παραδειγμάτων. Παρουσιάστηκε για πρώτη φορά από τον Quinlan το 1986 [11] και έχει γνωρίσει ευρεία αποδοχή.

Ο αλγόριθμος στοχεύει στην εξαγωγή ενός συνόλου κανόνων, σύμφωνα με τους οποίους, τα παραδείγματα θα κατατάσσονται σε κλάσεις (κατηγορίες) οι οποίες είναι εκ των προτέρων καθορισμένες. Το πρόβλημα αυτό ονομάζεται **κατάταξη** (classification) και είναι πολύ γνωστό από την βιβλιογραφία είτε σαν κύριο είτε σαν υποπρόβλημα ενός γενικότερου προβλήματος. Συχνά το πρόβλημα της κατάταξης των παραδειγμάτων σε κλάσεις συγχέεται με το πρόβλημα της ομαδοποίησης των παραδειγμάτων σε κλάσεις (clustering). Πρόκειται όμως για δύο εντελώς διαφορετικά προβλήματα. Στην πρώτη περίπτωση, είναι γνωστές οι κλάσεις και το ζητούμενο είναι η εύρεση των κανόνων που περιγράφουν τις ιδιομορφίες τους και διαφοροποιούν μεταξύ τους. Αντίθετα, στη δεύτερη περίπτωση, το ζητούμενο είναι η εύρεση των ομοιοτήτων μεταξύ παραδειγμάτων και η ομαδοποίηση τους σε κλάσεις.

Το σύνολο των κανόνων, για τον ID3, παίρνει τη μορφή ενός δέντρου απόφασης. Τα δέντρα απόφασης χρησιμοποιούνται συχνά για την αναπαράσταση των κανόνων κατάταξης των παραδειγμάτων σε κλάσεις [23, 24, 25]. Οι κόμβοι του δέντρου αντιστοιχούν σε ελέγχους των τιμών κάποιων ιδιοτήτων ενώ τα κλαδιά είναι οι δυνατές τιμές των αντίστοιχων ιδιοτήτων. Κάθε φύλλο αναπαριστά μία συγκεκριμένη κλάση. Κάθε μονοπάτι του δέντρου, από τη ρίζα του μέχρι κάποιο συγκεκριμένο φύλλο, αποτελεί έναν ανεξάρτητο κανόνα κατάταξης παραδειγμάτων σε μία από τις κλάσεις.

Πριν χρησιμοποιήσει κανείς τον ID3, θα πρέπει να εντοπίσει εκείνες τις ιδιότητες που χαρακτηρίζουν/περιγράφουν το πεδίο και να εξετάσει ποιές είναι οι δυνατές τιμές τους. Ο αλγόριθμος αναγνωρίζει μόνο ένα τύπο ιδιοτήτων, τις **κατηγορηματικές** (nominal). Για παράδειγμα, η ιδιότητα χρώμα με τιμές {άσπρο, μαύρο, μπλε} είναι μία

κατηγορηματική ιδιότητα.

Είσοδος του συστήματος είναι αυτή η **περιγραφή** του πεδίου εφαρμογής και ένα σύνολο παραδειγμάτων. Δημιουργούνται δηλαδή δύο αρχεία εισόδου στο σύστημα. Το ένα αρχείο αποτελείται από τον κατάλογο των ιδιοτήτων της εφαρμογής και τις δυνατές τιμές καθεμιάς και το άλλο περιέχει τα παραδείγματα. Τα παραδείγματα είναι διανύσματα τιμών ιδιοτήτων και αναπαριστούν μία συγκεκριμένη κατάσταση, αντικείμενο ή πρόβλημα ανάλογα με το αντικείμενο της εφαρμογής. Ειδικότερα, παράδειγμα για μια ιατρική εφαρμογή είναι το σύνολο των εξετάσεων ενός ασθενούς, για μια εφαρμογή αναγνώρισης αντικειμένων παράδειγμα είναι οι τιμές των γεωμετρικών ιδιοτήτων που περιγράφουν τα αντικείμενα κ.λ.π Τα παραδείγματα αυτά ονομάζονται εκπαιδευτικά παραδείγματα γιατί από αυτά θα εξαχθούν οι κανόνες κατάταξης. Τα εκπαιδευτικά παραδείγματα πρέπει, επιπλέον, να περιέχουν την πληροφορία για την κλάση στην οποία ανήκουν.

Στόχος του ID3 είναι η κατασκευή ενός δέντρου απόφασης από το σύνολο EX των παραδειγμάτων εκπαίδευσης. Αν το EX είναι κενό ή περιέχει παραδείγματα που ανήκουν σε μία κλάση, το απλούστερο δέντρο απόφασης είναι ένας κόμβος-φύλλο με την ετικέτα της κλάσης. Η παραπάνω περίπτωση είναι τετριμμένη. Στη συνηθισμένη περίπτωση ο κόμβος αναπαριστά έλεγχο της τιμής κάποιας ιδιότητας. Για κάθε παράδειγμα, ελέγχεται η τιμή της ιδιότητας-κόμβος. Αν οι δυνατές τιμές είναι π.χ A_1, A_2, \dots, A_v , ο έλεγχος αυτός έχει σαν αποτέλεσμα την δημιουργία EX_1, EX_2, \dots, EX_v αντίστοιχων υποσυνόλων. Αν η διαδικασία αυτή επαναληφθεί για κάθε EX_i , το αποτέλεσμα θα είναι ένα δέντρο απόφασης για το σύνολο των παραδειγμάτων EX . Εφόσον δύο ή περισσότερα υποσύνολα είναι μη κενά, κάθε EX_i είναι μικρότερο από το E (γνήσιο υποσύνολο). Στη χειρότερη περίπτωση, αυτή η πολιτική του 'διαίρει και βασίλευε' θα δημιουργήσει μονομελή υποσύνολα και επομένως το πρόβλημα ανάγεται στην τετριμμένη περίπτωση.

Η επιλογή της ιδιότητας που θα αποτελέσει κόμβο είναι το κρίσιμο ζήτημα για την παραγωγή απλών δέντρων απόφασης. Ο τρόπος επιλογής της ιδιότητας σχετίζεται με τη μεγιστοποίηση μιας συνάρτησης υπολογισμού της πληροφοριακής ισχύος της [23, 26, 27, 28, 18, 11, 12]. Πρακτικά, μία ιδιότητα έχει μεγάλη πληροφοριακή ισχύ αν οι τιμές της, στο σύνολο των παραδειγμάτων, είναι τέτοιες ώστε να επιτυγχάνεται

όσο το δυνατόν μεγαλύτερος διαχωρισμός μεταξύ των παραδειγμάτων που ανήκουν σε διαφορετική κλάση. Πιο τυπικά, η μεθοδολογία στηρίζεται σε δύο βασικές υποθέσεις. Για απλότητα θεωρούμε δύο μόνο κλάσεις και υποθέτουμε ότι το σύνολο EX των παραδειγμάτων περιέχει p παραδείγματα της κλάσης P και n παραδείγματα της κλάσης N . Οι υποθέσεις είναι:

1. Κάθε σωστό δέντρο απόφασης θα πρέπει να κατατάσσει τα παραδείγματα με την ίδια αναλογία που αυτά παρουσιάζονται στο σύνολο εκπαίδευσης. Ένα τυχαίο παράδειγμα ανήκει στην κλάση P με πιθανότητα $p/(p+n)$ και στην κλάση N με πιθανότητα $n/(p+n)$.
2. Το δέντρο απόφασης χρησιμοποιείται για την κατάταξη ενός παραδείγματος σε μία από τις κλάσεις P και N . Έτσι μπορεί να θεωρηθεί σαν μια γεννήτρια μηνυμάτων 'P' και 'N'. Η πληροφορία που απαιτείται για τη δημιουργία των μηνυμάτων αυτών είναι

$$I(p, n) = -\frac{p}{(p+n)} \log_2 \frac{p}{(p+n)} - \frac{n}{(p+n)} \log_2 \frac{n}{(p+n)} \quad (1)$$

Ο παραπάνω τύπος γενικεύεται για περισσότερες κλάσεις. Αν C_1, \dots, C_k οι δυνατές κλάσεις και e_1, \dots, e_k το πλήθος των παραδειγμάτων κάθε κλάσης, τότε:

$$I(C_1, \dots, C_k) = -\sum_{i=1}^k \frac{e_i}{\sum_{j=1}^k e_j} \log_2 \frac{e_i}{\sum_{j=1}^k e_j} \quad (2)$$

Ας υποθέσουμε πως η ιδιότητα A με τιμές A_1, A_2, \dots, A_v αποτελεί τη ρίζα του δέντρου. Τότε το σύνολο των εκπαιδευτικών παραδειγμάτων EX διαιρείται αντίστοιχα στα υποσύνολα EX_1, EX_2, \dots, EX_v , όπου το EX_i περιέχει εκείνα τα παραδείγματα του EX με τιμή A_i για την ιδιότητα A . Υποθέτουμε πως το EX_i περιέχει p_i παραδείγματα κλάσης P και n_i παραδείγματα κλάσης N . Η πληροφορία που απαιτείται για την κάλυψη του υποσυνόλου EX_i είναι $I(p_i, n_i)$. Για ένα δέντρο με ρίζα την A , χρησιμοποιείται η παρακάτω συνάρτηση εντροπίας για την μέτρηση της αβεβαιότητας του συστήματος μετά την επιλογή της A .

$$E(A) = \sum_{i=1}^v \frac{p_i + n_i}{p + n} I(p_i, n_i) \quad (3)$$

Από τις σχέσεις (1) και (3) προκύπτει πως η **πληροφοριακή ισχύς** της ιδιότητας A (το κέρδος από την επιλογή της) είναι:

$$gain(A) = I(p, n) - E(A) \quad (4)$$

Στην σχέση αυτή η ποσότητα $I(p, n)$ είναι σταθερή. Επομένως μεγάλη πληροφοριακή ισχύ έχει η ιδιότητα που επιφέρει, μετά την επιλογή της, την μεγαλύτερη μείωση στην εντροπία του συστήματος.

Ο αλγόριθμος υπολογίζει την πληροφοριακή ισχύ κάθε ιδιότητας και επιλέγει εκείνη με την μεγαλύτερη πληροφοριακή ισχύ. Η ιδιότητα αυτή ονομάζεται **πιο πληροφοριακή** ή **πιο διαχωριστική**. Στα πλαίσια αυτής της εργασίας θα χρησιμοποιούνται και οι δύο όροι. Η ιδιότητα αυτή γίνεται κόμβος του δέντρου, δημιουργούνται τόσα κλαδιά όσες είναι οι δυνατές τιμές της και τα παραδείγματα χωρίζονται στα αντίστοιχα υποσύνολα.

Η παραπάνω διαδικασία επιλογής της ιδιότητας-κόμβος επαναλαμβάνεται αναδρομικά, για κάθε υποσύνολο παραδειγμάτων. Όταν το σύνολο αυτό περιέχει παραδείγματα μίας μόνο κλάσης ή είναι κενό, ικανοποιείται η συνθήκη τερματισμού και δημιουργείται ένα φύλλο στο οποίο αντιστοιχίζεται η κλάση ή το ειδικό σύμβολο κενό (Null), αντίστοιχα. Σημειώνουμε πως οι τυχόν εμφανίσεις της κενής κλάσης οφείλονται στο γεγονός ότι ο ID3 δημιουργεί κλαδιά για όλες τις τιμές της ιδιότητας που επιλέχτηκε σαν κόμβος και όχι μόνο γι'αυτές που εμφανίζονται στα παραδείγματα.

Για κάθε σύνολο εκπαιδευτικών παραδειγμάτων υπάρχουν περισσότερα από ένα σωστά δέντρα απόφασης. Μία προσέγγιση είναι η κατασκευή όλων των δέντρων που κατατάσσουν σωστά τα εκπαιδευτικά παραδείγματα και η επιλογή του πιο απλού από αυτά. Ο αριθμός των δέντρων αυτών είναι πεπερασμένος αλλά πολύ μεγάλος κι έτσι η προσέγγιση αυτή μπορεί να χρησιμοποιηθεί μόνο για μικρές εφαρμογές. Ο ID3 προσαθεί να κατασκευάσει ένα αρκετά καλό δέντρο χωρίς μεγάλο υπολογιστικό κόστος. Για παράδειγμα, αν θεωρήσουμε το σύνολο παραδειγμάτων που περιέχονται στον πίνακα 1, ο ID3 δημιουργεί το δέντρο του σχήματος 1 και όχι αυτό του σχήματος 2 το οποίο όμως είναι επίσης σωστό. Γι'αυτό το παράδειγμα, το πρόβλημα που περιγράφεται συνίσταται

στη λήψη της απόφασης αν πάρει κάποιος μαζί του ομπρέλλα κρίνοντας από μερικές ενδείξεις του καιρού. Οι ιδιότητες που περιγράφουν τον καιρό είναι {Οψη Καιρού, Θερμοκρασία, Υγρασία, Αέρας} και οι κλάσεις είναι {Ναι, Οχι}. Ο ID3 κατασκευάζει απλά δέντρα αλλά δεν είναι εγγυημένα τα απλούστερα δυνατά.

	Ιδιότητες				
ΑΑ	Οψη Καιρού	Θερμοκρασία	Υγρασία	Αέρας	Κλάση
1	Λιακάδα	Υψηλή	Υψηλή	Οχι	Οχι
2	Λιακάδα	Υψηλή	Υψηλή	Ναι	Οχι
3	Συννεφιά	Υψηλή	Υψηλή	Οχι	Ναι
4	Βροχερός	Μέτρια	Υψηλή	Οχι	Ναι
5	Βροχερός	Χαμηλή	Κανονική	Οχι	Ναι
6	Βροχερός	Χαμηλή	Κανονική	Ναι	Οχι
7	Συννεφιά	Χαμηλή	Κανονική	Ναι	Ναι
8	Λιακάδα	Μέτρια	Υψηλή	Οχι	Οχι
9	Λιακάδα	Χαμηλή	Κανονική	Οχι	Ναι
10	Βροχερός	Μέτρια	Κανονική	Οχι	Ναι
11	Λιακάδα	Μέτρια	Κανονική	Ναι	Ναι
12	Συννεφιά	Μέτρια	Υψηλή	Ναι	Ναι
13	Συννεφιά	Υψηλή	Κανονική	Οχι	Ναι
14	Βροχερός	Μέτρια	Υψηλή	Ναι	Οχι

Πίνακας 1: Ένα μικρό σύνολο εκπαιδευτικών παραδειγμάτων

Σχήμα 1: Ένα απλό δέντρο απόφασης

Σχήμα 2: Ένα πολύπλοκο δέντρο απόφασης

Όσον αφορά την πολυπλοκότητα του αλγορίθμου, η πολυπλοκότητά της εύρεσης κάθε κόμβου είναι $O(|EX| \times |A|)$, όπου EX και A ο αριθμός των παραδειγμάτων και ιδιοτήτων αντίστοιχα. Έτσι η ολική πολυπλοκότητα, ανά αναδρομική κλήση, είναι ανάλογη του γινομένου του αριθμού παραδειγμάτων, του αριθμού των ιδιοτήτων και του αριθμού των κόμβων του δέντρου που δεν είναι φύλλα. Η ίδια σχέση ισχύει και για την ολική πολυπλοκότητα του αλγορίθμου ακόμα και για μεγάλο πλήθος αναδρομών. Από τα μέχρι τώρα πειράματα με τον ID3, δεν παρατηρήθηκε εκθετική αύξηση του χώρου ή

χρόνου που απαιτεί, ακόμα και για πολύ μεγάλες εφαρμογές (πολλά παραδείγματα και ιδιότητες).

Ένα από τα μεγαλύτερα πλεονεκτήματα του ID3 είναι ότι προβλέπονται περιπτώσεις όπου η τιμή κάποιας ιδιότητας είναι άγνωστη σε ένα παράδειγμα. Ο κατάλογος των ιδιοτήτων που περιγράφουν το πεδίο εφαρμογής πρέπει να είναι πλήρης. Για παράδειγμα, στις ιατρικές εφαρμογές, πρέπει να περιλαμβάνονται όλες οι δυνατές εξετάσεις που είναι δυνατόν να γίνουν προκειμένου να καταλήξει ο γιατρός στη διάγνωση. Ωστόσο δεν είναι δυνατόν να υποβληθεί ο ασθενής σε όλες αυτές τις εξετάσεις όταν δεν είναι αναγκαίο. Επειδή πρέπει το παράδειγμα να περιλαμβάνει τιμές για όλες τις ιδιότητες, ο αλγόριθμος επιτρέπει την χρησιμοποίηση μιας ειδικής τιμής που υποδηλώνει ακριβώς την μη διαθεσιμότητά της και ονομάζεται **άγνωστη τιμή** (unknown). Η διαδικασία δημιουργίας του δέντρου απόφασης έχει μεταβληθεί κατάλληλα για το χειρισμό αυτών των τιμών.

Συγκεκριμένα, ας υποθέσουμε πως έχει επιλεγεί σαν κόμβος η ιδιότητα A με τιμές A_1, A_2, \dots, A_v . Το σύνολο παραδειγμάτων πρέπει να χωριστεί σε υποσύνολα με βάση τις τιμές των παραδειγμάτων για την A. Υποθέτουμε πως στο σύνολο E υπάρχουν p_i παραδείγματα με τιμή A_i που ανήκουν στην κλάση P και n_i παραδείγματα με τιμή A_i που ανήκουν στην κλάση N. Ακόμη, έστω p_u και n_u τα παραδείγματα που ανήκουν στις κλάσεις P και N αντίστοιχα και έχουν άγνωστες τιμές για την ιδιότητα A. Η πληροφοριακή ισχύς της ιδιότητας A υπολογίζεται σαν ο πραγματικός αριθμός των p_i να ήταν:

$$p_i + p_u * ratio_i, \quad (5)$$

όπου

$$ratio_i = \frac{(p_i + n_i)}{\sum_{i=1}^v (p_i + n_i)} \quad (6)$$

και παρόμοια για τον αριθμό των n_i . Αυτή η έκφραση έχει το πλεονέκτημα πως η παρουσία αγνώστων τιμών μειώνει την πληροφοριακή ισχύ μιας ιδιότητας.

Αν, παρόλα αυτά, η ιδιότητα επιλεγεί, τα παραδείγματα γίνονται μέλη όλων των υποσυνόλων που δημιουργούνται αλλά αντιστοιχίζεται σε αυτά ένα βάρος που είναι το αντίστοιχο $ratio_i$.

3.2 Μάθηση μέσω Παραδειγμάτων και Εμπειρα Συστήματα

Ενας αλγόριθμος μάθησης μπορεί αρχικά να χρησιμοποιηθεί σαν πηγή γνώσης για ένα έμπειρο σύστημα. Στις περιπτώσεις που η γνώση αυτή δεν μπορεί να εφαρμοστεί σε μία δεδομένη κατάσταση στην οποία έχει περιέλθει το σύστημα, ενεργοποιείται ο αλγόριθμος μάθησης για να εντοπίσει τυχόν ιδιαιτερότητες, εξαιρέσεις και να εμπλουτίσει την υπάρχουσα γνώση.

Ειδικότερα για τον ID3, χρειάζεται μία μετατροπή της αναπαράστασης της πληροφορίας που περιέχεται στο δέντρο απόφασης. Η πληροφορία αυτή πρέπει να αναπαρασταθεί με κανόνες παραγωγής της μορφής

ΑΝ (Προϋποθέσεις) ΤΟΤΕ (Κλάση) (Μέτρο Βεβαιότητας)

Η μετατροπή αυτή κρίνεται αναγκαία για τρεις λόγους. Πρώτον, οι κανόνες παραγωγής είναι ο πιο διαδεδομένος τρόπος αναπαράστασης της γνώσης στα έμπειρα συστήματα. Δεύτερον, είναι δύσκολο για έναν ειδικό να κατανοήσει και να μεταβάλει ένα δέντρο απόφασης. Τέλος, η μετατροπή αυτή μπορεί να αυξήσει την ακρίβεια κατάταξης των παραδειγμάτων με την αφαίρεση από το δέντρο ιδιοτήτων που ανταποκρίνονται σε ιδιαιτερότητες του συνόλου των εκπαιδευτικών παραδειγμάτων. Το πρόβλημα της μετατροπής αυτής έχει περιγραφεί στην βιβλιογραφία όπου και προτείνονται τρόποι αποδοτικής επίλυσής του [12].

Η χρησιμοποίηση του ID3 κατά την ανάπτυξη και λειτουργία ενός έμπειρου συστήματος δεν αποτέλεσε ερευνητικό θέμα στα πλαίσια της παρούσας εργασίας.

3.3 Ο αλγόριθμος NewId

Το σύστημα NewId, που είχαμε στη διάθεσή μας [17], αποτελεί μία βελτιωμένη έκδοση του αρχικού αλγορίθμου ID3. Αναπτύχθηκε από το Turing Institute, στη Γλασκώβη, στα πλαίσια του MLT (Machine Learning Toolbox) (Esprit project). Ο κύριος αλγόριθμος εξαγωγής του δέντρου απόφασης παραμένει ο ίδιος και για το λόγο αυτό θα περιοριστούμε στον εντοπισμό των ιδιοτήτων του NewId που αφορούν κυρίως στο χειρισμό νέων τύπων ιδιοτήτων και τιμών που αυτές παίρνουν.

3.3.1 Τύποι Ιδιοτήτων-Τιμών του NewId

Ο ID3 προβλέπει μόνο ένα τύπο ιδιοτήτων, τις **κατηγορηματικές** (nominal). Σε πολλές εφαρμογές όμως υπάρχουν και **αριθμητικές** (linear) ιδιότητες. Στον NewId έχει αλλάξει η διαδικασία δημιουργίας του δέντρου για να αντιμετωπιστεί η περίπτωση επιλογής τέτοιου είδους ιδιότητας σαν κόμβος. Επιτρέπονται αριθμητικές ιδιότητες που παίρνουν ακέραιες ή και πραγματικές τιμές. Όταν μία αριθμητική ιδιότητα επιλέγεται σαν κόμβος του δέντρου, υπολογίζεται στατιστικά από τα παραδείγματα εκείνη η αριθμητική τιμή που διαχωρίζει καλύτερα τις κλάσεις. Η τιμή αυτή ονομάζεται **σημείο διαχωρισμού** (split point). Δημιουργείται τότε ένα κλαδί για τα παραδείγματα με τιμές ιδιότητας μεγαλύτερες του σημείου διαχωρισμού και ένα δεύτερο κλαδί με τιμές μικρότερες του σημείου διαχωρισμού.

Επιπλέον, ακόμα και για τις κατηγορηματικές ιδιότητες, μπορεί κανείς να ορίσει διάταξη μεταξύ των τιμών τους. Ας θεωρήσουμε την ιδιότητα επάγγελμα με τιμές {καθαρίστρια, γραμματέας, διευθυντής}. Είναι πιθανό να μας ενδιαφέρει να ορίσουμε πως η θέση του διευθυντή, για παράδειγμα, είναι υψηλότερη αυτής του γραμματέα. Δηλώνουμε πως το επάγγελμα είναι μία **διατεταγμένη** ιδιότητα και την διάταξη των τιμών της καθορίζει η σειρά με την οποία αυτές παρατίθενται. Ο NewId διαχειρίζεται τις ιδιότητες αυτές σαν να ήταν αριθμητικές. Συγκεκριμένα στο αρχείο περιγραφής της εφαρμογής η διατεταγμένη ιδιότητα ορίζεται:

επάγγελμα: (ORDERED) καθαρίστρια, γραμματέας, διευθυντής

Όσον αφορά τους τύπους ιδιοτήτων υπάρχουν ακόμα οι **ιεραρχικές** ιδιότητες. Τέτοιες ιδιότητες είναι αυτές των οποίων οι τιμές έχουν κάποια ιεραρχία, μπορούν δηλαδή να καθοριστούν σε διαφορετικά επίπεδα λεπτομέρειας. Για παράδειγμα η ιδιότητα οχήμα μπορεί να πάρει μία από τις τιμές {τετράγωνο, τρίγωνο}. Επιπλέον, το τρίγωνο μπορεί να αποτελεί μία διαφορετική ιδιότητα με τιμές {οκαληνό, ισόπλευρο, ορθογώνιο}. Η ύπαρξη των ιεραρχικών ιδιοτήτων μας επιτρέπει να ορίσουμε μία μόνο ιδιότητα με το όνομα οχήμα και να δηλώσουμε πως κάθε τιμή της είναι ουσιαστικά μία ιδιότητα με το δικό της σύνολο τιμών.

Όσον αφορά τις τιμές των ιδιοτήτων ο NewId, όπως και ο ID3, διαχειρίζεται άγνωστες τιμές αλλά εισάγει και μία επιπλέον ειδική τιμή, την **αδιάφορη** (don't care). Όταν σε ένα συγκεκριμένο παράδειγμα μία ιδιότητα παίρνει κάποια συγκεκριμένη τιμή είναι δυνατόν να γίνεται αδιάφορη η τιμή κάποιας άλλης ιδιότητας. Για να πληροφορήσουμε τον αλγόριθμο πως δεν θέλουμε να λάβει υπόψιν του αυτή τη δεύτερη ιδιότητα της δίνουμε την ειδική τιμή αδιάφορη. Για να γίνει κατανοητή η ανάγκη ύπαρξης της τιμής αυτής παραθέτουμε το εξής παράδειγμα. Για μια ιατρική εφαρμογή πρέπει να καταγράψουμε όλες τις δυνατές εξετάσεις που μπορεί να γίνουν για να καταλήξουμε σε διάγνωση. Προφανώς, ο ασθενής δεν υποβάλεται σε όλες αυτές τις εξετάσεις αν έχει ήδη αποφασιστεί η διάγνωση με βάση μερικά στοιχεία μόνο. Οι τιμές των άλλων ιδιοτήτων είναι μεν άγνωστες αλλά όχι γιατί τα στοιχεία δεν είναι διαθέσιμα αλλά γιατί δεν είναι σημαντικά. Σε αυτές τις περιπτώσεις χρησιμοποιείται η αδιάφορη τιμή.

Ο NewId θεωρεί πως, αν μία ιδιότητα έχει αδιάφορη τιμή, θα μπορούσε να πάρει οποιαδήποτε δυνατή τιμή. Κατά τον υπολογισμό της ιδιότητας με τη μεγαλύτερη πληροφοριακή ισχύ, ισχύει για τις αδιάφορες τιμές, η πολιτική που ισχύει και για τις άγνωστες. Συγκεκριμένα, κάθε παράδειγμα που έχει άγνωστη τιμή για την συγκεκριμένη ιδιότητα θεωρείται πως μπορεί να έχει οποιαδήποτε τιμή. Αρχικά το βάρος κάθε παραδείγματος είναι ίσο με τη μονάδα. Τα παραδείγματα που έχουν άγνωστες τιμές για την ιδιότητα αυτή γίνονται μέλη όλων των υποσυνόλων που δημιουργούνται αλλά αντιστοιχίζεται σε αυτά ένα βάρος. Όσον αφορά στα παραδείγματα με αδιάφορη τιμή, γίνονται κι αυτά μέλη όλων των υποσυνόλων χωρίς όμως να μεταβληθεί το βάρος τους. Το βάρος αυτό

θα είναι ίσο με τη μονάδα ή κάποια μικρότερη πραγματική τιμή αν το συγκεκριμένο παράδειγμα εκτός από την αδιάφορη εμπεριέχει και άγνωστες τιμές σε κάποιες ιδιότητες που έχουν επιλεγεί από τον NewId νωρίτερα.

Σε κάθε βήμα του αλγορίθμου ελέγχεται το τρέχον σύνολο παραδειγμάτων. Αν αυτό περιέχει παραδείγματα μίας μόνο κλάσης, ο κόμβος γίνεται φύλλο και παίρνει τον χαρακτηρισμό της κλάσης αυτής. Αν το σύνολο παραδειγμάτων είναι κενό τερματίζει και πάλι ο αλγόριθμος δημιουργώντας ένα φύλλο με τον χαρακτηρισμό **κενό** (null). Επιπλέον, υπάρχει το ενδεχόμενο να έχουν εξαντληθεί οι ιδιότητες αλλά τα παραδείγματα να ανήκουν σε διαφορετικές κλάσεις. Το σύνολο παραδειγμάτων δεν είναι δυνατόν να διαχωριστεί παραπέρα και έτσι δημιουργείται ένα φύλλο με τον χαρακτηρισμό **σύγκρουση** (clash).

3.3.2 Σχέση Διάταξης στον NewId

Κατά την ανάπτυξη του NewId προβλέφθηκε η ύπαρξη μιας ειδικής σχέσης εξάρτησης ιδιοτήτων. Επειδή ο τρόπος επιλογής ιδιοτήτων στηρίζεται σε στατιστικά μέτρα, μπορεί η σειρά επιλογής να παίζει σημαντικό ρόλο. Έτσι ορίζεται μία σχέση **διάταξης** (ordering) μεταξύ δύο ή περισσότερων ιδιοτήτων. Ο ορισμός της σχέσης είναι πολύ απλός. Στο τέλος του αρχείου δήλωσης των ιδιοτήτων και των τιμών τους, προσθέτουμε προτάσεις της μορφής:

A BEFORE B (Η ιδιότητα A ΠΡΙΝ ΤΗΝ B)

Με τον τρόπο αυτό μπορούμε να καθορίσουμε μία σχέση μερικής ή ολικής διάταξης στο σύνολο των ιδιοτήτων. Αν ορίσουμε πως μία ιδιότητα προηγείται όλων των άλλων, ουσιαστικά καθοδηγούμε τον αλγόριθμο να τοποθετήσει την ιδιότητα αυτή στη ρίζα του δέντρου απόφασης και επομένως να θεωρείται σημαντικότερη (πιο πληροφοριακή) από όλες τις άλλες.

Ο NewId χρησιμοποιεί τις σχέσεις διάταξης για να περιορίσει τον αριθμό των ιδιοτήτων που λαμβάνει υπόψην του κατά τη διαδικασία επιλογής της πιο πληροφοριακής ιδιότητας. Συγκεκριμένα, πριν υπολογίσει την συνάρτηση πληροφοριακής ισχύος εξετάζει

αν έχει οριστεί σχέση διάταξης σύμφωνα με την οποία η ιδιότητα αυτή έπεται κάποιας άλλης. Αν όχι, προχωρά στον υπολογισμό της συνάρτησης. Αν όμως έχει οριστεί πως η ιδιότητα αυτή, έστω B , έπεται κάποιας ιδιότητας A , υπάρχουν δύο ενδεχόμενα:

1. Για το τρέχον μονοπάτι, η ιδιότητα A έχει ήδη επιλεγεί σε κάποιο προηγούμενο βήμα. Στην περίπτωση αυτή υπολογίζεται κανονικά η συνάρτηση πληροφορίας και επιτρέπεται η επιλογή της B .
2. Για το τρέχον μονοπάτι, η ιδιότητα A δεν έχει επιλεγεί σε κάποιο προηγούμενο βήμα. Στην περίπτωση αυτή η ιδιότητα B απορρίπτεται έστω και αν, υπολογίζοντας την πληροφοριακή της ισχύ, προέκυπτε σαν η περισσότερο διαχωριστική από τις ιδιότητες.

Σε πολλές περιπτώσεις μπορεί να εμφανιστούν δύο ή περισσότερες ιδιότητες με την ίδια πληροφοριακή ισχύ. Αν δεν έχει οριστεί μεταξύ τους σχέση διάταξης, επιλέγεται η ιδιότητα που εμφανίζεται πρώτη στην περιγραφή της εφαρμογής. Η επιλογή αυτή είναι τυχαία και επομένως το αποτέλεσμα αλλάζει με την απλή αναδιάταξη των ιδιοτήτων. Η τυχειότητα αυτή αποφεύγεται με τον ορισμό της διάταξης. Αν η σχέση διάταξης δεν υπαγορεύεται από την εφαρμογή αλλά ορίζεται από τον χρήστη στην προσπάθειά του να εξάγει καλύτερα αποτελέσματα, υπάρχει πάντα το ενδεχόμενο ελάχιστες μεταβολές του συνόλου των παραδειγμάτων να μεταβάλουν κατά πολύ το τελικό αποτέλεσμα. Οι αλλαγές αυτές θα μεταβάλουν την πληροφοριακή ισχύ κάθε ιδιότητας με αποτέλεσμα η σχέση διάταξης να μην βελτιώνει ή ακόμα και να επιδρά αρνητικά στο τελικό αποτέλεσμα.

4 Ο αλγόριθμος IDDD

Ο βασικός στόχος της εργασίας αυτής είναι η ενίσχυση των επαγωγικών μηχανισμών μάθησης μέσω παραδειγμάτων με βάση τη γνώση για το πεδίο εφαρμογής. Η χρησιμότητα της γνώσης αυτής έχει τεκμηριωθεί στη βιβλιογραφία [16, 29, 20, 21, 15]. Οι αλγόριθμοι μάθησης μέσω παραδειγμάτων επιδιώκουν την εξάλειψη της ανάγκης εισαγωγής αυτής της γνώσης με σκοπό τον περιορισμό του ρόλου του ειδικού. Έτσι επιτυγχάνουν και γενική εφαρμογή τους σε οποιοδήποτε πρόβλημα γιατί δεν χρησιμοποιούν γνώση του ειδικού αλλά γενικές μεθόδους. Αν και η διαδικασία εξαγωγής γνώσης από τον ειδικό αποτελεί τις περισσότερες φορές το πιο δύσκολο κομμάτι της ανάπτυξης έμπειρων συστημάτων, η ακραία προσέγγιση των αλγορίθμων μάθησης μέσω παραδειγμάτων οδηγεί σε περιορισμό των δυνατοτήτων μάθησης. Αγνοείται και πληροφορία που εύκολα μπορεί να τυποποιηθεί και να υποβοηθήσει τη μάθηση. Η προσέγγιση που ακολουθείται σε αυτή την εργασία είναι πως η εισαγωγή τέτοιας γνώσης βελτιώνει τους αλγόριθμους χωρίς να τους επιβαρύνει όσον αφορά την πολυπλοκότητα και την ευκολία χρήσης τους. Στο πλαίσιο αυτό μελετήθηκαν οι σχέσεις που μπορεί να υπάρχουν μεταξύ δύο ιδιοτήτων. Σκοπός ήταν ο εντοπισμός και τυπικός ορισμός τέτοιων σχέσεων που διέπουν πραγματικές εφαρμογές αλλά είναι αρκετά γενικές ώστε να μην ισχύουν για μία μόνο εφαρμογή ή ένα είδος εφαρμογών.

Μία προσέγγιση, στο πρόβλημα εισαγωγής γνώσης για το πεδίο εφαρμογής, παρουσιάστηκε από τον Nunez [19]. Σύμφωνα με την προσέγγιση αυτή, η επιλογή κάθε ιδιότητας συνοδεύεται από ένα αντίστοιχο κόστος. Μία ιδιότητα μπορεί να παρέχει σημαντική πληροφορία αλλά να στοιχίζει πολύ η συλλογή της. Για παράδειγμα, στην ιατρική μία εξέταση μπορεί να οδηγεί άμεσα στη διάγνωση αλλά η εξέταση αυτή μπορεί επίσης να είναι επιβαρυντική για τον ασθενή, ακριβή ή ακόμη κι επικίνδυνη. Θα θέλαμε λοιπόν να κάνουμε την εξέταση αυτή μόνο όταν είναι αδύνατη η διάγνωση με βάση τα υπόλοιπα στοιχεία. Με αυτό το σκεπτικό, ο Nunez μετέβαλε την συνάρτηση υπολογισμού της πληροφορίας ώστε να λαμβάνεται υπόψιν και το κόστος επιλογής της ιδιότητας για την οποία υπολογίζεται. Σημειώνουμε πως η προσέγγιση αυτή χρησιμοποιεί μεν γνώση του πεδίου εφαρμογής αλλά θεωρεί κι αυτή κάθε ιδιότητα σαν αυτόνομη πηγή πληροφορίας

και δεν αντιμετωπίζει τις σχέσεις μεταξύ ιδιοτήτων.

Μια χαρακτηριστική τέτοια σχέση είναι η σχέση **εξάρτησης** δύο ιδιοτήτων. Συγκεκριμένα, σε πολλές εφαρμογές η χρησιμοποίηση μιας ιδιότητας σε ένα κανόνα εξαρτάται από την ύπαρξη κάποιας άλλης ιδιότητας ή ακόμα και από την συγκεκριμένη τιμή της. Η εξαρτημένη ιδιότητα ονομάζεται **Κόρη**, ενώ αυτή από την οποία εξαρτάται ονομάζεται **Μάνα**.

Οι σχέσεις αυτές είναι γνωστές από την βιβλιογραφία [15]. Ο Someren μελέτησε τις σχέσεις αυτές και ανέπτυξε ένα σύστημα που τις αναγνωρίζει και μεταβάλλει ανάλογα το σύνολο των ιδιοτήτων και τις τιμές τους. Η προσέγγισή του έγκειται στην αντικατάσταση, κάθε ζευγαριού εξαρτώμενων ιδιοτήτων, με μια άλλη ιδιότητα. Το πεδίο τιμών αυτής της ιδιότητας είναι το σύνολο των νόμιμων συνδυασμών των τιμών των αρχικών ιδιοτήτων. Συγκεκριμένα, έστω η ιδιότητα A με τιμές $\{a_1, a_2, \dots, a_v\}$ και η ιδιότητα B με τιμές $\{b_1, b_2, \dots, b_k\}$. Ας υποθέσουμε πως η B εξαρτάται από την τιμή a_j της A. Αυτό σημαίνει πως η B δεν πρέπει να εμφανίζεται σε ένα κανόνα σε συνδυασμό με την τιμή a_j για την A. Τότε δημιουργείται μία καινούργια ιδιότητα που αντικαθιστά τις A και B και έχει πεδίο τιμών το σύνολο: $\{a_1 \text{ AND } b_1, a_1 \text{ AND } b_2, \dots, a_j, a_{j+1} \text{ AND } b_1, a_{j+1} \text{ AND } b_2, \dots, a_v \text{ AND } b_1, a_v \text{ AND } b_2, \dots, a_v \text{ AND } b_k\}$. Με τον τρόπο αυτό δημιουργείται ένα σύνολο ιδιοτήτων που εκφράζει έμμεσα τις σχέσεις εξάρτησης. Ο αλγόριθμος μάθησης μέσω παραδειγμάτων χρησιμοποιεί αυτό το καινούργιο σύνολο και έτσι αποκλείεται να χρησιμοποιήσει οποιοδήποτε μη νόμιμο συνδυασμό τιμών ιδιοτήτων.

Η παραπάνω μεθοδολογία εκφράζει μεν τις σχέσεις εξάρτησης αλλά παρουσιάζει βασικά μειονεκτήματα:

1. Αν στο παράδειγμα που αναφέραμε παραπάνω η A εξαρτάται από την τιμή κάποιας τρίτης ιδιότητας, αυξάνεται πολύ ο αριθμός των δυνατών συνδυασμών. Αλγόριθμοι όπως ο ID3 τείνουν να διαλέγουν, σαν πιο πληροφοριακές, ιδιότητες με πολυπληθή σύνολα τιμών. Σε περιπτώσεις διαδοχικών εξαρτήσεων, η ιδιότητα που θα δημιουργηθεί θα εμφανίζεται να έχει μεγάλη πληροφοριακή ισχύ αν και αυτό μπορεί, στην πραγματικότητα, να μην ισχύει.

2. Η αντικατάσταση δύο εξαρτώμενων ιδιοτήτων από μία τρίτη μπορεί σε μερικές περιπτώσεις να σημαίνει περιορισμό των επιλογών για τον αλγόριθμο μάθησης. Η ιδιότητα που αντικαθιστά τις εξαρτώμενες μπορεί να έχει μικρή πληροφοριακή ισχύ. Τότε δεν υπάρχει περίπτωση να δημιουργηθούν κανόνες που να λαμβάνουν υπόψη τους τις ιδιότητες που αντικαταστάθηκαν. Η μεθοδολογία του Someren αποτυγχάνει να εκφράσει καταστάσεις όπου θέλουμε η σχέση εξάρτησης να ισχύει μόνο όταν έχουμε ήδη εξετάσει την τιμή της Μάνας.

Κατά την ανάπτυξη του IDDD, ακολουθήσαμε διαφορετική προσέγγιση από αυτή του Someren. Ο στόχος ήταν η χρησιμοποίηση γνώσης του πεδίου εφαρμογής κατά την διαδικασία κατασκευής του δέντρου απόφασης και όχι η μετατροπή του συνόλου ιδιοτήτων. Επιπλέον, ενδιαφερόμαστε για τον χειρισμό των σχέσεων εξάρτησης και όχι για το κόστος κατασκευής του δέντρου απόφασης. Έτσι το κριτήριο επιλογής της πιο πληροφοριακής ιδιότητας δεν μεταβάλλεται με τον τρόπο που πρότεινε ο Nunez.

Κατά τη μελέτη του NewId διαπιστώθηκαν ορισμένες αδυναμίες του όσον αφορά τον εντοπισμό και χειρισμό των σχέσεων εξάρτησης. Η μέτρηση της πληροφοριακής ισχύος κάθε ιδιότητας βασίζεται στα εκπαιδευτικά παραδείγματα. Αυτά αποτελούν και την μοναδική πηγή γνώσης. Ο αλγόριθμος προσπαθεί να εξαγάγει κανόνες που θα καλύψουν όλα τα παραδείγματα. Στην προσπάθειά του αυτή καταλήγει σε εξειδικευμένα δέντρα [27]. Αν τα παραδείγματα δεν αντιπροσωπεύουν καλά το πεδίο εφαρμογής, οι κανόνες δεν μπορεί να εκφράζουν τη διαδικασία λήψης της απόφασης.

Η προεπεξεργασία των παραδειγμάτων είναι δύσκολο να αυτοματοποιηθεί. Θα έπρεπε ο ειδικός να διαλέξει από τα παραδείγματα εκείνα που θεωρεί αντιπροσωπευτικά της εφαρμογής, να εντοπίσει και να εξαλείψει τον θόρυβο. Έτσι όμως το σύστημα κάνει το σοβαρότερο πλεονέκτημά του, την απλότητα. Αν απαιτεί τόσο κόπο από τον ειδικό γιατί να μην χρησιμοποιηθεί κάποιο άλλο σύστημα μάθησης που δεν στηρίζεται σε παραδείγματα αλλά στη γνώση του ειδικού;

Στην ιδανική περίπτωση, όπου τα παραδείγματα είναι αντιπροσωπευτικά των κλάσεων και δεν περιέχουν θόρυβο, η συνάρτηση πληροφορίας θα ήταν ικανή να εντοπίσει

κρυμμένες σχέσεις εξάρτησης και να τις χειριστεί σωστά. Ακόμα κι έτσι όμως, ο NewId σε περιπτώσεις ισοδύναμων ιδιοτήτων δεν μπορεί να εντοπίσει αυτή που είναι σημαντικότερη για την εφαρμογή. Διαλέγει τυχαία μία από αυτές γιατί όλες πετυχαίνουν το ίδιο καλό διαμερισμό των παραδειγμάτων με βάση τις κλάσεις. Με τη έννοια αυτή μία ιδιότητα μπορεί να έχει μεγάλη πληροφοριακή ισχύ αλλά ουσιαστικά να μην είναι ιδιαίτερα χρήσιμη για την επίλυση του προβλήματος. Για να αποδείξει κανείς την ύπαρξη τέτοιων περιπτώσεων μπορεί να εκτελέσει το παρακάτω πείραμα. Ας υποθέσουμε πως στις ιδιότητες που περιγράφουν την εφαρμογή προσθέτουμε μία ακόμα, που δεν σχετίζεται με την εφαρμογή αυτή. Σε όλα τα παραδείγματα της ίδιας κλάσης αντιστοιχούμε για αυτή την ιδιότητα μία μοναδική τιμή, ενώ παραδείγματα διαφορετικών κλάσεων παίρνουν διαφορετικές τιμές. Προφανώς, σύμφωνα με τη συνάρτηση υπολογισμού της πληροφορίας, η ιδιότητα αυτή είναι η πιο πληροφοριακή και επιλέγεται από τον αλγόριθμο στο πρώτο κιάλας βήμα σαν ρίζα του δέντρου απόφασης. Συνεπώς περιέχεται σε όλους τους κανόνες που παράγονται. Ο NewId δεν έχει κανένα τρόπο να εντοπίσει το γεγονός ότι η ιδιότητα αυτή ήταν άσχετη με την εφαρμογή και επομένως τη χρησιμοποίησε. Φυσικά το παράδειγμα αυτό είναι ακραίο. Σε φυσιολογικές καταστάσεις, το πεδίο εφαρμογής περιγράφουν ιδιότητες που είναι σχετικές με αυτό. Χρησιμοποιήθηκε όμως για να φανεί πόσο βασίζεται η συνάρτηση υπολογισμού της πληροφορίας στα δεδομένα ώστε να επιλέγει ακόμα και άσχετες ιδιότητες, αν αυτές έχουν μία καλή διασπορά τιμών ανάμεσα στις κλάσεις.

Στόχος της εργασίας αυτής ήταν η χρησιμοποίηση της γνώσης του ειδικού με τέτοιο τρόπο ώστε η εξαγωγή της γνώσης αυτής να είναι εύκολη. Όταν κανείς αναπτύσσει μία εφαρμογή παίρνει μία ‘συνέντευξη’ από τον ειδικό ώστε να εντοπίσει τα χαρακτηριστικά και τις ιδιότητες του πεδίου. Μέρος της γνώσης αυτής είναι οι εξαρτήσεις μεταξύ ιδιοτήτων. Μελετήθηκαν και υλοποιήθηκαν οι δυνατές μορφές εξάρτησης μεταξύ ιδιοτήτων με σκοπό την εισαγωγή δομής πάνω στην ‘επίπεδη’ μορφή αναπαράστασης δεδομένων που χρησιμοποιεί ο NewId. Η γνώση αυτή θα χρησιμοποιηθεί στην αντιμετώπιση καταστάσεων όπου δημιουργείται πρόβλημα λόγω δεδομένων. Με τον τρόπο αυτό θα εξάγεται το καλύτερο δυνατό σύνολο κανόνων.

Συνήθως το πρόβλημα έγκειται στη διατύπωση των σωστών ερωτήσεων που θα επιτρέψουν την εξαγωγή της πληροφορίας αυτής. Ο τυπικός ορισμός των σχέσεων εξάρτησης είναι το πρώτο βήμα στη διαδικασία αυτή. Ο εμπλουτισμός του NewId με μηχανισμούς χειρισμού τέτοιων εξαρτήσεων οδήγησε στην ανάπτυξη ενός νέου βελτιωμένου συστήματος, του IDDD(Inductive Domain Dependent Decision). Αν στην περιγραφή του πεδίου εφαρμογής δεν περιέχεται οποιοσδήποτε ορισμός εξάρτησης, η λειτουργία του IDDD ανάγεται σε αυτή του NewId. Από την μελέτη των σχέσεων που διέπουν πραγματικές εφαρμογές, εντοπίστηκαν δύο μορφές εξάρτησης μεταξύ ιδιοτήτων. Πρόκειται για την **απλή** εξάρτηση, η οποία ορίζεται με την πρόταση:

B εξαρτάται από την A (B depends_on A).

Η πρόταση αυτή δηλώνει πως η πληροφορία που αναπαριστά η ιδιότητα B είναι χρήσιμη μόνο σε συνδυασμό με την πληροφορία που αναπαριστά η A. Δηλαδή, η ιδιότητα B θα χρησιμοποιηθεί σε ένα κανόνα μόνο αν έχει ήδη χρησιμοποιηθεί η A.

Μία διαφορετική μορφή εξάρτησης είναι ο αποκλεισμός μιας ιδιότητας όταν στον κανόνα εμφανίζεται μία άλλη ιδιότητα με κάποια συγκεκριμένη τιμή. Για παράδειγμα, όταν $A = a_i$ απαγορεύεται η χρήση της ιδιότητας B. Ονομάζουμε τη σχέση αυτή **εξάρτηση αποκλεισμού** και την ορίζουμε χρησιμοποιώντας την ακόλουθη πρόταση:

B εξαρτάται από την A : a_i, \dots, a_j

όπου a_i, \dots, a_j , αναπαριστούν τιμές της Μάνας αν αυτή είναι κατηγορηματική και το διάστημα $[a_i, a_j]$ αν είναι αριθμητική ιδιότητα.

Ένας περιορισμός που τίθεται από τον IDDD είναι πως μία ιδιότητα που είναι Κόρη σε κάποια εξάρτηση δεν μπορεί να είναι Κόρη και σε κάποια άλλη εξάρτηση. Δεν μπορεί δηλαδή να έχει περισσότερες από μία Μάνες. Ο περιορισμός αυτός δεν ισχύει για τις Μάνες οι οποίες μπορεί να έχουν περισσότερες από μία Κόρες. Η ύπαρξη μίας μόνο Μάνας για κάθε Κόρη απλοποιεί τον χειρισμό της σχέσης εξάρτησης. Η Κόρη θα χρησιμοποιηθεί μόνο αν έχει ήδη χρησιμοποιηθεί η Μάνα. Αν η Μάνα δεν υπάρχει ήδη στον κανόνα ο αλγόριθμος οδηγείται στην υποχρεωτική επιλογή της. Όταν υπάρχουν περισσότερες από μία Μάνες πρέπει να εξεταστεί η σειρά επιλογής των ιδιοτήτων που

είναι Μάνες και δεν εμφανίζονται στον κανόνα. Η άρση αυτού του περιορισμού και η μελέτη των πρακτικών προβλημάτων που αυτή δημιουργεί αποτελεί αντικείμενο περαιτέρω έρευνας.

Ο NewId είναι ικανός να χειριστεί σχέσεις, όπως αυτές που ορίστηκαν παραπάνω, υπό ορισμένες προϋποθέσεις. Ο ορισμός διάταξης μεταξύ των ιδιοτήτων φαίνεται ότι εκφράζει την ίδια πληροφορία με αυτή της απλής εξάρτησης. Πραγματικά, αν ορίσει κανείς πως η ιδιότητα Μάνα πρέπει να προηγείται της Κόρης επιτυγχάνει την επιλογή της Κόρης μόνο όταν η Μάνα έχει ήδη επιλεγεί. Στα πλαίσια της εργασίας αυτής επιχειρήθηκε και η συγκριτική αξιολόγηση του ορισμού των δύο σχέσεων. Η απόδοση της σχέσης διάταξης εξακολουθεί να επηρεάζεται από τα εκπαιδευτικά παραδείγματα με αποτέλεσμα να αποτυγχάνει σε μερικές περιπτώσεις. Τα πειραματικά αποτελέσματα επιβεβαιώνουν τον παραπάνω ισχυρισμό.

Η χρησιμοποίηση της αδιάφορης τιμής για την Κόρη-ιδιότητα στα παραδείγματα που η Μάνα έχει την απαγορευτική τιμή θα μπορούσε να εκφράσει την εξάρτηση αποκλεισμού. Όμως, ακόμα και με τη χρησιμοποίηση αυτών των μηχανισμών, το δέντρο απόφασης δεν παίρνει πάντα υπόψην του τις σχέσεις εξάρτησης μεταξύ των ιδιοτήτων. Ο τρόπος ορισμού των σχέσεων αυτών είναι έμμεσος, μέσω των παραδειγμάτων εκπαίδευσης, τα οποία δεν είναι πάντα εύκολο να ελέγξουμε για την ορθότητά τους.

Για να γίνει πιο κατανοητή η παραπάνω παρατήρηση ας θεωρήσουμε την περίπτωση όπου όταν σε μία ιδιότητα εμφανίζεται κάποια συγκεκριμένη τιμή σημαίνει ότι γίνεται απαγορευτική η χρήση κάποιας άλλης ιδιότητας. Στα παραδείγματα όπου η Μάνα έχει την συγκεκριμένη τιμή δίνουμε στην Κόρη την αδιάφορη τιμή για να υποδηλώσουμε ότι δεν θέλουμε να χρησιμοποιηθεί. Η συχνή εμφάνιση της αδιάφορης τιμής έχει σαν αποτέλεσμα την μείωση της πληροφοριακής ισχύος της αντίστοιχης ιδιότητας. Σύμφωνα με τα παραπάνω δεν θα έπρεπε να δημιουργείται πρόβλημα. Δυστυχώς η σωστή λειτουργία εξαρτάται από τη σειρά επιλογής των ιδιοτήτων αλλά και από τη πληροφοριακή ισχύ της Μάνας σε κάποιο τυχαίο βήμα. Πραγματικά, αν επιλεγεί η Μάνα, τα παραδείγματα χωρίζονται με τέτοιο τρόπο ώστε, στο υποσύνολο που αντιστοιχεί στην απαγορευτική τιμή, η ιδιότητα Κόρη να έχει μικρή πληροφοριακή ισχύ. Το πρόβλημα εμφανίζεται στις πε-

ριπτώσεις όπου σε μερικά παραδείγματα η Μάνα έχει άγνωστη τιμή. Δεν είναι δυνατόν τότε να αντιστοιχίσουμε στην Κόρη την αδιάφορη τιμή. Παρουσιάζεται έτσι ένα σύνολο παραδειγμάτων που απαιτεί κάλυψη και τίποτα δεν εμποδίζει τον αλγόριθμο να οδηγηθεί σε επιλογή της Κόρης. Η ιδιότητα Μάνα μπορεί να έχει επιλεγεί σε προηγούμενο βήμα και να παίρνει κάποια από τις απαγορευτικές τιμές ή να μην έχει επιλεγεί καθόλου. Και στις δύο περιπτώσεις το αποτέλεσμα είναι ασύμβατο με την πραγματικότητα της εφαρμογής. Όσο αυξάνει το πλήθος των παραδειγμάτων με άγνωστη τιμή για τη Μάνα τόσο πιο έντονο είναι το πρόβλημα. Εύκολα μπορεί να φανταστεί κανείς τι συμβαίνει στις περιπτώσεις όπου η Μάνα έχει άγνωστη τιμή γιατί δεν ενδιαφέρει τον ειδικό. Η Κόρη τότε θα επιλεγεί άσχετα αν ο ρόλος της είναι ακόμα πιο δευτερεύων.

Επιπλέον, όταν η Μάνα είναι αριθμητική ιδιότητα, η χρησιμοποίηση της αδιάφορης τιμής για την Κόρη δεν μπορεί να εκφράσει την πληροφορία για την ύπαρξη της αποκλειστικής εξάρτησης. Κατά την επιλογή της Μάνας καθορίζεται κάποιο σημείο διαχωρισμού και τίποτα δεν εγγυάται πως το σημείο αυτό θα είναι κάποιο από τις απαγορευτικές τιμές. Δεν υπάρχει τρόπος να ξέρουμε τα υποσύνολα των παραδειγμάτων που θα σχηματιστούν και αν υπάρχει περίπτωση επιλογής της Κόρης σε κανόνα όπου η Μάνα παίρνει κάποιες απαγορευτικές τιμές.

Σε αντίθεση με τον NewId που χειρίζεται έμμεσα τις σχέσεις εξάρτησης, στον IDDD επιδιώκουμε την μορφοποίηση ενός μοντέλου που να τις εκφράζει. Οι σχέσεις αυτές δεν εκφράζονται μέσω των παραδειγμάτων αλλά μέσω γενικότερων προτάσεων. Αποτελούν ανεξάρτητη πηγή πληροφορίας, δεν εξαρτώνται από τα παραδείγματα και επιδιώκουν την εξαγωγή κανόνων που θα είναι πιο κοντά σε αυτούς που χρησιμοποιούν οι ειδικοί. Το μοντέλο των σχέσεων χρησιμοποιείται κατά την επιλογή της ιδιότητας που θα αποτελέσει κόμβο στο δέντρο απόφασης. Η πληροφοριακή ισχύς κάθε ιδιότητας εξακολουθεί να είναι το βασικό κριτήριο για την επιλογή της σε ένα κόμβο. Η επιλογή ιδιοτήτων που δεν μετέχουν σε κάποια σχέση εξάρτησης δεν επηρεάζεται. Ωστόσο τίθενται περιορισμοί στην επιλογή των ιδιοτήτων που αποτελούν παιδιά σε μία σχέση. Για τις ιδιότητες αυτές, η συνάρτηση υπολογισμού της πληροφοριακής ισχύος των ιδιοτήτων χρησιμοποιείται σαν κριτήριο ενεργοποίησης ή όχι της σχέσης. Αν σε κάποιο βήμα της κατασκευής του

δέντρου εμφανίζεται σαν πιο πληροφοριακή ιδιότητα μία Κόρη, η ικανοποίηση των περιορισμών εξάρτησης είναι αυτή που θα καθορίσει αν πραγματικά επιλεγεί η ιδιότητα αυτή ή όχι.

Μία πρώτη ματιά στις σχέσεις εξάρτησης οδηγεί στην παρατήρηση πως η απλή εξάρτηση αποτελεί ειδική μορφή της αποκλειστικής εξάρτησης. Στην απλή εξάρτηση μία ιδιότητα εξαρτάται από κάποια άλλη, αλλά η επιλογή της Κόρης δεν απαγορεύεται από κάποια τιμή της Μάνας. Έχει περισσότερο την έννοια της προτεραιότητας της Μάνας έναντι της Κόρης. Στην αποκλειστική εξάρτηση καθορίζονται επιπλέον μία ή περισσότερες απαγορευτικές τιμές και επομένως οι σχέσεις είναι πιο ισχυρές και πολύπλοκες. Οπως θα φανεί και από την αναλυτική περιγραφή των σχέσεων, στις παραγράφους 4.1 και 4.2, ο καθορισμός απαγορευτικών τιμών δίνει επιπλέον δυνατότητες χρησιμοποίησης της γνώσης στην διαδικασία κατασκευής του δέντρου απόφασης.

Υπενθυμίζουμε πως οι ιδιότητες διακρίνονται σε κατηγορηματικές και αριθμητικές. Ο ορισμός και υλοποίηση των σχέσεων εξάρτησης είναι συνάρτηση και του τύπου της ιδιότητας Μάνα. Στις παραγράφους 4.1 και 4.2 περιγράφονται οι σχέσεις αυτές και οι διαφοροποιήσεις μεταξύ των δύο τύπων που μπορεί να έχει η Μάνα-ιδιότητα.

4.1 Απλή Εξάρτηση

Ας θεωρήσουμε πάλι τη σχέση απλής εξάρτησης μεταξύ δύο ιδιοτήτων A και B.

B εξαρτάται από την A

Όπως έχουμε ήδη αναφέρει, η σχέση αυτή μπορεί να οριστεί στον NewId χρησιμοποιώντας μία σχέση διάταξης. Η ισοδύναμη έκφραση είναι:

Η ιδιότητα A ΠΡΙΝ ΤΗΝ B

Σύμφωνα με την προσέγγιση που ακολουθεί ο NewId, η ιδιότητα B δεν μπορεί να επιλεγεί αν προηγουμένως δεν έχει επιλεγεί η A. Ο περιορισμός αυτός ισχύει ακόμα και αν, σε κάποιο βήμα, η B εμφανίζεται σαν η πιο πληροφοριακή ιδιότητα.

Η προσέγγιση που ακολουθείται στον IDDD είναι διαφορετική. Η ιδιότητα B δεν μπορεί και πάλι να επιλεγεί αν σε προηγούμενο βήμα δεν έχει επιλεγεί η A. Ωστόσο δεν μπορεί να αγνοηθεί η πληροφορία πως η B έχει μεγάλη πληροφοριακή ισχύ. Με βάση αυτή τη παρατήρηση, περιγράφουμε την υλοποίηση της σχέσης διαφοροποιώντας ανάμεσα σε δύο περιπτώσεις:

1. Η A εμφανίζεται σαν η πιο πληροφοριακή ιδιότητα. Στην περίπτωση αυτή η σχέση εξάρτησης δεν χρειάζεται να ενεργοποιηθεί. Η διαδικασία δημιουργίας του κόμβου και χωρισμού των παραδειγμάτων σε υποσύνολα παραμένει αναλλοίωτη. Ωστόσο, η επιλογή της A σημαίνει ότι, για το παρόν μονοπάτι, μπορεί να επιλεγεί η B σε κάποιο επόμενο βήμα. Για τη διατήρηση της πληροφορίας αυτής, δημιουργούμε για κάθε μονοπάτι μία λίστα. Κάθε στοιχείο της λίστας περιέχει πληροφορία για την ιδιότητα που επιλέχτηκε, την πληροφοριακή της ισχύ και την τιμή που έχει για το μονοπάτι. Αν η ιδιότητα ήταν αριθμητική, η πληροφορία για την τιμή της αφορά στο συγκεκριμένο σημείο διαχωρισμού αλλά και στο αν βρισκόμαστε σε μονοπάτι για τιμές 'μικρότερες ή ίσες' ή 'μεγαλύτερες' από αυτό. Η λίστα αυτή ονομάζεται **ιστορικό** του μονοπατιού.
2. Η B εμφανίζεται σαν η πιο πληροφοριακή ιδιότητα. Στην περίπτωση αυτή ελέγχε-

ται το μέχρι εκείνη τη στιγμή ιστορικό του μονοπατιού. Αν, σύμφωνα με αυτό, η A έχει ήδη επιλεγεί επιτρέπουμε την επιλογή της B . Στην αντίθετη περίπτωση, λαμβάνοντας υπόψη την μεγάλη πληροφοριακή ισχύ της B οδηγούμε τον αλγόριθμο σε επιλογή της A . Το σκεπτικό είναι πως, αν όντως η B περιέχει μεγάλη πληροφορία, θα επιλεγεί σε επόμενο βήμα. Με την προσέγγιση αυτή ο αλγόριθμος είναι συνεπής με τη γνώση του πεδίου, όπως αυτή εκφράζεται από την σχέση εξάρτησης, και παράλληλα εξακολουθεί να διαλέγει ιδιότητες με μεγάλη πληροφοριακή ισχύ.

Σημειώνουμε εδώ πως, όσον αφορά την απλή εξάρτηση, δεν υπάρχει διαφορά στον χειρισμό κατηγορηματικών και αριθμητικών ιδιοτήτων.

4.2 Εξάρτηση Αποκλεισμού

Σε αντίθεση με την απλή εξάρτηση, ο NewId δεν διαθέτει μηχανισμούς χειρισμού της εξάρτησης αποκλεισμού. Ωστόσο, με τον ορισμό μιας τέτοιας σχέσης μπορεί κανείς να ορίσει πιο πολύπλοκες σχέσεις μεταξύ ιδιοτήτων. Ορίζοντας ότι

B εξαρτάται από την $A : a_i, \dots, a_j$

απαγορεύουμε την επιλογή της B όταν η A παίρνει τιμές από το υποσύνολο $\{a_i, \dots, a_j\}$.

Ο τρόπος υλοποίησης της εξάρτησης αποκλεισμού διαφέρει ανάλογα με τον τύπο της ιδιότητας Μάνα. Όταν η Μάνα είναι κατηγορηματική ιδιότητα, ο χειρισμός της σχέσης είναι όμοιος με αυτόν της απλής εξάρτησης αλλά προβλέπει και την ύπαρξη των απαγορευτικών τιμών. Έτσι, η σχέση εξάρτησης εξετάζεται όταν η B εμφανίζεται σαν η πιο διαχωριστική ιδιότητα. Η επιλογή της ιδιότητας B απαγορεύεται όταν η A παίρνει κάποια από τις τιμές a_i, \dots, a_j . Στο ιστορικό του μονοπατιού υπάρχει η πληροφορία για την επιλογή της Μάνας καθώς και της τιμής που έχει στο τρέχον μονοπάτι. Αν η Μάνα δεν έχει επιλεγεί, όπως και στην απλή εξάρτηση, δημιουργούμε ένα κόμβο για αυτή. Αν έχει ήδη επιλεγεί σε προηγούμενο βήμα, εξετάζουμε την τιμή που έχει για το μονοπάτι που κατασκευάζουμε. Αν αυτή η τιμή είναι μία από τις απαγορευτικές, η ιδιότητα B πάει να θεωρείται υποψήφια για τον κόμβο, αν όχι θα δημιουργηθεί ένας κόμβος γι' αυτήν καθώς είναι η πιο πληροφοριακή ιδιότητα.

Όταν η ιδιότητα Μάνα είναι αριθμητική και πρόκειται να αποτελέσει κόμβο του δέντρου, πρέπει να επιλεγεί το σημείο διαχωρισμού. Η επιλογή του σημείου αυτού μπορεί να υποβοηθηθεί από τη γνώση που περιέχεται στη δήλωση της εξάρτησης αποκλεισμού. Για τις αριθμητικές ιδιότητες, ο NewId χρησιμοποιεί στατιστικά μέτρα για τον υπολογισμό του σημείου διαχωρισμού. Με βάση αυτό το σημείο, δημιουργεί ένα κόμβο με δύο κλαδιά. Στο πρώτο κλαδί αντιστοιχίζεται εκείνο το υποσύνολο των παραδειγμάτων για τα οποία η τιμή της ιδιότητας είναι 'μικρότερη ή ίση' του σημείου διαχωρισμού. Στο δεύτερο κλαδί αντιστοιχίζονται τα παραδείγματα με τιμή ιδιότητας 'μεγαλύτερη' του σημείου διαχωρισμού. Ωστόσο τα στατιστικά μέτρα που χρησιμοποιεί ο NewId είναι ευαίσθητα στο 'θόρυβο' που υπάρχει στις τιμές των παραδειγμάτων. Τείνουν να κληρονομήν τον

‘θόρυβο’ αυτό με αποτέλεσμα τα σημεία διαχωρισμού για τις αριθμητικές ιδιότητες που εμφανίζονται στο δέντρο να μην έχουν νόημα για τους ειδικούς της εφαρμογής.

Στον IDDD, θεωρούμε πως ο ορισμός μίας συγκεκριμένης τιμής στην εξάρτηση αποκλεισμού σημαίνει πως η τιμή αυτή είναι σημαντική για την εφαρμογή. Συγκεκριμένα η εξάρτηση αποκλεισμού παίρνει, για αριθμητικού τύπου μάνες, μία από τις ακόλουθες μορφές:

1. B εξαρτάται από την A : $\leq a_i$

Αν η A εμφανίζεται στον κανόνα (μονοπάτι) με τιμή ‘μικρότερη ή ίση’ του a_i η B δεν μπορεί να αποτελέσει κόμβο.

2. B εξαρτάται από την A : $> a_i$

Αν η A εμφανίζεται στον κανόνα με τιμή ‘μεγαλύτερη’ του a_i η B δεν μπορεί να αποτελέσει κόμβο.

3. B εξαρτάται από την A : $[a_i, a_j]$

Αν η A εμφανίζεται στον κανόνα με τιμή στο διάστημα $[a_i, a_j]$ απαγορεύουμε την επιλογή της B.

Κάθε μία από τις παραπάνω μορφές της εξάρτησης αποκλεισμού απαιτεί διαφορετικό χειρισμό όσον αφορά στον καθορισμό του σημείου διαχωρισμού αλλά και στην διαδικασία δημιουργίας του δέντρου.

Στις δύο πρώτες μορφές, όταν η B εμφανίζει μεγάλη πληροφοριακή ισχύ δημιουργούμε ένα κόμβο για την A με σημείο διαχωρισμού το a_i και η επιλογή της B απαγορεύεται στο κλαδί που αντιστοιχεί σε τιμές ‘μικρότερες ή ίσες’, ή ‘μεγαλύτερες’, του a_i αντίστοιχα.

Στην τρίτη περίπτωση, όταν η B εμφανίζει μεγάλη πληροφοριακή ισχύ δημιουργούμε αρχικά ένα κόμβο για την A με σημείο διαχωρισμού το a_i και κατόπιν έναν ακόμη κόμβο για την A με σημείο διαχωρισμού το a_j . Με τον τρόπο αυτό δημιουργείται ένα μονοπάτι με τιμές της A ‘μεγαλύτερες’ του a_i και ‘μικρότερες ή ίσες’ του a_j . Στο μονοπάτι αυτό απαγορεύουμε, σε μετέπειτα βήματα, την επιλογή της B.

Οι αριθμητικές τιμές που ορίζονται στη σχέση εξάρτησης χρησιμοποιούνται μόνο στην περίπτωση που η Μάνα γίνεται κόμβος λόγω της σχέσης. Αυτή η προσέγγιση ακολουθείται με το σκεπτικό πως οι τιμές αυτές της Μάνας είναι μεν σημαντικές για την εφαρμογή αλλά μόνο σε συνδυασμό με την ιδιότητα κόρη κι όχι αυτόνομα. Όταν στο ιστορικό του μονοπατιού εμφανίζεται η Μάνα λόγω της πληροφοριακής της ισχύς, και όχι λόγω κάποιας εξάρτησης, δεν μπορούμε να αποφασίσουμε αν θα πρέπει να επιλέξουμε την κόρη ή όχι. Λόγω της απαγορευτικής τιμής που ορίζεται για την αποκλειστική εξάρτηση, δημιουργούνται δύο διαστήματα τιμών της Μάνας για τα οποία απαγορεύεται η επιτρέπεται η επιλογή της Κόρης. Όταν όμως η Μάνα έχει ήδη επιλεγεί λόγω της πληροφοριακής της ισχύος το σημείο διαχωρισμού της θα είναι, στη γενική περίπτωση, διαφορετικό της απαγορευτικής τιμής. Επομένως δημιουργούνται δύο άλλα διαστήματα τιμών της Μάνας για τα οποία δεν μπορεί να καθοριστεί αν επιτρέπεται, ή όχι, η επιλογή της Κόρης.

Για παράδειγμα, αν έχει οριστεί η εξάρτηση:

B εξαρτάται από την A : $> a_i$

και η A έχει ήδη επιλεγεί με σημείο διαχωρισμού a_j με $a_i \neq a_j$, εμφανίζονται τα εξής ενδεχόμενα:

1. $a_j \leq a_i$

Στο κλαδί που αντιστοιχεί σε τιμές 'μικρότερες ή ίσες' του a_j επιτρέπεται η επιλογή της B. Στο κλαδί για τιμές 'μεγαλύτερες' του a_j υπάρχουν κάποια παραδείγματα με τιμές στο διάστημα $[a_j, a_i]$ αλλά και κάποια άλλα με τιμές μεγαλύτερες του a_i . Για το πρώτο υποσύνολο θα μπορούσε να επιλεγεί η B, ενώ για το δεύτερο η επιλογή της απαγορεύεται λόγω της σχέσης εξάρτησης. Γενικότερα, όσο αφορά το κλαδί αυτό, δεν μπορούμε να πάρουμε απόφαση επιλογής ή όχι της B, αν θέλουμε να είμαστε συνεπείς με τον περιορισμό που θέτει η εξάρτηση. Στο σχήμα 3 φαίνεται η κατάσταση όπως διαμορφώνεται για την περίπτωση αυτή. Κάτω από τον άξονα είναι τα διαστήματα όπου απαγορεύεται η επιλογή της B λόγω σχέσης. Πάνω από τον άξονα φαίνονται τα διαστήματα που σχηματίζονται από την επιλογή του a_j σαν σημείο διαχωρισμού.

Σχήμα 3: B εξαρτάται από την A : $> a_i$ και $a_j \leq a_i$

2. $a_j > a_i$

Στο κλαδί που αντιστοιχεί σε τιμές ‘μεγαλύτερες’ του a_j είναι ξεκάθαρο πως δεν μπορούμε να επιλέξουμε την B. Στο κλαδί για τιμές ‘μικρότερες ή ίσες’ τα παραδείγματα παίρνουν τιμές για την A είτε στο διάστημα $[a_i, a_j]$ ή μικρότερες του a_i . Για το πρώτο υποσύνολο παραδειγμάτων απαγορεύεται η χρήση της B, ενώ για το δεύτερο δεν ισχύει τέτοιος περιορισμός. Επομένως σε αυτό το κλαδί δεν μπορεί να ληφθεί απόφαση επιλογής ή όχι της B. Στο σχήμα 4 φαίνεται η κατάσταση όπως διαμορφώνεται για την περίπτωση αυτή.

Σχήμα 4: B εξαρτάται από την A : $> a_i$ και $a_j > a_i$

Για να αποφύγουμε αυτή την αβεβαιότητα ακολουθούμε την εξής προσέγγιση: Όταν

η Μάνα έχει επιλεγεί σε προηγούμενο βήμα με άλλο σημείο διαχωρισμού και η Κόρη είναι η πιο πληροφοριακή ιδιότητα στο παρόν βήμα, αποτρέπουμε την επιλογή της Κόρης και δημιουργούμε ένα κόμβο για την Μάνα καθορίζοντας το σημείο διαχωρισμού αναλόγως με τη συγκεκριμένη μορφή της εξάρτησης αποκλεισμού.

4.3 Διαδοχικές Εξαρτήσεις

Οι ιδιότητες που είναι μέλη μιας συγκεκριμένης εξάρτησης είναι δυνατόν να μετέχουν και σε άλλες εξαρτήσεις. Για παράδειγμα, στην περιγραφή της εφαρμογής επιτρέπεται ο ορισμός σχέσεων της μορφής

A depends_on B

B depends_on C

Σύμφωνα με τις σχέσεις αυτές ορίζεται μία έμμεση εξάρτηση της A από την B μέσω της C. Οι σχέσεις αυτές είναι αποδεκτές εφόσον δεν δημιουργείται κύκλος. Δεν μπορεί δηλαδή μία ιδιότητα να εξαρτάται από κάποια άλλη η οποία με τη σειρά της εξαρτάται από αυτή. Παράδειγμα τέτοιου κύκλου αποτελούν οι σχέσεις

A depends_on B

B depends_on A

Επομένως μία ιδιότητα μπορεί να έχει σαν Μάνα μία άλλη ιδιότητα που είναι, με τη σειρά της, Κόρη κάποιας άλλης. Θεωρούμε ξανά το ζευγάρι εξαρτήσεων

A depends_on B

B depends_on C

Αν σε κάποιο βήμα εμφανίζεται η A σαν η πιο πληροφοριακή ιδιότητα, ελέγχεται η ύπαρξη στο μονοπάτι της B. Αν μεν αυτή έχει επιλεγεί σε προηγούμενο βήμα, η τελική επιλογή της A εξαρτάται μόνο από την ύπαρξη ή όχι απαγορευτικής τιμής. Αν όμως η B δεν έχει επιλεγεί, θα πρέπει κανονικά να δημιουργηθεί κόμβος γι'αυτήν. Επειδή η B εξαρτάται από την C, επαναλαμβάνουμε τον έλεγχο επιλογής της C. Συμπερασματικά σημειώνουμε πως η επιλογή μιας ιδιότητας προϋποθέτει την επιλογή όλων των ιδιοτήτων που είναι Μάνες της έμμεσα ή άμεσα.

4.4 Πολυπλοκότητα του αλγορίθμου IDDD

Κατά την ανάπτυξη του IDDD μεταβάλλεται η διαδικασία επιλογής κόμβου του NewId. Οι μεταβολές αυτές έχουν σαν αποτέλεσμα:

Όταν στο μονοπάτι έχει ήδη επιλεγεί η Μάνα με απαγορευτική τιμή δεν υπολογίζεται η πληροφοριακή ισχύς της Κόρης καθώς αυτή απαγορεύεται να επιλεγεί. Η μέτρηση της πληροφοριακής ισχύος απαιτεί την εξέταση των τιμών της ιδιότητας στα παραδείγματα και την δημιουργία μιας τεράστιας δομής δεδομένων. Αυτή η δομή είναι ένας πίνακας με μέγεθος ίσο με το γινόμενο: (Τιμές ιδιότητας) x (Αριθμός κλάσεων). Ο υπολογισμός της πληροφοριακής ισχύος αντικαθίσταται από έναν έλεγχο του ιστορικού του μονοπατιού, δηλαδή μιας λίστας μικρού πλήθους στοιχείων (συνήθως 3-4 στοιχεία).

Επιπλέον, όταν η Μάνα σε μία αποκλειστική εξάρτηση είναι αριθμητική, δεν υπολογίζεται το σημείο διαχωρισμού αλλά τίθεται ίσο με την απαγορευτική τιμή.

Ο αλγόριθμος επιβαρύνεται στις περιπτώσεις όπου η Κόρη είναι όντως η πιο πληροφοριακή ακόμα και μετά από την επιλογή της Μάνας. Τότε η παρεμβολή του κόμβου για τη Μάνα στοιχίζει στον αλγόριθμο τόσες αναδρομές όσες και οι τιμές της. Επιπλέον, πριν την επιλογή οποιασδήποτε Κόρης, ελέγχεται το ιστορικό του μονοπατιού για την ύπαρξη της Μάνας.

Συνολικά οι μεταβολές αυτές δεν επηρεάζουν σημαντικά το χρόνο εκτέλεσης του αλγορίθμου. Οι διαφορές είναι πολύ μικρές και εξαρτώνται από τα δεδομένα όπως άλλωστε και για τον αλγόριθμο NewId. Η πολυπλοκότητα είναι $O(|EX| \times |A| \times |\text{ΕσωτερικοίΚόμβοι}|)$, όπου $|EX|$ ο αριθμός παραδειγμάτων και $|A|$ ο αριθμός ιδιοτήτων.

5 Πειραματικά Αποτελέσματα

Στα πλαίσια του MLT (Machine Learning Toolbox) project, αναπτύχθηκε μία πραγματική εφαρμογή που αφορά στην θεραπευτική απόφαση για την Κρυψορχία. Στην παράγραφο 5.1 ακολουθεί η περιγραφή του πεδίου εφαρμογής και τα ιδιαίτερα χαρακτηριστικά του. Για την εφαρμογή αυτή χρησιμοποιήθηκαν οι αλγόριθμοι NewId και IDDD των οποίων τα πειραματικά αποτελέσματα παρατίθενται και αξιολογούνται στη παράγραφο 5.4.

5.1 Περιγραφή Πεδίου Εφαρμογής

Η Κρυψορχία είναι εκείνη η παθολογική κατάσταση των αρσενικών παιδιών στην οποία ο ένας, ή και οι δύο όρχεις, δεν έχουν κατέλθει στη φυσιολογική τους θέση μετά τη γέννηση. Αν δεν θεραπευτεί η κατάσταση αυτή, μπορεί να οδηγήσει σε ατροφικούς όρχεις οι οποίοι σε μεγαλύτερη ηλικία μπορεί να εξελιχτούν σε καρκινώματα. Αν η κατάσταση αυτή διαγνωστεί στους πρώτους έξι μήνες της ζωής του παιδιού, η ιδανική θεραπεία συνίσταται στη χορήγηση ορμονών μετά τον έκτο μήνα, παρακολούθηση για 6-8 εβδομάδες και εγχείρηση αν η ορμονοθεραπεία αποδειχτεί ανεπιτυχής. Στην περίπτωση που η Κρυψορχία διαγνωστεί αργότερα, δεν υπάρχει κάποιο γενικά αποδεκτό και συγκεκριμένο πρωτόκολο θεραπείας [30].

Δυστυχώς είναι πολλές οι περιπτώσεις μη έγκαιρης διάγνωσης λόγω κοινωνικών προκαταλήψεων. Η εμπειρία καθοδηγεί τον γιατρό στην αντιμετώπιση αυτών των περιπτώσεων. Ο στόχος ανάπτυξης αυτής της εφαρμογής ήταν η δημιουργία ενός συνόλου κανόνων σύμφωνα με τους οποίους θα αντιστοιχίζεται, σε κάθε περιστατικό, μία από τις δυνατές θεραπευτικές αποφάσεις (κλάσεις):

1. Παρακολούθηση (follow-up)
2. Ορμονοθεραπεία (hormonal)
3. Εγχείρηση (surgical)

Οι ιδιότητες που περιγράφουν κάθε περιστατικό (παράδειγμα) είναι (στις παρενθέσεις βρίσκονται τα ονόματά τους):

1. Ο αριθμός της επίσκεψης (visit). Κάθε ασθενής επισκέπτεται τον γιατρό περισσότερες από μία φορές και η θεραπευτική απόφαση είναι διαφορετική σε κάθε επίσκεψη. Επιπλέον, όσο περισσότερες φορές έχει έρθει ο ασθενής τόσο χρόνιο γίνεται το πρόβλημά του και απαιτεί πιο δραστική αντιμετώπιση. Η ιδιότητα αυτή αναπαριστάται σαν μία κατηγορηματική και όχι αριθμητική μεταβλητή γιατί δεν επιθυμούμε να υπάρχουν στους κανόνες σημεία διαχωρισμού για την επίσκεψη όπως π.χ 3,5.
2. Η ηλικία του παιδιού (age). Η ιδιότητα αυτή είναι επίσης κατηγορηματική και παίρνει σαν τιμές διαστήματα ηλικιών, π.χ 18-24μηνών.
3. Αιτιολογία (etiology). Η Κρυπορχία μπορεί να είναι μία εκ γενετής κατάσταση η οποία δεν αντιμετωπίστηκε ή να προηγήθηκε εγχείρηση η οποία όμως δεν ήταν επιτυχής. Η ιδιότητα αυτή είναι κατηγορηματικού τύπου.
4. Η κατάσταση του ενός όρχη (one). (Κατηγορηματική ιδιότητα).
5. Η κατάσταση του άλλου όρχη (two). (Κατηγορηματική ιδιότητα).
6. Το μέγεθος του ενός όρχη (size-of-one). (Αριθμητική ιδιότητα).
7. Το μέγεθος του άλλου όρχη (size-of-two). (Αριθμητική ιδιότητα).

Για τη λήψη οποιασδήποτε απόφασης είναι αναγκαία η γνώση για την κατάσταση των όρχεων στην προηγούμενη επίσκεψη καθώς και η πληροφορία για το αν ο ασθενής έχει ήδη πάρει ορμόνες ή όχι. Επειδή δεν υπάρχει η δυνατότητα σύνδεσης των διαφορετικών επισκέψεων ενός ασθενή, προσθέτουμε στην περιγραφή της εφαρμογής και τις ακόλουθες ιδιότητες:

1. Προηγούμενη κατάσταση του ενός όρχη (pre-one).
2. Προηγούμενη κατάσταση του άλλου όρχη (pre-two).

3. Λήψη ορμονών (treatment).

Στη συνέχεια δίνουμε ένα παράδειγμα του αρχείου περιγραφής της εφαρμογής το οποίο δέχονται σαν είσοδο οι NewId και IDDD. Δίπλα σε κάθε ιδιότητα υπάρχει η λίστα των δυνατών τιμών που μπορεί να πάρει.

```
**ATTRIBUTE FILE**
```

```
age : (ORDERED) "0-6months" "6-12months" "12-18months" "18-24months"  
"2-4years" "4-6years" "6-8years" "8-10years" "10-12years" "12-14years";
```

```
visit: (ORDERED) first second third fourth fifth sixth;
```

```
etiology : congenital postoperative;
```

```
one : normal not-palpable inguinal sliding retractile;
```

```
two : normal not-palpable inguinal sliding retractile;
```

```
pre-one : normal not-palpable inguinal sliding retractile;
```

```
pre-two : normal not-palpable inguinal sliding retractile;
```

```
size-of-one : (FLOAT)
```

```
size-of-two : (FLOAT)
```

```
treatment : no hormonal;
```

```
class : follow-up surgical hormonal;
```

5.2 Προεπεξεργασία των Παραδειγμάτων Εκπαίδευσης

Η ποιότητα των παραδειγμάτων αποτελεί τον βασικότερο παράγοντα για την καλή λειτουργία οποιουδήποτε αλγόριθμου μάθησης μέσω παραδειγμάτων. Τα παραδείγματα αυτά καθορίζουν ποιες ιδιότητες και σε ποια σειρά θα συνθέσουν τους κανόνες. Για το λόγο αυτό είναι βασικό να ελαχιστοποιηθεί ο θόρυβος που περιέχουν. Ο θόρυβος αυτός προέρχεται από λανθασμένη μέτρηση/εκτίμηση της τιμής μιας ιδιότητας αλλά και από την ύπαρξη εξωτερικών παραγόντων που επηρεάζουν τα δεδομένα της εφαρμογής χωρίς να σχετίζονται άμεσα με αυτή.

Η διόρθωση μιας λανθασμένης μέτρησης είναι μία πολύπλοκη διαδικασία η οποία απαιτεί την εξέταση των στοιχείων από τον ειδικό και επομένως είναι πολύ δύσκολο να αυτοματοποιηθεί.

Τα δεδομένα μιας εφαρμογής μπορεί να επηρεάζονται από παράγοντες που δεν σχετίζονται με αυτή και επομένως δεν μπορούν να καταγραφούν. Για παράδειγμα, στην εφαρμογή που μελετήσαμε στη παράγραφο 5.1 η απόφαση του γιατρού επηρεάζεται από την εμπιστοσύνη που έχει στους γονείς του παιδιού. Ενδεχομένως η σωστή θεραπευτική απόφαση σε μία περίπτωση είναι η παρακολούθηση αλλά ο γιατρός επιλέγει τελικά την εγχείρηση όταν βλέπει ότι οι γονείς δυσανασχετούν με την προοπτική μιας νέας επίσκεψης στο Νοσοκομείο. Στόχος μας είναι η εξαγωγή ενός συνόλου κανόνων που θα περιγράφουν την εφαρμογή κι επομένως δεν θα θέλαμε οι παραπάνω παράγοντες να αντανakλώνται στο σύνολο αυτό.

Κατά την ανάπτυξη μιας εφαρμογής, είναι δυνατόν ο χρήστης να μην έχει την δυνατότητα ελέγχου των παραδειγμάτων εκπαίδευσης είτε γιατί αυτά έχουν συγκεντρωθεί από κάποια αρχεία, η πρόσβαση στα οποία είναι πια αδύνατη, είτε γιατί ο έλεγχος τους απαιτεί την εξέταση χιλιάδων στοιχείων και επομένως χρόνο και χρήμα.

Έτσι, στο σύνολο των παραδειγμάτων μπορεί να υπάρχει το ίδιο παράδειγμα περισσότερες από μία φορές ή στη χειρότερη περίπτωση, παραδείγματα που έχουν τις ίδιες ακριβώς τιμές για όλες τις ιδιότητες αλλά ανήκουν σε διαφορετική κλάση. Θα μπορούσε κανείς να θεωρήσει πως η πρώτη περίπτωση επιβαρύνει τον αλγόριθμο με ένα

επιπρόσθετο κόστος χωρίς να του προσφέρει ουσιαστική πληροφορία και για το λόγο αυτό πρέπει να αποφευχθεί. Εφόσον στις πολλαπλές εμφανίσεις του ίδιου παραδείγματος έχει αντιστοιχηθεί η ίδια κλάση, οι εμφανίσεις αυτές δεν αποτελούν πρόβλημα αλλά αντίθετα υποδεικνύουν πως το παράδειγμα αυτό είναι χαρακτηριστικό της κλάσης του και επομένως, έστω και έμμεσα, αποτελούν επιπλέον πληροφορία. Η δεύτερη περίπτωση όμως θα αποτελούσε πρόβλημα ακόμα και για την ανθρώπινη νοημοσύνη και φυσικά δημιουργεί και στον αλγόριθμο ο οποίος δεν μπορεί, σε αυτή την περίπτωση, να βρει ποια είναι αυτή η ιδιότητα που διαφοροποιεί τα δύο παραδείγματα.

Οι καταστάσεις αυτές κάθε άλλο παρά ασυνήθιστες είναι. Στην Ιατρική, για παράδειγμα, οι ασθενείς παρακολουθούνται από διαφορετικούς γιατρούς οι οποίοι, λόγω διαβαθμίσεων στις ικανότητες αλλά και στην εμπειρία τους, παίρνουν διαφορετικές διαγνωστικές ή θεραπευτικές αποφάσεις για ίδια περιστατικά. Οι διαβαθμίσεις στην ικανότητα των γιατρών αποτελούν έναν ακόμα εξωτερικό παράγοντα που επηρεάζει την εφαρμογή.

Γίνεται πια φανερή η ανάγκη ενός συστήματος το οποίο θα ελέγχει τα παραδείγματα και θα αποτρέπει την εμφάνιση του ίδιου παραδείγματος σαν μέλος διαφορετικής κλάσης. Η διαδικασία αυτή ονομάζεται **Προεπεξεργασία** των παραδειγμάτων (**Preprocessing of the Example Set**).

Στα πλαίσια της εργασίας αυτής αναπτύχθηκε ένα τέτοιο σύστημα προεπεξεργασίας και ενσωματώθηκε στον αλγόριθμο IDDD. Το σύστημα προεπεξεργασίας ενεργοποιείται από τον χρήστη μέσα από το περιβάλλον του IDDD. Ελέγχει τα παραδείγματα που έχουν ήδη δοθεί σαν είσοδος στο σύστημα. Μετά το τέλος της διαδικασίας προεπεξεργασίας το σύστημα μπορεί να εξάγει το δέντρο απόφασης από τα διορθωμένα πια παραδείγματα.

Όταν το σύστημα εντοπίσει δύο παραδείγματα που είναι ίδια όσον αφορά τις τιμές ιδιοτήτων αλλά ανήκουν σε άλλη κλάση δεν μπορεί να αποφασίσει ποια από τις δύο είναι η σωστή και ξεκινά ένα διάλογο με τον χρήστη όπου, αφού του παρουσιάσει το πρόβλημα, ζητά από αυτόν να καθορίσει το παράδειγμα που εμπεριέχει την λανθασμένη πληροφορία και το διορθώνει.

Για παράδειγμα θεωρούμε την περίπτωση όπου μία εφαρμογή περιγράφεται από δύο ιδιότητες A και B με τιμές a_i και b_i αντίστοιχα και έχουμε τις κλάσεις C_1, \dots, C_4 . Στο σύνολο παραδειγμάτων εκπαίδευσης μπορεί να εμφανίζονται ζεύγη παραδειγμάτων της μορφής:

$$a_1, b_3, C_1$$

$$a_1, b_3, C_4$$

ο χρήστης μπορεί είτε να διορθώσει την κλάση στο ένα παράδειγμα (C_1, C_4), είτε ακόμα να αποφασίσει πως η τιμή μιας ιδιότητας έχει εκχωρηθεί λανθασμένα.

Αρχικά προσδιορίζει ποιο από τα δύο παραδείγματα εμπεριέχει την λανθασμένη πληροφορία και:

1. Αν θέλει να αλλάξει την κλάση, το σύστημα του παρουσιάζει ένα κατάλογο με τις δυνατές κλάσεις του πεδίου απ' όπου θα διαλέξει τη σωστή.
2. Αν θέλει να αλλάξει την τιμή μιας ιδιότητας, πρέπει να προσδιορίσει ποιά είναι η ιδιότητα αυτή και από τον κατάλογο των δυνατών τιμών της, που θα του παρουσιάσει το σύστημα, διαλέγει την επιθυμητή. Με τον τρόπο αυτό δεν αυτοματοποιείται βέβαια η διαδικασία αφαίρεσης του θορύβου αλλά διευκολύνεται η εξέταση των δεδομένων από τον γιατρό καθώς παρουσιάζονται σε αυτόν αυτά που φαίνονται να είναι προβληματικά.

Όταν ο χρήστης δεν είναι σίγουρος για το τι ακριβώς θέλει να διορθωθεί ή δεν ξέρει πως να το διορθώσει μπορεί απλά να αναβάλει την απόφασή του για αργότερα αφήνοντας τα παραδείγματα όπως έχουν και διορθώνοντας κάποια άλλα που θα του εντοπιστούν παρακάτω.

5.3 Σχεδιασμός Πειραμάτων

Όπως έχουμε ήδη αναφέρει, ο NewId χρησιμοποιεί την αδιάφορη τιμή για τις ιδιότητες των οποίων η επιλογή δεν έχει νόημα σε ορισμένες περιπτώσεις. Για τη συγκεκριμένη εφαρμογή, όταν ένα παράδειγμα αναφέρεται στη πρώτη επίσκεψη του ασθενή, η εξέταση της τιμής των ιδιοτήτων `pre-one`, `pre-two` και `treatment` δεν έχει προφανώς νόημα. Έτσι, στα αντίστοιχα παραδείγματα χρησιμοποιούμε την αδιάφορη τιμή. Η τιμή αυτή έχει το χαρακτηριστικό πως μειώνει την πληροφοριακή ισχύ της ιδιότητας για την οποία χρησιμοποιείται. Περιμέναμε λοιπόν πως, κατά το σχηματισμό του δέντρου, δεν θα είχαμε κανόνες που θα περιείχαν τις ιδιότητες αυτές όταν η επίσκεψη ήταν η πρώτη. Ωστόσο ο NewId παράγει κανόνες όπως ο ακόλουθος:

```
AN (one=retractile)
KAI (pre-two=normal)
KAI (two=normal)
KAI (visit=first)
TOTE follow up.
```

Ο λόγος παραγωγής του παραπάνω κανόνα είναι ότι επιλέγεται η ιδιότητα `pre-two` πριν από την `visit`. Τα παραδείγματα που αναφέρονται στην πρώτη επίσκεψη, έχουν αδιάφορη τιμή για την ιδιότητα `pre-two`. Μετά την επιλογή της `pre-two` τα παραδείγματα αυτά γίνονται μέλη όλων των υποσυνόλων που δημιουργούνται για κάθε κλαδί. Όταν σε επόμενα βήματα ο NewId καλείται να διαχωρίσει μεταξύ των παραδειγμάτων με βάση την κλάση, βρίσκει πως μεγαλύτερη πληροφορία δίνει η ιδιότητα για τον αριθμό επίσκεψης. Έτσι δημιουργείται ένα κλαδί (κανόνας) που αντιστοιχεί στη πρώτη επίσκεψη.

Οι ιδιότητες που αναφέρονται στην προηγούμενη επίσκεψη (`pre-one`, `pre-two`) είναι χρήσιμες μόνο σε συνδυασμό με την αντίστοιχη πληροφορία για την τρέχουσα επίσκεψη (ιδιότητες `one`, `two`). Προστέθηκαν ώστε να δίνεται η δυνατότητα σύγκρισης της παρούσας και της προηγούμενης κατάστασης. Φυσικά, τίποτα δεν αναγκάζει τον αλγόριθμο να χρησιμοποιήσει τις ιδιότητες αυτές. Όταν όμως δεν μπορεί να διαχωρίσει τα παραδείγματα με βάση μόνο τα υπόλοιπα στοιχεία, ελπίζουμε πως θα λάβει υπόψην του και

αυτές. Αυτή είναι εξάλλου και η πρακτική που θα ακολουθούσε ένας γιατρός. Ωστόσο εξετάζοντας το σύνολο κανόνων που παράγει ο NewId παρατηρεί κανείς κανόνες όπως ο ακόλουθος:

AN (one=inguinal)
KAI (visit=third)
KAI (pre-two=normal)
TOTE surgical

Ο κανόνας αυτός δεν μπορεί να χρησιμοποιηθεί γιατί δεν περιέχει την εξέταση της τιμής της ιδιότητας που αφορά στην παρούσα κατάσταση (two). Κατά τη δημιουργία του κανόνα αυτού εμφανίστηκε η pre-two σαν η πιο πληροφοριακή ιδιότητα και φυσικά ήταν αυτή που επιλέχτηκε από τον αλγόριθμο καθώς δεν υπήρχε περιορισμός για την επιλογή της. Αν ο στόχος μας είναι η παραγωγή κανόνων σαν αυτούς που χρησιμοποιούν οι ειδικοί της εφαρμογής, θα επιθυμούσαμε την εξάλειψη κανόνων όπως ο παραπάνω.

Τέλος, ένα από τα χαρακτηριστικά της εφαρμογής είναι πως στην περίπτωση που οι ιδιότητες one και two παίρνουν την τιμή not-palpable το αντίστοιχο μέγεθος δεν μπορεί να υπολογισθεί και επομένως δεν χρησιμοποιείται στη διαδικασία λήψης της απόφασης. Για τις περιπτώσεις αυτές χρησιμοποιούμε στα παραδείγματα την ειδική άγνωστη τιμή. Ωστόσο ο NewId παράγει κανόνες όπως:

AN (one=not palpable)
KAI (size of two>0.60)
KAI (two=not palpable)
TOTE hormonal

Αν κατά τη δημιουργία του κανόνα είχε επιλεγεί πρώτα η ιδιότητα two στο αντίστοιχο υποσύνολο παραδειγμάτων θα είχαμε παραδείγματα με άγνωστη τιμή για το μέγεθος και επομένως δεν θα είχε ποτέ μεγάλη πληροφοριακή ισχύ η ιδιότητα size of two.

Η διόρθωση/εξάλειψη κανόνων όπως αυτών που παρουσιάστηκαν παραπάνω θα οδηγούσε τον αλγόριθμο στη παραγωγή ενός καλύτερου συνόλου κανόνων.

Η χρήση του IDDD ο οποίος διαθέτει μηχανισμούς χειρισμού σχέσεων εξάρτησης θα μπορούσε να οδηγήσει σε βελτίωση του συνόλου των παραγόμενων κανόνων. Για την αξιολόγηση της απόδοσης του IDDD σε σχέση με αυτή του NewId εκτελέστηκαν τέσσερα διαφορετικά πειράματα.

Για το πρώτο πείραμα, χρησιμοποιήθηκε ο NewId στην παραγωγή κανόνων χωρίς την οποιαδήποτε πληροφορία για το πεδίο εφαρμογής. Το σύνολο των κανόνων που προέκυψε ονομάζεται ‘**αρχικό**’ σύνολο κανόνων. Για παράδειγμα, οι τρεις κανόνες που αναφέρονται παραπάνω ανήκουν σε αυτό το σύνολο.

Στους προβληματικούς κανόνες που παρουσιάσαμε μπορεί κανείς να παρατηρήσει πως η αιτία του προβλήματος μπορεί να είναι η σειρά επιλογής των ιδιοτήτων. Έτσι επιχειρήσαμε να ορίσουμε μία σχέση διάταξης μεταξύ των ιδιοτήτων χρησιμοποιώντας τον μηχανισμό διάταξης που προσφέρει ο NewId [17]. Σύμφωνα με την αναπαράσταση που απαιτεί ο NewId, ορίζουμε:

1. visit BEFORE pre-one

visit BEFORE pre-two

(Που σημαίνει ότι οι ιδιότητες, που αφορούν στην προηγούμενη επίσκεψη, δεν μπορούν να χρησιμοποιηθούν αν προηγουμένως δεν έχει χρησιμοποιηθεί η ιδιότητα για τον αριθμό της επίσκεψης. Ο σκοπός είναι, προφανώς, να μην εξετάζονται οι ιδιότητες αυτές για την πρώτη επίσκεψη.)

2. one BEFORE pre-one

two BEFORE pre-two

(Οι ιδιότητες, που αφορούν στην προηγούμενη επίσκεψη, δεν μπορούν να χρησιμοποιηθούν παρά μόνο όταν έχει ήδη χρησιμοποιηθεί η αντίστοιχη ιδιότητα για την παρούσα κατάσταση.)

3. one BEFORE size of one

two BEFORE size of two

(Αν επιλέγεται πρώτα η ιδιότητα της κατάστασης του όρχη δεν θα εμφανίζονται περιπτώσεις όπου αυτή είναι not-palpable και εξετάζεται το μέγεθος.)

Στο δεύτερο πείραμα χρησιμοποιήθηκαν οι παραπάνω σχέσεις διάταξης και οι παραγόμενοι κανόνες ονομάστηκαν **‘διατεταγμένο’** σύνολο κανόνων. Οι κανόνες του συνόλου αυτού έχουν βελτιωθεί σε σχέση με τους κανόνες του **‘αρχικού’** συνόλου αλλά υπάρχουν ακόμη κανόνες της μορφής:

AN (one=not palpable)

KAI (size of one \leq 0.60)

KAI (visit=second)

TOTE surgical

Η ύπαρξη τέτοιων κανόνων δείχνει ότι, παρά τον ορισμό διάταξης μεταξύ των ιδιοτήτων one και size of one, δεν αντιμετωπίζεται το πρόβλημα που δημιουργείται λόγω της εμφάνισης του μεγέθους σε περιπτώσεις που αυτό είναι αδύνατο να υπολογιστεί.

Για το τρίτο πείραμα, χρησιμοποιήθηκε ο αλγόριθμος IDDD. Οι σχέσεις διάταξης, μεταξύ των ιδιοτήτων one, two και των αντίστοιχων μεγεθών, αντικαταστάθηκαν με σχέσεις εξάρτησης. Σημειώνουμε πως η χρήση της αποκλειστικής εξάρτησης μας επιτρέπει να ορίσουμε επιπλέον την συγκεκριμένη τιμή/τιμές της μάνας που απαγορεύουν την επιλογή της κόρης. Η αντικατάσταση των σχέσεων διάταξης επιτυγχάνεται με τον ορισμό των παρακάτω σχέσεων αποκλειστικής εξάρτησης (παράγρ. 4.2):

1. size of one DEPENDS_ON one:not palpable
2. size of two DEPENDS_ON two:not palpable

Το σύνολο κανόνων που δημιουργήθηκε από αυτό το πείραμα ονομάζεται σύνολο **‘διάταξης-εξάρτησης’** γιατί, για τη δημιουργία του, δεν αντικαταστάθηκαν όλες οι σχέσεις διάταξης από τις αντίστοιχες εξαρτήσεις.

Τέλος, εκτελέστηκε ένα τέταρτο πείραμα για το οποίο οι σχέσεις διάταξης αντικαταστάθηκαν, όπου αυτό ήταν δυνατό, με σχέσεις εξάρτησης. Κάθε ιδιότητα που αναφέρεται στην τρέχουσα κατάσταση (one, two) είναι μάνα και της ιδιότητας για την προηγούμενη

επίσκεψη (pre-one, pre-two) αλλά και του αντίστοιχου μεγέθους. Στον IDDD, μία κόρη-ιδιότητα μπορεί να έχει μόνο μία μάνα. Έτσι οι σχέσεις “visit BEFORE pre-one” και “visit BEFORE pre-two” δεν μπορούν να αντικατασταθούν από εξαρτήσεις. Ο λόγος είναι ότι τότε θα έπρεπε να οριστούν αντίστοιχα οι σχέσεις “pre-one DEPENDS_ON visit” και “pre-two DEPENDS_ON visit”. Τότε όμως η ιδιότητα pre-one θα είχε σαν μάνες τις ιδιότητες one και visit. Για το λόγο αυτό οι παραπάνω σχέσεις διάταξης διατηρούνται και γι’αυτό το πείραμα. Έτσι γι’αυτό το τέταρτο πείραμα, του οποίου το παραγόμενο σύνολο κανόνων ονομάζουμε σύνολο ‘εξάρτησης’, χρησιμοποιήθηκαν οι ακόλουθες σχέσεις:

1. size of one DEPENDS_ON one : not palpable
2. size of two DEPENDS_ON two : not palpable
3. pre-one DEPENDS_ON one
4. pre-two DEPENDS_ON two
5. visit BEFORE pre-one
6. visit BEFORE pre-two

5.4 Αξιολόγηση Αποτελεσμάτων

5.4.1 Ποσοτική Αξιολόγηση

Σε κάθε ένα από τα τέσσερα πειράματα που περιγράφονται στην παράγραφο 5.3, αντιστοιχεί μία διαφορετική περιγραφή του πεδίου εφαρμογής. Τα μέτρα αξιολόγησης του συνόλου κανόνων που παράγεται από έναν αλγόριθμο μάθησης με παραδείγματα είναι ο αριθμός των κανόνων και η ακρίβεια κατάταξης των παραδειγμάτων στις κλάσεις. Το αποτέλεσμα κάθε περιγραφής της εφαρμογής αξιολογήθηκε με βάση αυτά τα μέτρα.

Όσο μεγαλύτερος είναι ο αριθμός των κανόνων που παράγονται τόσο πιο πιθανό είναι ότι ο αλγόριθμος, προκειμένου να καλύψει το σύνολο των παραδειγμάτων, δημιούργησε πολλούς εξειδικευμένους κανόνες. Ο στόχος αντίθετα είναι η παραγωγή γενικών κανόνων ώστε να καλύπτουν και άλλα παραδείγματα της εφαρμογής που δεν ανήκουν στο εκπαιδευτικό σύνολο.

Οι αλγόριθμοι μάθησης με παραδείγματα προσπαθούν να βρουν κανόνες ώστε να καλύψουν όλα τα εκπαιδευτικά παραδείγματα. Όταν τα παραδείγματα αυτά δεν περιέχουν άγνωστες τιμές, η ακρίβεια κατάταξης είναι 100%. Δεν έχει νόημα λοιπόν να θεωρούμε μέτρο αξιολόγησης την ακρίβεια κατάταξης για αυτά τα παραδείγματα. Αντίθετα, είναι σημαντικό μέτρο αξιολόγησης η ακρίβεια κατάταξης παραδειγμάτων που δεν ανήκουν στο σύνολο των εκπαιδευτικών. Τα παραδείγματα αυτά αποτελούν το σύνολο των παραδειγμάτων **ελέγχου**.

Για την συγκεκριμένη ιατρική εφαρμογή είχαμε στη διάθεσή μας 265 παραδείγματα. Από τα παραδείγματα αυτά 172 ανήκουν στην κλάση ‘παρακολούθηση’, 31 στην κλάση ‘εγχείρηση’ και 62 στη κλάση ‘ορμονοθεραπεία’. Τα νούμερα αυτά αντιπροσωπεύουν την πραγματική συχνότητα εμφάνισης των αντίστοιχων θεραπευτικών αποφάσεων.

Για τον υπολογισμό των μέτρων αξιολόγησης, το σύνολο των παραδειγμάτων διαιρείται σε δύο υποσύνολα. Το πρώτο υποσύνολο περιέχει 232 παραδείγματα και αποτελεί

το εκπαιδευτικό σύνολο ενώ το δεύτερο υποσύνολο περιέχει τα υπόλοιπα 33 παραδείγματα και χρησιμοποιείται για τον υπολογισμό της ακρίβειας κατάταξης του δέντρου. Για λόγους αντικειμενικότητας ο διαμερισμός του συνόλου παραδειγμάτων έγινε 8 φορές και υπολογίστηκαν οι μέσοι όροι του αριθμού των κανόνων που παράγονται κάθε φορά και της ακρίβειας κατάταξης. Η διαδικασία αυτή αξιολόγησης είναι γνωστή στην βιβλιογραφία [31]. Ο πίνακας 2 περιέχει τα αποτελέσματα αυτών των μετρήσεων.

Κριτήρια Αξιολόγησης	Περιγραφή Εφαρμογής			
	Αρχική	Διατεταγμένη	Διάταξη-Εξάρτηση	Εξάρτηση
Μέσος όρος κανόνων	146	84	84.5	84.5
Μέση Ακρίβεια Κατάταξης	76.9	67.4	66.5	68.4

Πίνακας 2: Μέσοι Όροι Αριθμού Κανόνων και Ακρίβειας Κατάταξης.

Τα αποτελέσματα που φαίνονται στον πίνακα 2 δείχνουν καθαρά πως η εισαγωγή γνώσης για την εφαρμογή μειώνει την ακρίβεια κατάταξης νέων παραδειγμάτων αλλά επιφέρει και σημαντικότερη μείωση στον αριθμό των κανόνων που παράγονται. Η μείωση της ακρίβειας κατάταξης είναι αναμενόμενη. Ο αλγόριθμος προσπαθεί να δημιουργήσει κανόνες για να καλύψει ακόμα και τα παραδείγματα που περιέχουν ‘θόρυβο’. Με την εισαγωγή γνώσης επιχειρείται η εξαγωγή κανόνων που δεν θα επηρεάζονται από αυτό το θόρυβο. Είναι λοιπόν επόμενο τα παραδείγματα που περιέχουν ‘θόρυβο’ να μην κατατάσσονται αρκετά καλά. Αναλυτικότερα, στη συγκεκριμένη εφαρμογή, η απόφαση του γιατρού επηρεάζεται από παράγοντες που δεν εμπίπτουν στην Ιατρική. Είναι δυνατόν μία απόφαση για εγχείρηση να μεταβληθεί λόγω μη διαθεσιμότητας του χειρουργείου, ή αντίθετα να γίνει εγχείρηση και όχι Ορμονοθεραπεία λόγω κόστους. Με τις σχέσεις εξάρτησης προσπαθούμε να εξάγουμε τους Ιατρικούς κανόνες οι οποίοι προφανώς δεν συμπεριλαμβάνουν περιπτώσεις όπως οι προηγούμενες.

Ωστόσο παρατηρούμε πως η αντικατάσταση των σχέσεων διάταξης από σχέσεις εξάρτησης έχει σαν αποτέλεσμα την αύξηση της ακρίβειας κατάταξης. Η σύγκριση μεταξύ των συνόλων Διάταξης, Διάταξης-Εξάρτησης και Εξάρτησης ενισχύει αυτή την παρατήρηση. Όσο μειώνεται ο αριθμός των σχέσεων διάταξης αυξάνεται η ακρίβεια κατάταξης.

Το γεγονός αυτό οφείλεται στον τρόπο υλοποίησης της διάταξης στον NewId. Ιδιότητες με μεγάλη πληροφοριακή ισχύ δεν επιλέγονται αν έπονται κάποιας άλλης. Επιπλέον δεν δημιουργείται κόμβος για την ιδιότητα που προηγείται με αποτέλεσμα να μην μπορεί ποτέ να επιλεγεί η ιδιότητα που έπεται. Αυτή η μείωση της ακρίβειας κατάταξης κληρονομείται στα σύνολα Διάταξης-Εξάρτησης και Εξάρτησης λόγω της ύπαρξης σχέσεων διάταξης στις περιγραφές της εφαρμογής σύμφωνα με τις οποίες δημιουργήθηκαν αυτά τα σύνολα.

5.4.2 Ποιοτική Αξιολόγηση

Ο στόχος κατά την ανάπτυξη του IDDD και του ορισμού των σχέσεων εξάρτησης ήταν η παραγωγή καλύτερων ποιοτικά κανόνων. Οι κανόνες θεωρούνται ποιοτικά καλοί αν προσεγγίζουν το μοντέλλο που χρησιμοποιεί ο ειδικός της εφαρμογής. Στην συγκεκριμένη εφαρμογή, όπως και σε πολλές άλλες, τα παραδείγματα περιέχουν ‘θόρυβο’. Στοχεύουμε στην παραγωγή καλών κανόνων έστω και με αυτά τα παραδείγματα. Είναι ενδεχόμενο, οι κανόνες αυτοί να κατατάσσουν τα παραδείγματα σε κλάσεις με μικρότερη ακρίβεια. Επειδή οι παράγοντες που προκαλούν ‘θόρυβο’ είναι ανεξάρτητοι της ιατρικής γνώσης δεν θέλουμε αυτοί να επηρεάζουν τους κανόνες.

Την αξιολόγηση των κανόνων έκαναν γιατροί της Πανεπιστημιακής Κλινικής Ηρακλείου. Για κάθε σύνολο κανόνων που παράγεται από τα τέσσερα πειράματα (παράγραφος 5.3), τους ζητήθηκε να αντιστοιχίσουν σε κάθε κανόνα έναν από τους ακόλουθους ποιοτικούς χαρακτηρισμούς.

1. **Χρήσιμος** είναι ένας κανόνας που χρησιμοποιείται και από τον γιατρό προκειμένου να καταλήξει στην θεραπευτική απόφαση.
2. **Πολύ γενικός** είναι ένας κανόνας που θα μπορούσε να θεωρηθεί χρήσιμος εάν λάμβανε υπόψη κάποιο επιπλέον στοιχείο (πληροφορία).
3. **Πολύ ειδικός** είναι ένας κανόνας που λαμβάνει υπόψη επιπλέον πληροφορία.

4. **Αχρηστος** είναι ένας κανόνας όταν οι προϋποθέσεις που θεωρεί οδηγούν κανονικά σε διαφορετική απόφαση και όχι σε αυτή στην οποία καταλήγει.

Σύμφωνα με την κλίμακα αυτή, το σύνολο κανόνων είναι καλό όταν περιέχει μεγάλο ποσοστό χρησιμων κανόνων. Θεωρείται προτιμότερη η ύπαρξη στο σύνολο κανόνων που είναι πολύ γενικοί παρά πολύ ειδικοί. Και τα δύο είδη κανόνων προσεγγίζουν χρησιμους κανόνες. Η διαφορά είναι πως ένας πολύ ειδικός κανόνας δημιουργήθηκε από την προσπάθεια του αλγορίθμου να καλύψει παραδείγματα πολύ ειδικά που αποτελούν ενδεχομένως εξαίρεση λόγω ύπαρξης θορύβου. Σε ένα πολύ γενικό κανόνα η έλλειψη κάποιας πληροφορίας μπορεί να οφείλεται στο γεγονός πως η αντίστοιχη ιδιότητα έχει την ίδια τιμή μεταξύ παραδειγμάτων διαφορετικής κλάσης. Έτσι, αν και η ιδιότητα αυτή είναι **χαρακτηριστική** της κλάσης, δεν είναι διαχωριστική μεταξύ των κλάσεων και επομένως δεν έχει πληροφοριακή ισχύ για τον αλγόριθμο.

Οι πίνακες 3, 4, 5 και 6 περιέχουν τα στοιχεία της συγκριτικής αξιολόγησης των αποτελεσμάτων για τα τέσσερα πειράματα.

Σε κάθε πίνακα συγκρίνουμε ανά δύο τα αποτελέσματα των πειραμάτων. Η διαδικασία της σύγκρισης είναι η εξής: Αν A και B τα σύνολα κανόνων, εντοπίζουμε τους κανόνες που είναι κοινοί και για τα δύο σύνολα. Αφαιρούμε τους κανόνες αυτούς από τα A και B και συνεπώς σχηματίζονται τα σύνολα $A-B$ και $B-A$. Περιοριζόμαστε στην αξιολόγηση των μη κοινών κανόνων γιατί αυτοί είναι που προκύπτουν λόγω των διαφορών στην περιγραφή της εφαρμογής. Έτσι σε κάθε ένα από τους πίνακες 3, 4, 5 και 6 υπάρχει μία γραμμή που καταγράφει, για κάθε σύνολο τον αριθμό των διαφορετικών κανόνων που περιέχονται σε αυτό. Τα ποσοστά για την κλίμακα αξιολόγησης υπολογίζονται με βάση αυτούς τους αριθμούς.

Μαζί με την ποιοτική αξιολόγηση των κανόνων, στους πίνακες αυτούς περιέχεται και ένα μέτρο ποσοτικής αξιολόγησής τους. Επειδή όμως είναι περισσότερο διαισθητικό μέτρο περιλαμβάνεται στην ποιοτική αξιολόγηση. Σε μερικές περιπτώσεις ένας κανόνας είναι δυνατόν να καλύπτει, εκτός από ορισμένα παραδείγματα της κλάσης στην οποία καταλήγει, και μερικά άλλα διαφορετικής κλάσης. Τα παραδείγματα αυτά είτε έχουν

μικρό βάρος είτε είχαν εξαντληθεί οι ιδιότητες και δεν ήταν δυνατός ο παραπέρα διαχωρισμός μεταξύ κλάσεων. Αυτά τα παραδείγματα διαφορετικής κλάσης ονομάζουμε αρνητικά παραδείγματα του κανόνα. Και στις δύο περιπτώσεις ο αλγόριθμος τερματίζει. Ωστόσο η ύπαρξη τέτοιων κανόνων στο δέντρο αποτελεί, ενδεχομένως, αιτία μείωσης της ακρίβειας κατάταξης όταν το δέντρο προσπαθεί να κατατάξει νέα παραδείγματα. Θα θέλαμε λοιπόν να ελαχιστοποιήσουμε, κατά το δυνατόν, τον αριθμό τους.

Στον πίνακα 3 παρατίθενται τα στοιχεία για τα δέντρα (σύνολα κανόνων) Αρχικό και Διατεταγμένο. Η σύγκριση των συνόλων αυτών αποσκοπεί στον έλεγχο της ποιοτικής βελτίωσης που επιτυγχάνεται με τη χρησιμοποίηση σχέσεων διάταξης.

	Σύνολα Κανόνων	
	Αρχικό	Διατεταγμένο
Κριτήρια Αξιολόγησης		
Αχρηστοι	117 (66%)	31 (49%)
Πολύ γενικοί	34 (19%)	15 (24%)
Πολύ ειδικοί	16 (9%)	14 (22%)
Χρήσιμοι	9 (5%)	3 (5%)
Αριθμός Διαφορετικών Κανόνων	176	63
Κανόνες που καλύπτουν Αρνητικά Παραδείγματα	28 (16%)	8 (13%)

Πίνακας 3: Σύγκριση μεταξύ Αρχικού και Διατεταγμένου Συνόλου Κανόνων.

Όπως φαίνεται από τον παραπάνω πίνακα, η χρησιμοποίηση σχέσεων διάταξης δεν βελτιώνει κατά πολύ το τελικό αποτέλεσμα. Το ποσοστό των χρήσιμων κανόνων παραμένει σταθερό (5%). Ωστόσο έχει μειωθεί σημαντικά ο αριθμός των διαφορετικών κανόνων αλλά και ο αριθμός των κανόνων που καλύπτουν αρνητικά παραδείγματα.

Στο πίνακα 4 φαίνεται η σύγκριση μεταξύ των δέντρων Αρχικό και Διάταξης-Εξάρτησης.

Όπως φαίνεται από τον πίνακα, ακόμα και η μερική αντικατάσταση των σχέσεων διάταξης από εξαρτήσεις, επιφέρει σημαντική βελτίωση στο τελικό αποτέλεσμα. Παρατηρείται σημαντική μείωση των ποσοστών αλλά και των απόλυτων αριθμών για όλες τις κατηγορίες κανόνων εκτός από αυτή των χρήσιμων. Ο αριθμός των χρήσιμων κανό-

Κριτήρια Αξιολόγησης	Σύνολα Κανόνων	
	Αρχικό	Διάταξης-Εξάρτησης
Αχρηστοι	119 (64%)	7 (10%)
Πολύ γενικοί	41 (22%)	23 (32%)
Πολύ ειδικοί	17 (9%)	9 (12%)
Χρήσιμοι	9 (5%)	34 (47%)
Αριθμός Διαφορετικών Κανόνων	186	73
Κανόνες που καλύπτουν Αρνητικά Παραδείγματα	30 (16%)	7 (10%)

Πίνακας 4: Σύγκριση μεταξύ Αρχικού και Διάταξης-Εξάρτησης Συνόλου Κανόνων.

νων σημείωσε μεγάλη αύξηση. Μειώθηκε επίσης δραστικά ο αριθμός των κανόνων που καλύπτουν αρνητικά παραδείγματα.

Στον πίνακα 5 φαίνεται η σύγκριση μεταξύ των δέντρων Αρχικό και Εξάρτησης.

Κριτήρια Αξιολόγησης	Σύνολα Κανόνων	
	Αρχικό	Εξάρτησης
Αχρηστοι	119 (63%)	6 (8%)
Πολύ γενικοί	42 (22%)	23 (31%)
Πολύ ειδικοί	18 (10%)	6 (8%)
Χρήσιμοι	9 (5%)	38 (52%)
Αριθμός Διαφορετικών Κανόνων	188	73
Κανόνες που καλύπτουν Αρνητικά Παραδείγματα	28 (15%)	7 (10%)

Πίνακας 5: Σύγκριση μεταξύ Αρχικού και Εξάρτησης Συνόλου Κανόνων.

Η χρήση όλων των δυνατών εξαρτήσεων επιφέρει ακόμα μεγαλύτερη βελτίωση στην ποιότητα των παραγόμενων κανόνων (πίνακας 5). Η παρατήρηση αυτή γίνεται περισσότερο προφανής αν συγκρίνει κανείς τους πίνακες 4 και 5.

Ο ορισμός των σχέσεων διάταξης στον NewId στοχεύει στον εντοπισμό και χειρισμό 'κρυμμένων' σχέσεων μεταξύ ιδιοτήτων. Είναι δυνατόν να εκφράσει κανείς σχέσεις απλής εξάρτησης χρησιμοποιώντας μία σχέση διάταξης μεταξύ των εξαρτημένων ιδιοτήτων. Αν

και σύμφωνα με τα πειραματικά αποτελέσματα η αντικατάσταση εξαρτήσεων από διατάξεις δεν αποδίδει αρκετά καλά, υπάρχει το επιπλέον μειονέκτημα ότι είναι αδύνατος ο ορισμός της αποκλειστικής εξάρτησης με τη χρησιμοποίηση της διάταξης. Όταν οι σχέσεις διάταξης αντικαταστάθηκαν από σχέσεις εξάρτησης παρατηρήθηκε βελτίωση στην ποιότητα των κανόνων αλλά και αύξηση της ακρίβειας κατάταξης του δέντρου. Για την ενίσχυση της παραπάνω παρατήρησης με βάση τα πειραματικά αποτελέσματα παραθέτουμε τον πίνακα 6 ο οποίος περιέχει την σύγκριση μεταξύ του Διατεταγμένου και του Διάταξης-Εξάρτησης δέντρου. Το δέντρο Διάταξης-Εξάρτησης περιέχει πολύ περισσότερους χρήσιμους κανόνες ενώ είναι μειωμένοι οι αριθμοί όλων των άλλων κατηγοριών κανόνων.

Κριτήρια Αξιολόγησης	Σύνολα Κανόνων	
	Διατεταγμένο	Διάταξης-Εξάρτησης
Αχρηστοι	19 (58%)	3 (9%)
Πολύ γενικοί	4 (12%)	5 (15%)
Πολύ ειδικοί	9 (27%)	8 (24%)
Χρήσιμοι	1 (3%)	17 (52%)
Αριθμός Διαφορετικών Κανόνων	33	33
Κανόνες που καλύπτουν Αρνητικά Παραδείγματα	7 (21%)	3 (9%)

Πίνακας 6: Σύγκριση μεταξύ Διατεταγμένου και Διάταξης-Εξάρτησης Συνόλου Κανόνων.

6 Συμπεράσματα και Δυνατότητες Μελλοντικής Έρευνας

6.1 Συμπεράσματα

Οι αλγόριθμοι επαγωγικής μάθησης με παραδείγματα στοχεύουν στην εξαγωγή ενός συνόλου κανόνων της εκάστοτε εφαρμογής ελαχιστοποιώντας την ανάγκη εισαγωγής γνώσης από τον ειδικό. Ωστόσο είναι χρήσιμη η ανάπτυξη συστημάτων που, διατηρώντας την απλότητά τους, θα μπορούν να διαχειρίζονται επιπλέον γνώση. Η γνώση αυτή μπορεί να χρησιμοποιηθεί στην διαδικασία μάθησης για να οδηγήσει τον αλγόριθμο στην παραγωγή κανόνων που ακολουθούν την ίδια διαδικασία λήψης μιας απόφασης με αυτή που εφαρμόζουν οι ειδικοί.

Στα πλαίσια της εργασίας αυτής μελετήθηκε μία ειδική μορφή που μπορεί να έχει η γνώση του ειδικού, οι σχέσεις εξάρτησης μεταξύ ιδιοτήτων. Μελετήθηκε ο τυπικός ορισμός των σχέσεων αυτών και παρουσιάστηκε ο αλγόριθμος IDDD, ο οποίος βασίζεται στον γνωστό από τη βιβλιογραφία αλγόριθμο ID3. Για την ανάπτυξη του IDDD χρησιμοποιήθηκε το σύστημα NewId, που υλοποιήθηκε από το Turing Institute. Ο IDDD αποτελεί μία βελτιωμένη έκδοση του NewId γιατί υλοποιεί και χειρίζεται σχέσεις εξάρτησης μεταξύ ιδιοτήτων. Ο αλγόριθμος αποδείχτηκε αποτελεσματικός στην εισαγωγή γνώσης και την εφαρμογή της στη διαδικασία δημιουργίας του δέντρου απόφασης. Η αποτελεσματικότητα του αποδείχτηκε μέσω προσεκτικά σχεδιασμένων πειραμάτων με την χρησιμοποίηση μιας ιατρικής εφαρμογής.

Η εισαγωγή των σχέσεων εξάρτησης πιθανόν να επιφέρει μείωση της ακρίβειας κατάταξης. Ο σκοπός του NewId είναι η εύρεση κανόνων για κατάταξη των παραδειγμάτων κατά 100%. Οι κανόνες αυτοί κρίνονται, πολλές φορές, άχρηστοι από τους γιατρούς. Η εισαγωγή των σχέσεων εξάρτησης αποσκοπεί στην παραγωγή χρήσιμων κανόνων. Για να εξαχθούν οι χρήσιμοι κανόνες τίθενται περιορισμοί στην χρήση της συνάρτησης πληροφο-

ρίας σαν μέτρο επιλογής της ιδιότητας-κόμβος. Η συνάρτηση πληροφορίας υπολογίζεται με βάση τα παραδείγματα και στοχεύει στη μεγιστοποίηση της ακρίβειας κατάταξης τους. Επομένως, οποιοσδήποτε περιορισμός της περιορίζει, συνήθως, την ακρίβεια κατάταξης. Η ακρίβεια κατάταξης εξαρτάται από το σύνολο των παραδειγμάτων που προσπαθούμε να κατατάξουμε. Αν τα παραδείγματα αυτά περιέχουν ‘θόρυβο’ παρόμοιο με αυτόν που περιέχουν τα εκπαιδευτικά παραδείγματα, θα έχουμε μείωση της ακρίβειας κατάταξης. Στην αντίθετη περίπτωση είναι πιθανή ακόμα και αύξηση της ακρίβειας κατάταξης.

Το μέγεθος των κανόνων, δηλαδή το πλήθος των προποθέσεων που θέτουν για τον καθορισμό μιας κλάσης, θεωρείται μέτρο αξιολόγησης της απόδοσης ενός αλγορίθμου. Ωστόσο, θεωρούμε πως κανόνες που απορρίπτονται από τους ειδικούς δεν μπορεί να αποτελούν μέλη του παραγόμενου συνόλου κανόνων άσχετα αν είναι μικρού μεγέθους. Σε μερικές περιπτώσεις η χρησιμοποίηση των σχέσεων εξάρτησης είναι δυνατόν να αυξήσει το μέγεθος των κανόνων. Για να γίνει κατανοητή η παραπάνω παρατήρηση παραθέτουμε το εξής παράδειγμα:

Θεωρούμε μία εφαρμογή που περιγράφεται από τις ιδιότητες A και B με $\{a_1, \dots, a_k\}$ και $\{b_1, \dots, b_v\}$ τα αντίστοιχα πεδία τιμών τους. Αν το σύνολο των εκπαιδευτικών παραδειγμάτων αποτελείται από τα:

$$a_1 b_1 C_1$$

$$a_2 b_1 C_2$$

τότε, προφανώς, η πιο πληροφοριακή ιδιότητα είναι η A, δημιουργείται ένας κόμβος γι' αυτήν και στο κλαδί με τιμή a_1 αντιστοιχίζεται η κλάση C_1 , ενώ στο κλαδί με τιμή a_2 ή κλάση C_2 . Αν, σύμφωνα με τον ειδικό, η A δεν είναι όντως σημαντική για την εφαρμογή αλλά αντίθετα σημαντική είναι η B, δεν υπάρχει κανένας τρόπος επιλογής της B με βάση τα παραπάνω παραδείγματα. Η ιδιότητα B έχει την ίδια τιμή και για τις δύο κλάσεις. Επομένως, δεν πρόκειται ποτέ να εμφανιστεί με μεγάλη πληροφοριακή ισχύ. Για να χρησιμοποιηθεί η B θα πρέπει να δοθούν στο σύστημα παραδείγματα στα οποία η ιδιότητα έχει τιμές τέτοιες ώστε να επιτρέπουν τον διαχωρισμό των κλάσεων. Στον IDDD κάτι τέτοιο δεν είναι απαραίτητο. Αν ορίσουμε πως η A εξαρτάται από την B ο IDDD θα δημιουργήσει ένα κόμβο για την B και στο κλαδί με τιμή b_1 θα δημιουργηθεί

ένα υποδέντρο ίδιο με την περίπτωση μη ύπαρξης της εξάρτησης. Με τον τρόπο αυτό μεγαλώνουν οι κανόνες κατά μία συνθήκη αλλά παράλληλα ολοκληρώνονται και αντανακλούν την πρακτική που ακολουθείται από τους ειδικούς. Στο σχήμα 5 φαίνονται τα αντίστοιχα δέντρα.

Σχήμα 5: (a) Το δέντρο του NewId, (b) Το δέντρο του IDDD

Επιπλέον, όπως φαίνεται από τα πειραματικά αποτελέσματα, ακόμα και οι σχέσεις διάταξης μεταξύ των ιδιοτήτων, είναι προτιμότερο να οριστούν σαν σχέσεις εξάρτησης. Έτσι αυξάνεται όχι μόνο η ποιότητα των κανόνων αλλά και η ακρίβεια κατάταξης του δέντρου απόφασης.

6.2 Δυνατότητες Μελλοντικής Έρευνας

Παρά την αποτελεσματικότητα του IDDD, θα μπορούσε κανείς να παρατηρήσει ότι θέτει περιορισμούς στον ορισμό των σχέσεων εξάρτησης. Κάθε Κόρη ιδιότητα μπορεί να έχει μόνο μία Μάνα. Ετσι δεν μπορούν να οριστούν πολύπλοκες εξαρτήσεις μίας ιδιότητας από άλλες. Συχνά σε πραγματικές εφαρμογές, οι σχέσεις εξάρτησης είναι πιο πολύπλοκες από τις εξαρτήσεις που ορίστηκαν στον IDDD. Μία ιδιότητα εξαρτάται από περισσότερες από μία άλλες ιδιότητες. Μία προφανής επέκταση του IDDD είναι ο ορισμός αλλά και η υλοποίηση τέτοιων πολύπλοκων σχέσεων.

Όταν μία ιδιότητα είναι η πιο πληροφοριακή αλλά είναι Κόρη περισσότερων από μίας ιδιοτήτων θα μπορούσε να δημιουργείται ένας κόμβος στο δέντρο για την πιο πληροφοριακή από τις ιδιότητες Μάνες. Η Κόρη ιδιότητα μπορεί να αποτελέσει κόμβο μόνο όταν έχουν ήδη επιλεγεί, για το μονοπάτι αυτό, όλες οι μάνες της.

Στις δυνατές κατευθύνσεις περαιτέρω έρευνας περιλαμβάνεται και η μελέτη εξαρτήσεων μεταξύ των ιδιοτήτων και της ειδικής ιδιότητας που αναπαριστά την κλάση.

Εκτός από την ειδική σχέση εξάρτησης, ιδιαίτερο ενδιαφέρον παρουσιάζει η μελέτη και ορισμός κι άλλων σχέσεων μεταξύ ιδιοτήτων. Τέτοιες σχέσεις έχουν περιγραφεί στην βιβλιογραφία [15], αλλά δεν είναι προφανής η χρήση τους στη διαδικασία σχηματισμού του δέντρου απόφασης. Αναφέρουμε ενδεικτικά τη σχέση του αμοιβαίου αποκλεισμού, σύμφωνα με την οποία η τιμή μιας ιδιότητας αποκλείει την ανάθεση μίας ή περισσότερων τιμών σε μία άλλη ιδιότητα, και αντιστρόφως. Η μελέτη των χαρακτηριστικών ενός ευρύτερου φάσματος πεδίων εφαρμογών θα μπορούσε να οδηγήσει στον εντοπισμό κι άλλων σχέσεων μεταξύ ιδιοτήτων.

Οι αλγόριθμοι, που περιγράφονται στα πλαίσια της εργασίας αυτής, θεωρούν πως σε κάθε ένα παράδειγμα αντιστοιχίζεται μία μόνο κλάση. Ωστόσο, η υπόθεση αυτή δεν ισχύει για ένα μεγάλο φάσμα εφαρμογών. Ένας ασθενής μπορεί να πάσχει από περισσότερες της μίας ασθένειες. Τα συμπτώματα της μιας επηρεάζουν τα συμπτώματα της άλλης και οδηγούν σε λανθασμένες επιλογές στοιχείων που είναι άσχετα με τη συγκεκριμένη κλάση. Μεγάλο ερευνητικό ενδιαφέρον παρουσιάζει η άρση του περιορισμού αυτού.

Η παραγωγή δέντρων απόφασης αποτελεί μία απλή αλλά αρκετά καλή τεχνική για την παραγωγή κανόνων. Αυτή η εργασία στοχεύει να συμβάλει στην τυπική και άμεση χρήση της γνώσης της εφαρμογής κατά τη διαδικασία δημιουργίας του δέντρου. Με αυτή την έννοια συμπληρώνει παρόμοιες προσπάθειες των Someren [15], Nunez [19] και άλλων. Η αποτελεσματικότητα των σχέσεων που υλοποιεί ο IDDD, δείχνει πόσο απαραίτητη είναι η έρευνα για τον εντοπισμό, τυπικό ορισμό και υλοποίηση της γνώσης του ειδικού στην διαδικασία επαγωγικής μάθησης.

Παραπομπές

- [1] B.G Buchanan. Can Machine Learning Offer Anything to Expert Systems? *Machine Learning*, 4:251–254, December 1989.
- [2] Louis Anthony Cox. Pragmatic Information-Seeking Strategies for Expert Classification Systems.
- [3] D.B Lenat. When Will Machines Learn? *Machine Learning*, 4:255–257, December 1989.
- [4] R.S Michalski, J.G Carbonell, and T.M Mitchell. *Machine Learning:An Artificial Intelligence Approach*, volume 1. Tioga Publishing, 1983.
- [5] R.S Michalski, J.G Carbonell, and T.M Mitchell. *Machine Learning:An Artificial Intelligence Approach*, volume 2. Morgan Kaufmann, 1986.
- [6] Yves Kodratoff and R.S Michalski. *Machine Learning:An Artificial Intelligence Approach*, volume 3. Morgan Kaufmann, 1990.
- [7] Yves Kodratoff. *Introduction to Machine Learning*. Pitman Publishing, 1988.
- [8] I. Mozetic. Hierarchical Model-Based Diagnosis. Technical Report MLI/89/1, George Mason University, 1989.
- [9] S. Amarel. Program Synthesis as a Theory Formation Task:Problem Representations and Solution Methods. In J.G Carbonell R.S Michalski and T.M Mitchell, editors, *Machine Learning:An Artificial Intelligence Approach*, volume 2, pages 499–570. Morgan Kaufmann, 1986.
- [10] Pat Langley. Areas of Application for Machine Learning. In *Proceedings of the Fifth International Symposium on Knowledge Engineering*, Sevilla, 1992.
- [11] J.R. Quinlan. Induction of Decision Trees. *Machine Learning*, 1:81–106, 1986.
- [12] J.R. Quinlan. Generating Production Rules from Decision Trees. In *IJCAI Proceedings*, pages 304–307, Milan, August 1987. Morgan Kaufmann.

- [13] P. Clark and T. Niblett. Induction in Noisy Domains. In *Progress in Machine Learning, EWSL Proceedings*, pages 11–30, Bled, Yugoslavia, May 1987. Sigma Press.
- [14] John Mingers. An Empirical Comparison of Pruning Methods for Decision Tree Induction. *Machine Learning*, 4:227–243, 1989.
- [15] M.W van Someren. Using Attribute Dependencies for Rule Learning. In K Morik, editor, *Knowledge Representation and Organization in Machine Learning*, pages 192–210. Springer-Verlag, Berlin, 1989.
- [16] M. Pazzani and D. Kibler. The Utility of Knowledge in Inductive Learning. *Machine Learning*, 9:57–94, 1992.
- [17] R.A. Boswell. Manual for Newid version 2.0. Technical report, Turing Institute, January 1990.
- [18] J.R Quinlan. Learning Logical Definitions from Relations. *Machine Learning*, 5:239–266, August 1990.
- [19] Marlon Nunez. The Use of Background Knowledge in Decision Tree Induction. *Machine Learning*, 6:231–250, May 1991.
- [20] Harish Ragavan and Larry Rendell. Relations, Knowledge and Empirical Learning. In *Proceedings of the Eighth International Workshop on Machine Learning*, pages 188–192. Morgan Kaufmann, 1991.
- [21] G. Silverstein and M.J Pazzani. Relational Cliches: Constraining Constructive Induction during Relational Learning. In *Proceedings of the Eighth International Workshop on Machine Learning*, pages 203–207. Morgan Kaufmann, 1991.
- [22] L. Gaga, V. Moustakis, G. Charissis, and S. Orphanoudakis. IDDD: An Inductive, Domain Dependent Decision Algorithm. In *ECML Proceedings*, Vienna, 1993.

- [23] R.M Goodman and P. Smyth. Decision Tree Design using Information Theory. *Knowledge Acquisition*, 2:1–1, 1990.
- [24] J.R. Quinlan. Decision Trees as Probabilistic Classifiers. In *Proceedings of the Fourth International Workshop on Machine Learning*, pages 31–37, University of California, June 1987.
- [25] S.L Crawford, R.M Fung, L.A Appelbaum, and R.M Tong. Classification Trees for Information Retrieval. In *Proceedings of the Eighth International Workshop on Machine Learning*, pages 245–249. Morgan Kaufmann, 1991.
- [26] I. Kononenko and I. Bratko. Information-Based Evaluation Criterion for Classifier’s Performance. *Machine Learning*, 6:67–80, January 1991.
- [27] J. Cheng, U.M Fayyad, K.B Irani, and Z. Qian. Improved Decision Trees: A Generalized Version of Id3. In *Proceedings of the Fifth International Conference on Machine Learning*, pages 100–106, University of Michigan, June 1988.
- [28] J.R Quinlan. Knowledge Acquisition from Structured Data. *IEEE EXPERT*, pages 32–37, December 1991.
- [29] Judea Pearl. Learning Hidden Causes from Empirical Data. In *IJCAI Proceedings*, pages 567–572, Los Angeles, August 1985.
- [30] W.E Grupe. Abnormalities of the Genital Tract. In ME Avery and HW Taeusch Jr., editors, *Schaffer’s Diseases of the Newborn*, pages 401–411. WB Saunders Company, Philadelphia, 1984.
- [31] L. Breiman, J.H Friedman, R.A Olshen, and C.J Stone. *Classification and Regression Trees*. Wadsworth International, 1984.