



ΕΘΝΙΚΟ ΜΕΤΣΟΒΙΟ ΠΟΛΥΤΕΧΝΕΙΟ
ΣΧΟΛΗ ΗΛΕΚΤΡΟΛΟΓΩΝ ΜΗΧΑΝΙΚΩΝ
ΚΑΙ ΜΗΧΑΝΙΚΩΝ ΥΠΟΛΟΓΙΣΤΩΝ
ΤΟΜΕΑΣ ΤΕΧΝΟΛΟΓΙΑΣ ΠΛΗΡΟΦΟΡΙΚΗΣ ΚΑΙ
ΥΠΟΛΟΓΙΣΤΩΝ

Machine Learning Techniques In Categorical Time Series Analysis Of Manufacturing Process

ΔΙΠΛΩΜΑΤΙΚΗ ΕΡΓΑΣΙΑ

ΧΑΡΑΛΑΜΠΟΣ ΜΙΧΑΗΛΙΔΗΣ, ΙΣΙΔΩΡΑ ΧΑΡΑ ΤΟΥΡΝΗ

Επιβλέπων : Νεκτάριος Κοζύρης
Καθηγητής Ε.Μ.Π.

Αθήνα, Ιούλιος 2016



ΕΘΝΙΚΟ ΜΕΤΣΟΒΙΟ ΠΟΛΥΤΕΧΝΕΙΟ
ΣΧΟΛΗ ΗΛΕΚΤΡΟΛΟΓΩΝ ΜΗΧΑΝΙΚΩΝ
ΚΑΙ ΜΗΧΑΝΙΚΩΝ ΥΠΟΛΟΓΙΣΤΩΝ
ΤΟΜΕΑΣ ΤΕΧΝΟΛΟΓΙΑΣ ΠΛΗΡΟΦΟΡΙΚΗΣ ΚΑΙ
ΥΠΟΛΟΓΙΣΤΩΝ

Machine Learning Techniques In Categorical Time Series Analysis Of Manufacturing Process

ΔΙΠΛΩΜΑΤΙΚΗ ΕΡΓΑΣΙΑ

ΧΑΡΑΛΑΜΠΟΣ ΜΙΧΑΗΛΙΔΗΣ, ΙΣΙΔΩΡΑ ΧΑΡΑ ΤΟΥΡΝΗ

Επιβλέπων : Νεκτάριος Κοζύρης
Καθηγητής Ε.Μ.Π.

Εγκρίθηκε από την τριμελή εξεταστική επιτροπή την 20η Ιουλίου 2016.

.....
Νεκτάριος Κοζύρης
Καθηγητής Ε.Μ.Π.

.....
Γεώργιος Γκούμας
Λέκτορας Ε.Μ.Π.

.....
Κώστας Καρπούζης
Ερευνητής Α' Ε.Π.Ι.Σ.Ε.Υ.

Αθήνα, Ιούλιος 2016

.....
Χαράλαμπος Μιχαηλίδης, Ισιδώρα Χαρά Τουρνή

Διπλωματούχοι Ηλεκτρολόγοι Μηχανικοί και Μηχανικοί Υπολογιστών Ε.Μ.Π.

Copyright © Χαράλαμπος Μιχαηλίδης, Ισιδώρα Χαρά Τουρνή, 2016.

Με επιφύλαξη παντός δικαιώματος. All rights reserved.

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

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

Abstract

The complexity of modern industrial processes and the constant innovations in production monitoring technologies and data collection, strongly outline the need for advancements in production data analysis. Data Mining is a rapidly growing field, aiming in understanding data and extracting previously unknown information, with the use of Machine Learning techniques, in order to optimize production.

Our cooperation with Johnson & Johnson, enabled us to obtain and explore real case production data. The exploitation of them focused on two distinct goals. The first one was the visualization of the data and the graphical representation of all variables characterizing the mixing process, for an improved data overview. The second one was the Machine Learning algorithms' modification and application on the properly pre-processed production data, to look into the possibilities of these techniques in the enterprise space.

Machine Learning, by definition, sets to represent data as objects in space, utilizing labels and distances between them. In this direction, data were grouped into objects and vectorized with various techniques. Classification and Clustering algorithms were parameterized and implemented, investigating unique attributes of the provided data. Distance calculation methods for each algorithm were examined in depth, and each experiment was assessed through different evaluation metrics, in order to examine the performance of the algorithms and the result's accuracy, compared to the initial data. Our conclusions indicate that Machine Learning can drive important business decisions, through process quantification, and further research can be done in each specific case.

Key words

Machine Learning, Variable Length Classification, Categorical Time Series, Manufacturing Process, Transition Matrix

Περίληψη

Η πολυπλοκότητα των σύγχρονων βιομηχανικών διεργασιών και οι συνεχείς καινοτομίες στον τομέα των τεχνολογιών παρακολούθησης της παραγωγής και της συλλογής δεδομένων, τονίζουν έντονα την ανάγκη για πρόοδο στον τομέα ανάλυσης δεδομένων παραγωγής. Το Data Mining είναι ένας ταχέως αναπτυσσόμενος τομέας, έχει ως κύριο στόχο την κατανόηση δεδομένων και την εξαγωγή προηγουμένως άγνωστης πληροφορίας, με τη χρήση τεχνικών μηχανικής μάθησης, με στόχο τη βελτιστοποίηση της παραγωγής.

Η συνεργασία μας με την Johnson & Johnson, μας έδωσε τη δυνατότητα να εξερευνήσουμε δεδομένα παραγωγής μιας πραγματικής περίπτωσης. Η ανάλυσή τους επικεντρώθηκε σε δύο διαφορετικούς στόχους. Ο πρώτος ήταν η οπτικοποίηση των δεδομένων και η γραφική παράσταση όλων των μεταβλητών που χαρακτηρίζουν την διαδικασία ανάμιξης, για να έχουμε μια βελτιωμένη επισκόπησή τους. Ο δεύτερος ήταν η τροποποίηση και εφαρμογή αλγορίθμων Μηχανικής Μάθησης, στα κατάλληλα προ-επεξεργασμένα δεδομένα της παραγωγής, καθώς επίσης και η εξέταση των δυνατοτήτων των τεχνικών αυτών στον χώρο της Βιομηχανίας.

Η Μηχανική Μάθηση, εξ ορισμού, θέλει μια αντιπροσώπευση των δεδομένων ως αντικείμενα στο χώρο, χρησιμοποιώντας τις ετικέτες και τις αποστάσεις μεταξύ τους για να τα αναλύσει. Σε αυτή την κατεύθυνση, τα δεδομένα ομαδοποιήθηκαν σε αντικείμενα και διανυσματοποιήθηκαν με διάφορες τεχνικές. Αλγόριθμοι Classification και Clustering παραμετροποιήθηκαν και υλοποιήθηκαν, για τη διερεύνηση μοναδικών χαρακτηριστικών των παρεχόμενων δεδομένων. Επιπρόσθετα, μέθοδοι υπολογισμού αποστάσεων για κάθε αλγόριθμο εξετάστηκαν σε βάθος, και κάθε πείραμα εκτιμήθηκε χρησιμοποιώντας διαφορετικές μεθόδους αξιολόγησης, προκειμένου να εξεταστεί η απόδοση των αλγορίθμων και η ακρίβεια του αποτελέσματος, σε σύγκριση με τα αρχικά δεδομένα.

Τα συμπεράσματά μας δείχνουν ότι οι αλγόριθμοι Μηχανικής Μάθησης μπορούν να αξιοποιηθούν και να βοηθήσουν στο να παρθούν σημαντικές επιχειρηματικές αποφάσεις, μέσω της ποσοτικοποίησης των παραγωγικών διαδικασιών, καθώς επίσης, περαιτέρω έρευνα μπορεί να γίνει σε κάθε εξειδικευμένη περίπτωση.

Λέξεις κλειδιά

Machine Learning, Variable Length Classification, Categorical Time Series, Manufacturing Process, Transition Matrix

Acknowledgements

As this Thesis Project comes to an end, we remember all the people that helped us make it through and we feel the need to express our gratitude to them, because this would not be possible without their guidance and support.

We would like to thank all members of the Computing Systems Laboratory of National Technical University of Athens, especial Professor Koziris who enabled and trusted us to fulfill such a complex Project, involving many stakeholders. Of course we would like to thank Dr. Konstantinou how was our guide in this trip, supervising us patiently throughout the process. We also can not forget the help of Giannis Giannakopoulos whenever needed, regarding the infrastructure where we ran our experiments.

This Project wouldn't be a reality, and wouldn't be so interesting if it wasn't for Michalis Augoulis who connected us with Johnson & Johnson Hellas, providing us with real data which lead us into facing a real-case scenario. Moreover we deeply appreciate the openness and willingness of Johnson & Johnson Hellas' stuff throughout the process.

In addition to the above, kudos for the crucial support and inspiration goes to Dr. Anagnostopoulos and Mr. Cotsikis from the company Mentat. To Mr. Cotsikis for initially believing in our vision and to Dr. Anagnostopoulos for the truly inspiring sessions we had and his crystal-clear way of addressing and solving complex problems.

Last but certainly not least, this is dedicated to our families and friends, who not only helped us the past months but the last 6 years.

Athens, July 20, 2016

This thesis is also available as Technical Report CSD-SW-TR-42-14, National Technical University of Athens, School of Electrical and Computer Engineering, Department of Computer Science, Software Engineering Laboratory, July 2016.

URL: <http://www.softlab.ntua.gr/techrep/>
FTP: <ftp://ftp.softlab.ntua.gr/pub/techrep/>

Ευχαριστίες

Δεδομένου ότι αυτή η Διπλωματική Εργασία, έρχεται στο τέλος της, θυμόμαστε όλους εκείνους τους ανθρώπους που μας βοήθησαν να πραγματοποιηθεί και αισθανόμαστε την ανάγκη να εκφράσουμε την ευγνωμοσύνη μας σε αυτούς, καθώς η εργασία αυτή δεν θα είχε πραγματοποιηθεί χωρίς την καθοδήγηση και υποστήριξη τους.

Θα θέλαμε να ευχαριστήσουμε όλα τα μέλη του Εργαστηρίου Υπολογιστικών Συστημάτων του Εθνικού Μετσόβιου Πολυτεχνείου της Αθήνας, ιδιαίτερα τον Καθηγητή κ. Κοζύρη που μας εμπιστεύτηκε για να αναλάβουμε ένα τόσο σύνθετο έργο, με τη συμμετοχή πολλών ενδιαφερομένων μελών. Φυσικά θα θέλαμε να ευχαριστήσουμε τον Δρ. Κωνσταντίνου ο οποίος ήταν οδηγός μας στο ταξίδι αυτό και μας επέβλεπε υπομονετικά καθ' όλη τη διάρκειά του. Επίσης, δεν μπορούμε να ξεχάσουμε τη βοήθεια του Γιάννη Γιαννακόπουλου, ο οποίος μας βοήθησε όλες τις στιγμές που χρειαστήκαμε, όσον αφορά την υποδομή στην οποία τρέξαμε τα πειράματά μας.

Η Διπλωματική αυτή δεν θα ήταν πραγματικότητα, και δεν θα είχε τόσο μεγάλο ενδιαφέρον, αν δεν ήταν ο Μιχάλης Αυγουλής ο οποίος μας έφερε σε επαφή με την Johnson & Johnson Hellas, παρέχοντάς μας πραγματικά δεδομένα τα οποία και αξιοποιήσαμε για να αντιμετωπίσουμε ένα πραγματικό σενάριο παραγωγής. Επιπλέον εκτιμάμε βαθύτατα την προθυμία των εργαζομένων της Johnson & Johnson Hellas καθ' όλη τη διάρκεια της συνεργασίας μας.

Εκτός από τα παραπάνω, θέλουμε να ευχαριστήσουμε τον Δρ. Αναγνωστόπουλο και τον κ. Κοτσίκη από την εταιρεία Mentat, για την κρίσιμη στήριξη και την έμπνευση που μας έδωσαν. Τον κ. Κοτσίκη για την αρχική πίστη στο όραμά μας και τον Δρ. Αναγνωστόπουλο για τις συναντήσεις που είχαμε και μας ενέπνευσαν πραγματικά καθώς και για τον πεντακάθαρα τρόπο αντιμετώπισης και επίλυσης σύνθετων προβλημάτων.

Τέλος, αυτή η εργασία είναι αφιερωμένη στις οικογένειες και στους φίλους μας, που δεν μας βοήθησε μόνο τους τελευταίους μήνες, αλλά τα τελευταία 6 χρόνια.

Χαράλαμπος Μιχαηλίδης, Ισιδώρα Χαρά Τουρνή,

Αθήνα, 20η Ιουλίου 2016

Η εργασία αυτή είναι επίσης διαθέσιμη ως Τεχνική Αναφορά CSD-SW-TR-42-14, Εθνικό Μετσόβιο Πολυτεχνείο, Σχολή Ηλεκτρολόγων Μηχανικών και Μηχανικών Υπολογιστών, Τομέας Τεχνολογίας Πληροφορικής και Υπολογιστών, Εργαστήριο Τεχνολογίας Λογισμικού, Ιούλιος 2016.

URL: <http://www.softlab.ntua.gr/techrep/>

FTP: <ftp://ftp.softlab.ntua.gr/pub/techrep/>

Περιεχόμενα

Abstract	5
Περίληψη	7
Acknowledgements	9
Ευχαριστίες	11
Περιεχόμενα	13
1. Introduction	17
1.1 Data Analysis in Manufacturing	17
1.2 Problem Motivation	17
1.3 Proposed and Implemented Solutions	18
1.4 Thesis Structure	20
2. Data and Problem Description	21
2.1 Problem Description	21
2.1.1 Case Description	21
2.1.2 Dataset Description	22
2.1.3 Definition of sub-Problems	24
2.2 Data Pre-Processing	25
2.2.1 Flattening - Cleaning	25
2.2.2 Labelling	25
2.2.3 Object Creation for Machine Learning	25
2.3 Data Visualization	27
2.3.1 Infrastructure and methodology	27
2.3.2 Chart Explanation	27
3. Machine Learning	33
3.1 Data Mining and Machine Learning	33
3.2 Classification	34
3.2.1 Instance - Based Learning	34
3.3 Clustering	35
3.3.1 Partitioning Methods-k-Means algorithm	36
3.4 Distances	37
3.4.1 Minkowski Metric	37
3.4.2 Cosine Distance	37
3.4.3 Kullback-Leibler Divergence	37
3.4.4 Kolmogorov-Smirnov Test	38
3.4.5 Infinite Norm	38
3.5 Evaluation	38
3.5.1 Classification Evaluation	38

3.5.2	Clustering Evaluation	39
4.	Implementation	41
4.1	Classification	41
4.1.1	Nearest Centroid Classifier	42
4.1.2	K - Nearest Neighbors Classifier	42
4.2	Clustering	42
4.2.1	K - Means Implementation	42
4.3	Distances	43
4.3.1	Implementation of Distance Methods	44
4.3.2	Evaluation of Distance Methods	44
4.4	Outcome Evaluation	45
4.4.1	Implementation of Classification Evaluations	45
4.4.2	Implementation of Clustering Evaluations	45
5.	Results	47
5.1	Classification	47
5.1.1	Baseline	48
5.1.2	Distance Method Evaluation and Selection	48
5.1.3	Nearest Centroid and k-Nearest Neighbors Algorithms Comparison	51
5.1.4	k-Nearest Neighbors Algorithm for Different k	54
5.2	Clustering	55
5.2.1	Baseline	55
5.2.2	Distance Algorithms	55
5.2.3	Centroids	58
6.	Epilogue	61
6.1	Conclusions	61
6.1.1	Classification	61
6.1.2	Clustering	61
6.2	Future Work	62
7.	Εισαγωγή	65
7.1	Ανάλυση Δεδομένων στην Παραγωγή	65
7.2	Ορισμός Προβλήματος	66
7.3	Προτεινόμενες Λύσεις και Υλοποίηση	67
7.4	Δομή Διπλωματικής Εργασίας	68
8.	Περιγραφή του Προβλήματος και των Δεδομένων	69
8.1	Περιγραφή Προβλήματος	69
8.1.1	Περιγραφή Θέματος	69
8.1.2	Περιγραφή του αρχικού Συνόλου Δεδομένων	70
8.1.3	Ορισμός των υποπροβλημάτων	72
8.2	Προεπεξεργασία των Δεδομένων	73
8.2.1	Flattening - Καθαρισμός	73
8.2.2	Τιτλοφόρηση	73
8.2.3	Δημιουργία Αντικειμένων για τη Μηχανική Μάθηση	74
8.3	Οπτικοποίηση των Δεδομένων	75
8.3.1	Υποδομές και Μεθοδολογία	76
8.3.2	Επεξήγηση του Διαγράμματος	76

9. Μηχανική Μάθηση	81
9.1 Εξόρυξη δεδομένων και Μηχανική Μάθηση	81
9.2 Classification	82
9.2.1 Instance - Based Learning	82
9.3 Clustering	83
9.3.1 Partitioning Methods- K - Means algorithm	84
9.4 Αποστάσεις	85
9.4.1 Minkowski Metric	85
9.4.2 Cosine Distance	85
9.4.3 Kullback- Leibler Divergence	85
9.4.4 Kolmogorov- Smirnov Test	85
9.4.5 Infinite Norm	86
9.5 Αξιολόγηση	86
9.5.1 Αξιολόγηση Classification	86
9.5.2 Αξιολόγηση Clustering	87
10. Υλοποίηση	89
10.1 Classification	89
10.1.1 Nearest Centroid Classifier	90
10.1.2 K - Nearest Neighbors Classifier	90
10.2 Clustering	91
10.2.1 Υλοποίηση του k-Means	91
10.3 Αποστάσεις	91
10.3.1 Υλοποίηση Μεθόδων Αποστάσεων	92
10.3.2 Αξιολόγηση των Μεθόδων Αποστάσεων	93
10.4 Αξιολόγηση των Αποτελεσμάτων	93
10.4.1 Υλοποίηση της Αξιολόγησης του Classification	93
10.4.2 Υλοποίηση της Αξιολόγησης του Clustering	94
11. Αποτελέσματα	95
11.1 Classification	95
11.1.1 Baseline	95
11.1.2 Αξιολόγηση Μεθόδου Απόστασης	96
11.1.3 Σύγκριση Nearest Centroid και K-Nearest Neighbors	98
11.1.4 Αλγόριθμος K-Nearest Neighbors για διαφορετικά k	101
11.2 Clustering	102
11.2.1 Baseline	102
11.2.2 Μέθοδοι Υπολογισμού Αποστάσεων	102
11.2.3 Αρχικά κέντρα	104
12. Επίλογος	107
12.1 Συμπεράσματα	107
12.1.1 Classification	107
12.1.2 Clustering	107
12.2 Μελλοντική Εργασία	108
Παράρτημα	111
A. Scripts	111
A.1 Classification	111
A.1.1 Nearest Centroid Classifier	111

A.1.2	K - Nearest Neighbors Classifier	111
A.2	Clustering	112
A.2.1	K - Means Implementation	112
A.3	Distances	113
A.3.1	2D to Vector Conversion	113
A.3.2	Distance Algorithms Implementation	113
A.3.3	Evaluation of Distance Methods in Classification	117
A.3.4	Evaluation of Distance Methods in Clustering	117
	Bibliography	119
	Κατάλογος σχημάτων	121

Κεφάλαιο 1

Introduction

1.1 Data Analysis in Manufacturing

Modern industrial systems and processes have become very complex. The information technologies' development, and the application of high - performance computers in all branches of an industrial procedure, result in a vast amount of data being produced continuously. An industrial process is being monitored by different sensors at various times, frequencies and resolutions, thus having an increasing complexity in all of its stages, and standing as a challenging problem for nowadays scientific disciplines. The technological advancements have improved production in all its aspects, yet they have made processes complex and nonlinear anymore, so it is rather hard, or impossible, to gather precise and direct information from measurement equipment only. Classical Analysis methods need to become more advanced, in order to lead the process effectively, and new techniques need to be integrated into old production methods or even replace them. [10]

In this direction, industries adopt and develop data analysis and exploitation strategies, in order to increase production quality, process security and human safety. This approach is not new, as computational and business intelligence ideas are being implemented since early 70s, including the application of artificial neural networks and predictive and adaptive system control, in an attempt to help human operators in a production line, and, if possible, eliminate them. This is where data mining comes to the spotlight, as a way to extract valid, previously unknown and comprehensible information from a process, and use it to make important business decisions. Like schematic drawings and mathematical equations were formerly used, to gain insightful knowledge, monitor, understand and optimize procedures, data analysis methods play a similar role, enhancing and boosting currently applied ones. The multidisciplinary nature of the field, which incorporates, between others, machine learning, image processing and statistics, aims in detecting regularities and patterns, invisible to humans or other analysis' methods.

Industries have the ability to use historical process data, identify relations among discrete process steps, as well as the determinants of a procedure's performance, and then optimize the factors that prove to have the greatest effect on production. The same rules apply when real - time data are examined, where more complex techniques should be developed, yet the immediate analysis provides the ability for short- term planning and actions, crucial to the improvement and effectiveness of the procedure. Visualization techniques are a crucial tool for pattern detection and prioritize data collection, as they are widely applied here to enhance the result, with the use of distribution graphs and clustering diagrams.

These needs and possibilities were the starting point of the current thesis, which begun as an attempt to identify production monitoring's problems and systemize data analysis and exploitation.

1.2 Problem Motivation

The multinational production environment of Johnson & Johnson Hellas fits in as an ideal field to experiment and implement all of the above. From the initial conversations, the interest from both sides regarding the project was high. The field of data analysis is vast and the applications in such an

environment are tremendous.

In this plant, each production line consists of many machines, a really important is the vessel in which the mixing takes place. In contrast with other production methodologies, this plant produces its products in a batch manner. This means that in each vessel, a great variety of products can be made, depending on the demand of the market. Mixing is a crucial step of production, because this is the part where the main product is made. As it will be explained in the next chapter in more detail, this vessel is being monitored and controlled by a PLC (Programmable Logic Controller), of which the raw output is the data-set provided to us.

The first need that was communicated to us, was that of the visualization of the production data. The data were logged on a daily basis, but there was no further analysis and utilization of them. Implementing a graphical depiction of data enabled us to have a better understanding of it and proceed with Machine Learning techniques.

The second one had many sub-goals, which included the assessment of the mixing process and the comparison towards a golden standard, having always in mind the optimization of the production process. Each product has the same very strict quality limits which are easily tracked through chemical analysis of the final product. The part, which is difficult to be measured and evaluated, is the execution of the recipe, while making a certain product. This is particularly interesting, because all the products produced, are within the quality limits, but there are no two identical time-series of actions. To provide an order of magnitude, each batch is completed in around 400 actions.

In order to be able to present a case of evaluating and comparing procedures, we proceeded with classification and clustering methods. We believe that this was a great first step of evaluating and experimenting on how data-series could be converted and used in machine learning implementations.

The main challenges, which derived during the process of utilization of the given data-set, were the following:

- **Data Cleansing**

The provided data-series from the PLC were almost perfect from an integrity point of view. On the contrary, the logged data of each vessel were written by humans, which made them pretty inaccurate.

- **Feature Selection**

In the data-set, for any given moment we had information from 17 different variables, regarding the state of the vessel, which increased their correlation complexity and made it difficult to come to a meaningful conclusion, just by observing their values.

- **Unequal length series comparison**

One of the most complex problems we had to overcome was the comparison of unequal length data-series. Each batch consists of actions, represented by their codes (MsgNum) in our data-set. Due to the fact that two same products can be produced by different actions, varying in sequence of actions and in their total number, the comparison could not be done one-to-one.

- **Distance Evaluation**

Most of the implementations of Machine Learning Algorithms manipulate elements that have a vector of attributes. We needed to decide on the elements, of which the distance we would calculate for each application, and also use the proper metrics for our data.

1.3 Proposed and Implemented Solutions

In order to tackle the above defined challenges and make best use of the given data, we followed the steps stated below, which form our solution:

- **Pre-processing and Data cleansing:**

Combining the initial PLC dataset and the human-logged data of each process, and, based on known company internal rules, we grouped separate actions together, to create objects representing an independent procedure that takes place. This was either a Production or Cleaning Process. We decided to only use one of the variables of the data in this direction, the Message Number (MsgNum), as, after a few experiments, we concluded that it was the best and simplest approach, which eliminated all complexity due to the large number of process variables.

- **Visualization:**

A visualization tool was implemented, to depict all different variables and their values with time, in order to create a graphic and interactive display of the initial data. This was an approach to better understand and correlate the data, and observe probable visible patterns of the initially unmanageable and incomprehensible consecutive timeseries.

- **Unequal length series comparison:**

We cast out the time relevance of the actions. This was done, because of the repeatability of them and the fact that many of the actions can be performed in consequent iterations, without interfering with the outcome. Having that in mind, our initial approach was to create a normalized frequency vector for all the messages in a batch. Yet we decided to include more information from the initial dataset. The final approach was to construct a normalized transition matrix for each batch. Each cell of this matrix indicates the probability of transitioning from one message code to another. To state it in a simple manner, cell[A][B] indicates the probability of executing the action B after the action A. This approach turned to be really useful and helped us through the Machine Learning part.

- **Distance Evaluation:**

A few methods to calculate distance between different 2D objects (matrices) was designed and implemented, in order to compare the objects we created, and find their relative distance. The first approach was to transform the matrix into a vector, and the second one was to create our own distance method and apply the selected algorithms through that. This was implemented for many known distance metrics (Euclidean, Cosine, KL- Divergence, KS-Test, Infinity Norm) and was further used in the Machine Learning algorithm's application. As our experiments evolved, we observed great differences in the results, depending on the distance method we had chosen.

- **Classification:**

We experimented with two Classification algorithms on the Production objects created after the Data pre-processing, Nearest Centroid Classifier and k- Nearest Neighbours Classifier. Each Product has different attributes (Product Group, Product Cleaning Group it belongs to), already known from the previous labelling, thus the classification was performed on these characteristics. The distance of the objects, needed for classification, was a variable in each execution. For the k-NN algorithm, the value of k was also a parameter examined. The data were trained and a testing data set was used to identify the performance of the training procedure, which was measured by two metrics, Accuracy and Kappa Coefficient. The experiments gave interesting results, with performance measurements of a maximum 85%.

- **Clustering:**

We customized the k - Means Clustering algorithm, and based again on the Production objects' different attributes, we assigned them to different clusters. The parameterization of the algorithm focused on two things, the selection of the initial centroids and the proper distance metric between the data objects used for clustering. Two metrics were utilized in the evaluation of the experiments, V-Measure and Rand-Index, and the performance of them proved to be quite poor, with a maximum score of 35%.

1.4 Thesis Structure

In Chapter 2, we define the problem we will work on, with all its different parameters, describe all steps taken towards data processing and understanding, and present the data visualization tool created for the purposes of J&J company from the production data.

In Chapter 3, we focus on the theoretical background needed for our Data Mining and Machine Learning applications. We define the above concepts, and analyze those of Classification and Clustering, focusing on the algorithms used in the current work. We also present the concepts of Distances and Evaluation in Machine Learning and metrics and techniques used for these purposes.

In Chapter 4, we present the Implementation of our solutions. Classification and Clustering are analytically parameterized and explained, as are Distance and Evaluation metrics applied in the experiments executed.

In Chapter 5, we present the results and diagrams from the experiments run for all Classification and Clustering algorithms, for different initial setups, and their implementations.

Finally, in Chapter 6, we summarize our conclusions from the results, and refer to issues we did not have the chance to work on and which could be regarded as future work.

Κεφάλαιο 2

Data and Problem Description

In this Chapter, we describe all characteristics of the industry problem we deal with in the current thesis, and the initial steps we took in order to utilize the production data provided to us, in an effective way.

In Section 2.1, we first describe the mixing process we are focusing on, by giving its most significant points. Then, we present all the parameters in the initial dataset, with a screenshot from it for better understanding, and we shortly define the steps followed in its processing and the main challenges tackled.

In Section 2.2, we describe the pre-processing techniques applied, which are divided in three categories, Flattening, Labelling and Object Creation.

Finally, in Section 2.3, we analyze the way the dataset was depicted graphically and present the visualization tool we created for the purposes of the company.

2.1 Problem Description

In this Section, the main parameters and variables of the problem will be presented.

2.1.1 Case Description

The case we studied is a real-life problem. The data have been provided by Johnson & Johnson Hellas, specifically, the data which were used in the Machine Learning Section, are from a single Vessel which is used in the batch production line of a certain factory in Greece. Each Vessel is used for a specific number of actions and can produce a great variety of products. We present below more details for each element.

- **Vessel**

This vessel has some automation mechanisms but most of the processes are being executed by the operator. There is no use in pointing exactly which the automated functions are and which are being executed by the user. The total data set, explained, can be found in the next sub-section.

- **Mixing/Production**

The main function of this vessel is the mixing of different materials in order to produce a certain product. The products vary from lotions to oils and creams. Some standard materials are being provided by a network of fluids and other, more rare, are being imported by hand by the operator.

- **Possible Actions**

The possible actions which can be performed are: Mixing, Heating, Cooling, Import of Fluids from the network, Import of materials by a trap door, Creation of vacuum pressure. Many of them can be performed at the same time, for example, mixing while heating.

- **Recorded Values**

As it will be described in the following subsection, the actions which were performed, are stored by a PLC as a time-series of messages with the appropriate measurements and time-stamps. There is one special characteristic regarding the format of the output CSV, which is the following. A set of values is being recorded only at the time of a specific action. For every variable we have two values, one is the actual, measured, value of the respective attribute and the other is the Set-Point of the attribute. Meaning that at the time that the Heating process has started, we have two values of temperature. The one is the measured, ex. 50, 87 and the other is the goal of the heating process, ex. 75. Those two sets are represented in the initial excel with two different rows. Most of the time it is easy to distinguish them, because Set-Point rows have mostly integer values, in contrast to the measurements, which have decimals.

2.1.2 Dataset Description

This is a really important part of our thesis, which is going to help the reader better understand the data pre-processing section. The data-set which was provided to us, is a raw output of the PLC controlling this particular vessel. The format is CSV and the data, which were used in the Machine Learning part, represent the total data from a whole year.

• Data-Set Information

- Format: CSV
- Time period: 16/02/2015 – 13/02/2016
- Rows of Data: 132005
- Filesize: 25Mb
- Attributes (Columns):
 - * StateAfter
 - * MsgNumber
 - * Temp
 - * Pressure
 - * Agitation1
 - * Agitation2
 - * Homogen
 - * Pump Power
 - * Raw Meterials
 - * TimeString

• Screenshot from the initial CSV

	A	B	C	D	E	F	G	H	I	J
1	StateAfter	MsgNumber	Temp	Pressure	Agitation1	Agitation2	Homogen	Pump Power	Raw Materials	TimeString
2	1	552	78.62	0	14.84028	24.75579	0	0.8101852		16.02.2015 09:41:08
3	1	564	78.62	0	14.84028	24.75579	0	0.8101852		16.02.2015 09:41:08
4	1	553	75	-500	15	30	750	60		16.02.2015 09:41:08
5	1	565	75	-500	15	30	750	60		16.02.2015 09:41:08
6	1	552	78.89	1.15741	14.84028	24.75579	0	0.8391203		16.02.2015 09:43:39
7	1	564	78.89	1.15741	14.84028	24.75579	0	0.8391203		16.02.2015 09:43:39
8	0	553	75	-500	15	30	750	60		16.02.2015 09:43:39
9	0	565	75	-500	15	30	750	60		16.02.2015 09:43:39
10	1	522	50.715	1.15741	6.534722	25.62182	0	0.8101852	21.20999	16.02.2015 10:25:43
11	1	523	75	-500	15	31	750	60	500	16.02.2015 10:25:43
12	0	536	77.1	0	14.83333	24.77691	0	0.8101852		16.02.2015 10:25:43
13	0	562	77.1	0	14.83333	24.77691	0	0.8101852		16.02.2015 10:25:43
14	1	537	75	-500	15	30	750	60		16.02.2015 10:25:43
15	1	563	75	-500	15	30	750	60		16.02.2015 10:25:43

Σχήμα 2.1: Initial Data-Set format in CSV format

- **Explanation of important data-set variables**

- *StateAfter*: This variable can be either 0 or 1. When the value is 1, it means that the action corresponding to the message number is starting. When it is 0, the process is ending.
- *MsgNumber*: Through a mapping file we can map the message numbers to specific actions performed in the vessel. For example, some of the messages and their meaning can be seen below.

	A	B
1	500	ΑΝΑΚΥΚΛΟΦΟΡΙΑ ΑΠΌ ΠΑΝΩ
2	501	ΑΝΑΚΥΚΛΟΦΟΡΙΑ ΑΠΌ ΠΑΝΩ (Set Points)
3	504	ΑΔΕΙΑΣΜΑ ΣΕ TNT
4	505	ΑΔΕΙΑΣΜΑ ΣΕ TNT (Set Points)
5	522	ΕΙΣΑΓΩΓΗ ΑΠΙΟΝΙΣΜΕΝΟΥ ΝΕΡΟΥ
6	523	ΕΙΣΑΓΩΓΗ ΑΠΙΟΝΙΣΜΕΝΟΥ ΝΕΡΟΥ (Set points)

Σχήμα 2.2: Example MsgNumbers

- *TimeString*: This is the time-stamp, down to seconds accuracy, at the beginning or end of each action.

- **Splitting**

Having a continuous time-series of action messages did not provide us any information regarding the beginning and end of each batch. This was an important step, that enabled us to apply the Machine Learning techniques, as explained in the Implementation Chapter.

The splitting of the data-set into "chunks" was done by applying certain rules that were explained to us by senior members of the production management department.

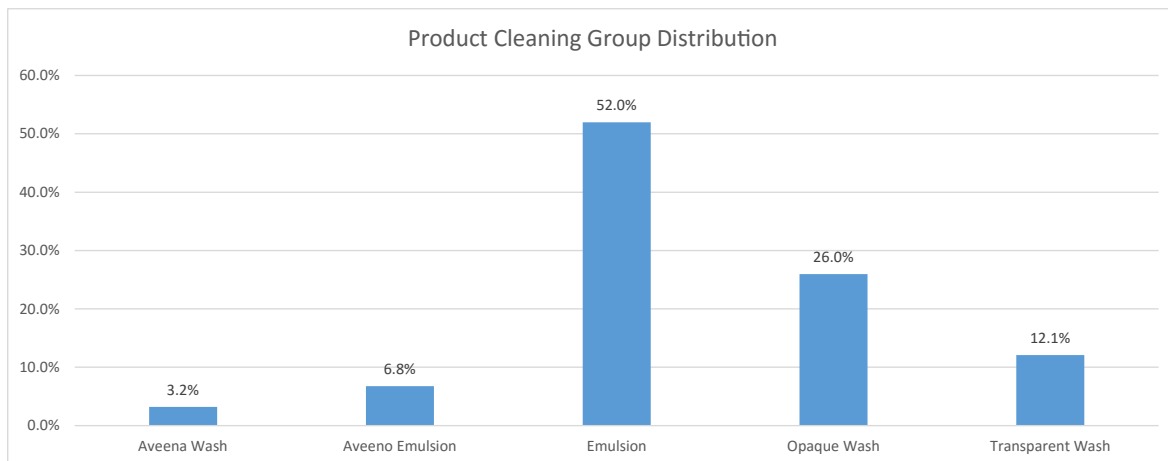
- **Labelling**

Another really important step of the process was the Labelling of the data-set. In the initial data-set there was no information regarding the product that was being produced at any given moment. In order to label the splitted chunks, two files were provided to us.

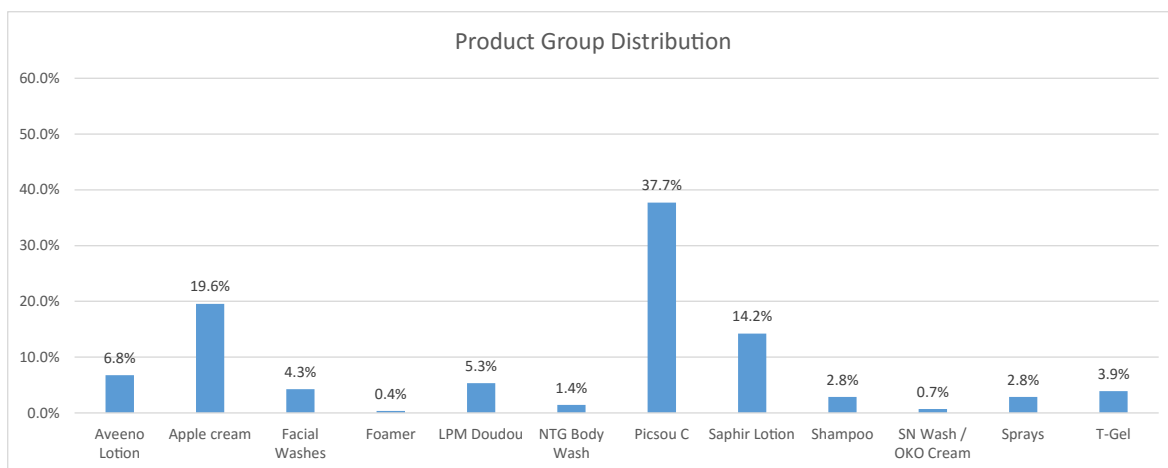
The first one was the logging file of the operators, where each one submitted the end time and the product code of each batch. The second one was a mapping file, which enabled us to retrieve based on the product code, two crucial elements of our analysis, the Product Cleaning Group and the Product Group.

Those two categorizations will play a major role in both classification and clustering.

- The Product Cleaning Group consists of 5 groups and all the products belong to one of them. The exact distribution can be seen on Figure 8.3
- The Product Group is a different categorization which consists of 12 groups, like the previously mentioned, all products belong to one of them. The exact distribution can be seen on Figure 8.4



Σχήμα 2.3: Product Cleaning Group Distribution



Σχήμα 2.4: Product Group Distribution

2.1.3 Definition of sub-Problems

Two of the most important problems we faced, and tackled, were the following.

- **Variable selection**

The initial data-set provided us many variables to work with. After some first tests we decided to move forward utilizing only the message number. In Section 6.2 we mention that a great deal of work can be done in the sector of exploring Machine Learning Techniques using different variables or combination of them.

- **Comparison of non-equal length data-series**

As one understands, each batch can be produced by different combinations of actions. This is due to the fact that each operator can fine-tune the process according to his beliefs. Each process consists of around 400 messages, the instructions that are being provided by the chemists of the company are not so detailed and the materials used have some tolerance limits. Those factors allow the operator to take initiatives in the production process of each batch.

This characteristic of the process makes every batch unique, regarding the length of the time-series and the exact order of the actions. As it will become clear in Chapter 4, the compared elements must have same number of attributes in order to be classified or clustered. The idea

we implemented is to construct a transitions matrix for each chunk. This matrix has always size $M \times M$, where M is the number of unique message codes, which is the same for every batch. In our case, $M = 45$.

2.2 Data Pre-Processing

In this section all the methods used for data pre-processing will be analyzed. The role of this part of the thesis is to clean the initial data-set and bring it into the proper format in order to be easier to visualize the data-set as well as to apply the Classification and Clustering techniques.

2.2.1 Flattening - Cleaning

As explained in subsection 2.1.1 the data come in pairs of measurement and Set-Point, so we developed a flattening procedure in order to have one row for each action, containing both all the measurements and all the set-points. Alongside, we cleansed the data from test values that were entered by operators. This part was done in order to help the visualization process.

2.2.2 Labelling

As previously mentioned, Labelling was an important and necessary step, because through the labels we were able to evaluate the performance of both classification and clustering.

The files we used in this step were three: the logging file from the operators, the mapping file between product codes and product attributes, and the chunk series containing all the initial data, splitted into chunks.

The first step was to assign the proper attributes to the chunks found on the log file. This was easy, because we had the unique product code from the log file, and retrieved all the necessary information from the mapping file. The only problem found was the rare case of some typos in the production codes, as found on the log file, which did not enable us to find its attributes. The attributes retrieved were the Production Cleaning Group and the Production Group.

The second step was to assign the products as seen in the log file with the chunks from the initial data-set. While the task seemed easy to undertake, it proved to be more demanding than expected. The way that this assignment was done, was to find the end time of a certain batch in the log file and then locate, in the data-set, the chunk which had the nearest end time. Three were the main problems in this procedure.

- The inconsistent way the time-stamps were inserted in the log file. This is a process executed by hand which means that there were a lot of typos and non-usable information.
- Because the splitting of the initial data-set was done by non-perfect rules, some chunks were merged together and others split into two.
- The inconsistency between the two lists. In order to make the assignment we implemented an algorithm, which assigned the labels to the nearest chunk. Because many outliers were found, we decided to use, for the Machine Learning part, only the chunks which had a time difference of less than 7 *hours*, between the end time found in the log and the one found from the data-set.

2.2.3 Object Creation for Machine Learning

The last step of the data pre-processing was to construct a new data-set containing the previously labeled chunks in a usable format. We decided to create the **Chunk class** and make each batch an instance of this class. When all of them were converted, we stored them in a json file, for easier access and to ensure integrity of our procedure.

• Chunk Class description

```
1 class Chunk(object):
2
3     def __init__(self, **entries):
4         super(Chunk, self).__init__()
5
6         #The start time of this chunk as found on the initial data-set
7         self.start_time = None
8         #The end time of this chunk as found on the initial data-set
9         self.end_time = None
10        #The chunk_type can be either "Production" or "Cleaning"
11        self.chunk_type = None
12
13        #The product code as found from the log file.
14        self.pr_code = None
15        #The end time as found from the log file
16        self.pr_logged_end_time = None
17
18        #The following attributes were found from the referance file through the
19        #production code
20        self.pr_name = None
21        self.pr_group = None
22        self.pr_cl_group = None
23
24        #This is used for both classification and clustering
25        #This attribute is set to the name of the class/cluster it belongs
26        #and is being evaluated at the end of each experiment
27        self.cluster = None
28
29        #If this chunk represents the center of a cluster,
30        #this variable is set to its name
31        self.name = None
32
33        #This is the Transition Matrix of this chunk
34        self.TM = None
35
36        #This function is used by create_TM in order to initiate the Transition
37        #Matrix
38        def init_TM(self, num_msgs):
39            self.TM = [[0 for _ in range(num_msgs)] for _ in range(num_msgs)]
40
41        #This function creates the Transition Matrix of the chunk
42        #Given the message sequence and a message dictionary
43        #1)initialises the Transitions Matrix
44        #2)for each transition from message A to message B increases by 1 the
45        #respective cell
46        #3)normalises the matri by row
47        def create_TM(self, msg_sequence, msgs_dict):
48            self.init_TM(len(msgs_dict))
49            cur_msg = msg_sequence[0]
50
51            for i in range(len(msg_sequence)-1):
52                past_msg = cur_msg
53                cur_msg = msg_sequence[i+1]
54                self.TM[self.msgs_dict[past_msg]][self.msgs_dict[cur_msg]] += 1
```

```

53     for row_id, row in enumerate(self.TM, 0):
54         row_sum = sum(row)*1.0
55         self.TM[row_id] = [0.0 if row_sum == 0.0 else cell/row_sum for cell
in row]

```

Listing 2.1: Chunk Class Implementation

- **Chunk List creation and export to JSON**

Each chunk found in the initial splitted data-set, was converted into a chunk object. The chunks are of two main types "cleaning chunk" and "production chunk". As one can see in the implementation of the chunk class, certain attributes were filled during the creation of each object, and others left blank for further use. One of the most significant attributes of the chunks is the Transition Matrix, which was created from the action message sequence of each chunk, after that was split from the initial data set.

Once all the message sequences were labeled and converted into chunk objects, we saved that list in order to have a solid data-set for our further experimentation and not to waist time re-running the algorithms. The output was a JSON file, which proved to be very useful later on.

2.3 Data Visualization

In the initial conversation with Johnson & Johnson, it was clearly communicated that a crucial need for them, was the ability to visualize the data from the production process. This task was very difficult and sometimes impossible through the raw output, due to the fact that each batch might be composed from up to 800 rows of data. For this purpose, we created a visualization tool, to graphically display and associate all data variables provided. This tool is already in use and has been proven very helpful in monitoring the production process.

2.3.1 Infrastructure and methodology

The visualization of data variables was implemented in charts with the use of a Javascript Library. We combined various chart features to create and customize our visualization tool. In order to achieve the visualization of the raw data file, many steps had to be completed.

The steps that the user has to do are, the upload of the CSV file and the beginning of the pre-processing application from the Landing page of the Tool, as seen in the next figures 8.5 . This interface was created with PHP and HTML. From this step on, the process was automated and the user only has to wait a couple of seconds before he is able to see and interact with the chart.

The pre-processing begins with reading the input CSV file, the file is being cleaned from missing values and re-arranged in a more suitable format for the visualization, as explained in flattening Subsection 2.2.1. The previous part was implemented only with the use of Python. After the file is ready in memory, it is stored in a database table. For this part we used the WAMP server, which enabled us to have a mySQL database.

Once everything is ready and stored in the database, we prompt the user to go to the chart page. This page is the end result of our visualization interface, made possible by Highcharts, with the appropriate modifications from the initial templates. You can find more details and screenshots of the tool in the next subsection.

2.3.2 Chart Explanation

- **Description**

Our goal was to visualize all variables in one chart, so that one will be able to view all correlations and dependencies between them, choose which ones to depict in the charts and hide the others, and obtain information for all possible variable combinations. The chart provides the facility of showing or hiding each separate variable, by clicking its name at the bottom of the screen, which then displays or hides both the chart and its respective values' vertical axis. Moreover, the user can zoom in the specific time frame he is interested in, either by typing the exact dates to be drawn in the upper right boxes, or by choosing one of the options in the upper left corner, or, finally, by moving the bar at the bottom of the chart accordingly. Mousing over the chart, at the time points that were in the initial file, there appears a tooltip table, showing the information for the current timestamp, meaning the variables and their values, and also the message with the process description taking place at that specific time. We also added an extra variable, in the form of flags, which is placed at the bottom of the chart and on mouseover displays the message of the process that occurred at that moment. Some screenshots of charts viewable in the tool are presented below.

• Screenshots

Below can be found the previously mentioned screenshots, which can help the reader better understand the tool we created for visualizing the data from the production process.

Visualisation Tool

NTUA Thesis Project in cooperation with Johnson & Johnson Hellas

Developed by:

Haris Michailidis haris.michailidis@gmail.com

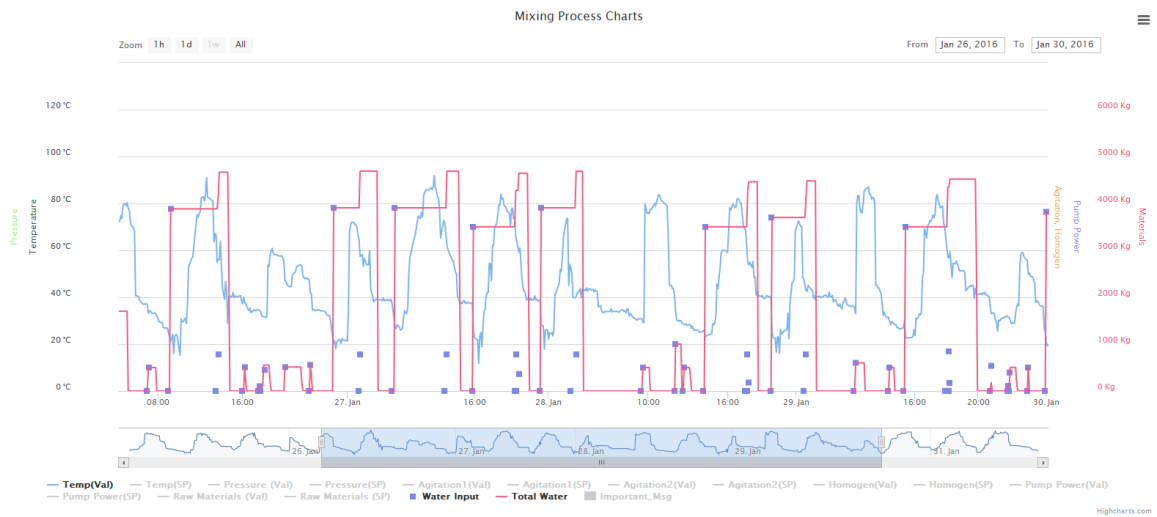
Isidora Tourni isidora.tourni@gmail.com

Instructions

- 1) Copy the proper CSV to the VT_input folder
- 2) Rename it to "input.csv"
- 3) Click on the "Run Main Script" button
- 4) Wait for ~1min
- 5) The Log and an extra button will appear
- 6) Click on the "Go to Chart" button

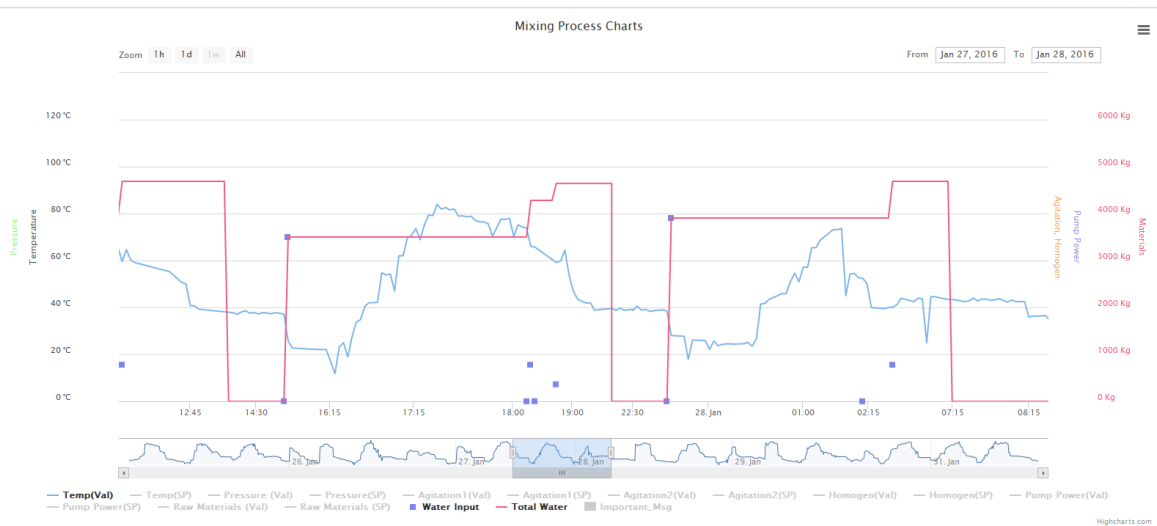
Σχήμα 2.5: Visualization Tool - Landing Page

This screenshot demonstrates the initial page of the tool, where the user can read the instructions. When he has finished placing the file in the correct folder he clicks the "Run Main Script" button.



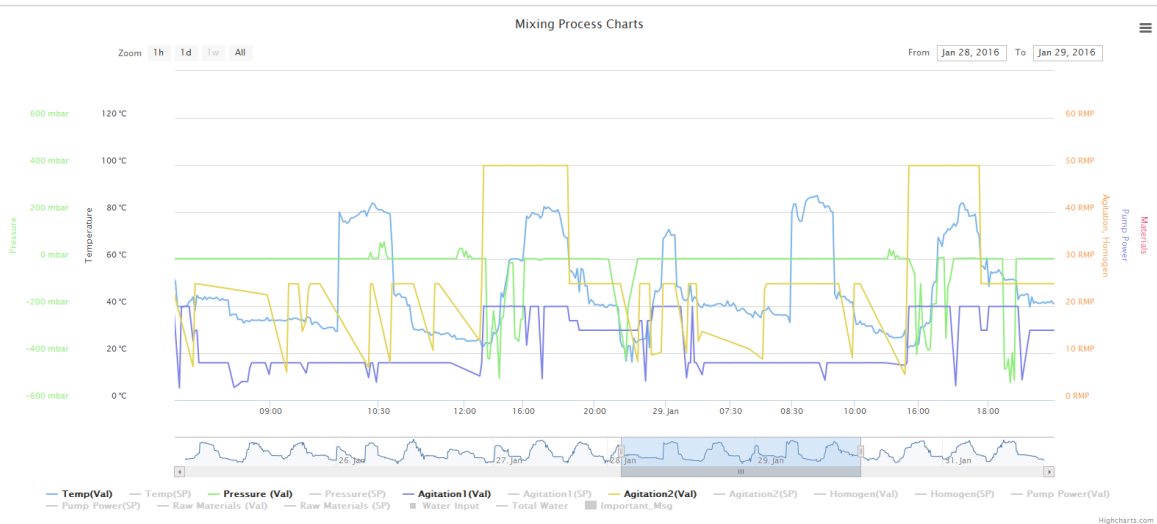
Σχήμα 2.6: Visualization Tool - Temperature, Water - 4 days

In this screen, the user is interested in tracking the temperature and the total water in the vessel. The time window is relatively large, it has an overview of 4 days.



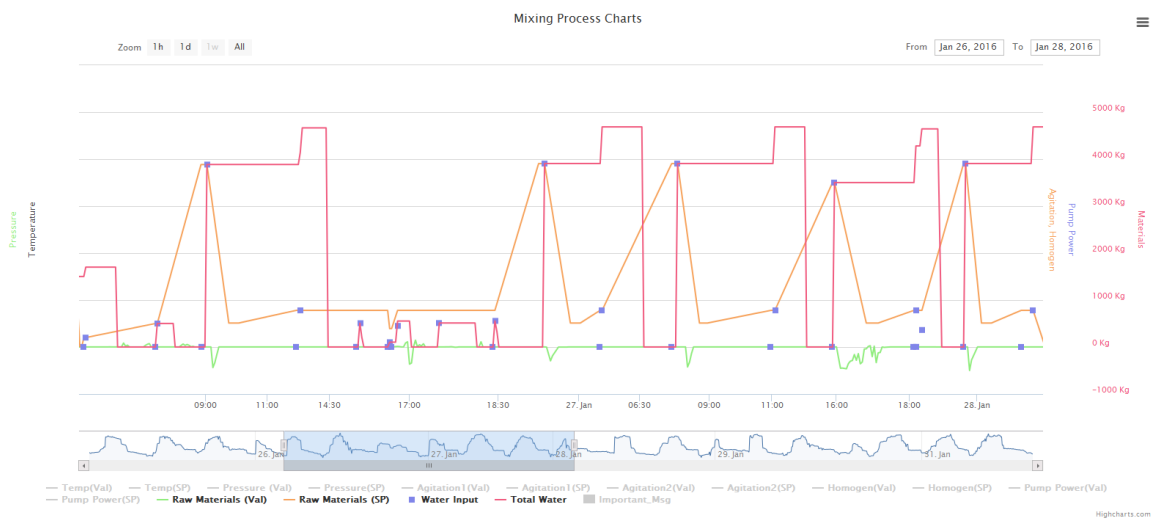
Σχήμα 2.7: Visualization Tool - Temperature, Water - 1 day

This is the same screen as in the previous figure, except for the time window which, in this case, is more narrow and displays only one day of activity.



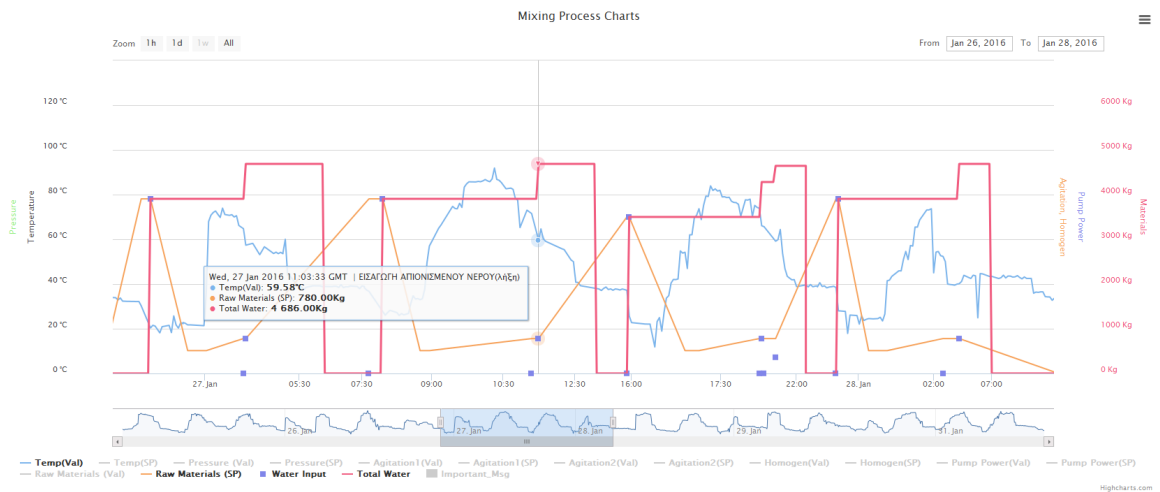
Σχήμα 2.8: Visualization Tool - Four Variables Correlation

In this screen, the user has enabled four separate graphs to be displayed, temperature, Pressure, Agitation 1 and Agitation 2. Through this graph, the user can better understand the relation between those variables during the mixing procedure.



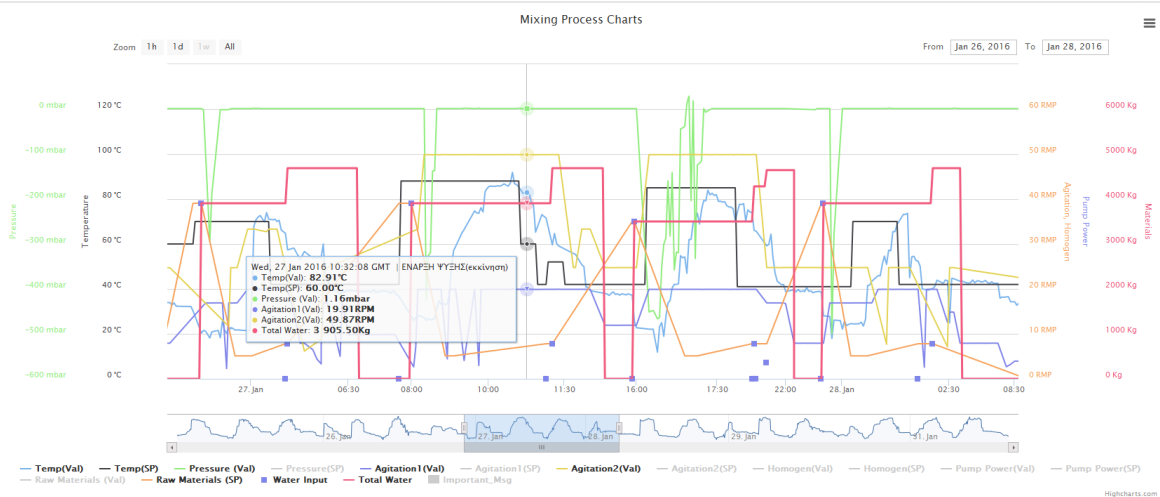
Σχήμα 2.9: Visualization Tool - Raw Materials Addition

One of the most important elements in the mixing process is the water addition. In this graph we can see with purple squares the absolute water addition is each moment. With the red line we observe the sum of all added water, till the moment it is drained from the vessel, where the red line goes to zero.



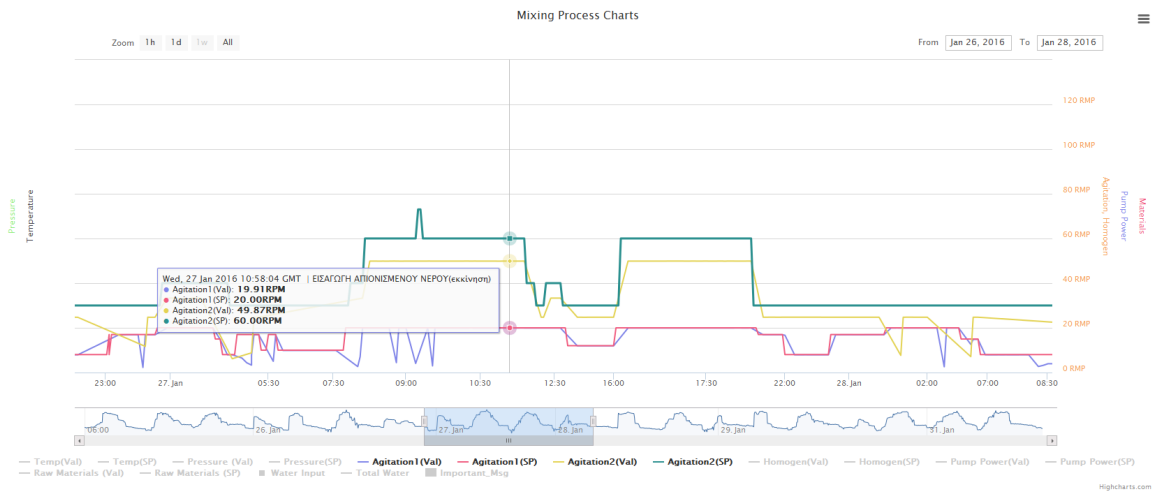
Σχήμα 2.10: Visualization Tool - Action Message Presentation

In order for the user to be able to examine the process in detail, we have implemented a tooltip which includes all the values of the current active variables of the chart for any given moment. The user can hover the mouse over the graph and see the exact values of each variable, along with the action message of that moment.



Σχήμα 2.11: Visualization Tool - Variable Values Presentation

As presented in the previous screenshot, in the box that appears aside of the measurements, all the active variables can be displayed, enabling the user to have a crystal clear overview of each moment of the process.



Σχήμα 2.12: Visualization Tool - Values and Set Points

In this screenshot, the concept of set-points is being presented. The set-point is a value, set by the user, which acts as a guide for the process. As one can see, the SP values are always ahead of the measured ones.

Κεφάλαιο 3

Machine Learning

In this Chapter, we provide the theoretical background of all Machine Learning techniques applied, and their parameters, and define the concepts that are needed to understand the implementation methods.

In Section 3.1, we introduce the reader to the fields of Data Mining and Machine Learning, and name the different categories they consist of, and the ones we will focus on in the current thesis.

In Section 3.2, we further define Classification Learning and its sub-categories, and analyze the Instance-Based Learning techniques, specifically the k- Nearest Neighbour algorithm, which we apply later on our dataset.

In Section 3.3, we describe the concept of Clustering, naming the groups it is divided into, and dedicating the rest of the section to k- Means algorithm, as a Partitioning Clustering method, which will be further implemented.

In Section 3.4, we explain the role of distance in the previously defined Learning Methods, and analyze some important distance metrics.

Finally, in Section 3.5, we examine some Evaluation approaches and techniques utilized in Machine Learning and theoretically describe the ones which will be further used.

3.1 Data Mining and Machine Learning

[11] [3] The process of discovering patterns in data is widely defined as Data Mining. It is an automatic or semiautomatic procedure, applied on substantial data quantities, in order to discover desired kinds of combinations and obtain meaningful information. This pattern identification is extremely useful, as it allows us to make nontrivial predictions on new data. [2]

[8] When patterns, that are mined, are represented in terms of a structure, that can be examined, reasoned about and used to inform future decisions, the patterns are called structural, as they are considered to capture the decision structure in an explicit way. Machine Learning is defined as a collection of techniques for finding and describing those structural data patterns, a tool for helping to explain the data and make predictions from it.

Data take then the form of a set of examples or situations, generally characterized as instances, and the output of the process takes the form of predictions and sets of rules about new examples, under given circumstances. So, learning can be considered as having two separate definitions: the acquisition of knowledge and the ability to use it for further purposes. The thing to be learned is defined as concept, and the output produced by a learning scheme is the concept's definition.

In Data Mining applications, the learning procedures can be categorized in four different fields:

- **Classification Learning**, where the learning scheme is presented with a set of classified examples, from which it is expected to learn a way of classifying unseen examples
- **Association Learning**, where any association among features is sought, not just ones that predict a particular class value
- **Clustering**, where groups of examples that belong together are sought

- **Numeric Prediction**, where the outcome to be predicted is not a discrete class, but a numeric quantity

In the current thesis, two of the above methods will be thoroughly examined and applied on our data, Classification and Clustering, which are more specific definitions for the two larger categories of Machine Learning: Supervised and Unsupervised Learning.

3.2 Classification

The task of classification covers a wide range of human activity. In its broadest sense, the term involves any decision or forecast made on the basis of currently available information, and a classification procedure is defined as a formal method for repeatedly making such judgments in upcoming situations. In more specific terms, the problem concerns the construction of a procedure that will be applied to a continuing sequence of cases, in which each case must be assigned to a predefined class, depending on observed features or attributes. [1]

Classification is called Supervised Learning, because, in a sense, the scheme operates under supervision, by being provided with the actual outcome for each of the training examples. This outcome is called the class of the example.

Some of the most urgent problems arising in science, industry and commerce, demanding complex and often extensive data, can be regarded as classification or decision problems, such as the preliminary diagnosis of a patient's disease, whilst awaiting definitive test results, in order to select immediate treatment, or the assignment of individuals to credit status on the basis of their financial and personal information.

[5] Supervised learning is one of the tasks most frequently carried out by so-called Intelligent Systems. A large number of techniques have been developed, which are nominally divided in the following fields:

- **Logical and Symbolic techniques**, such as Decision Trees and Learning Rulesets
- **Perception-based techniques**, analyzed in Single Layer Perceptrons, Multilayered Perceptrons and Radial Basis Function (RBF) Networks
- **Statistics**, which include Naive Bayes Classifiers and Bayesian Networks
- **Instance - Based Learning**
- **Support Vector Machines**

In the next section we will focus on Instance- Based techniques, some of which were applied in the current thesis project.

3.2.1 Instance - Based Learning

[9] Instance - Based Learning algorithms are lazy-learning algorithms, as they delay the induction or generalization process, until classification is performed. Once a set of training instances has been memorized, on encountering a new instance, the memory is searched for the training instance that most strongly resembles the new one. In other words, the known instances are being stored and new instances, whose class is unknown, are being related to existing ones. Thus, all the real work is done when the time comes to classify a new instance, rather than when the training set is processed, and so the algorithms require less computation time during the training phase than eager learning algorithms, but more computation time during the classification process.

The above procedure is called the Nearest-Neighbor classification method. The absolute position of the instances within this space is not as significant, as the relative distance between them. Using

a suitable distance method, which ideally minimizes the distance between two similarly classified instances, while maximizing the distance between instances of different classes, the closest existing instance is used to assign the new one to the class. Sometimes, more than one nearest neighbor is used, and the majority class of the closest k - Neighbors (or the distanceweighted average, if the class is numeric) is assigned to the new instance. This is termed the k - Nearest Neighbor method. The selection of k strongly affects the performance of the k - NN algorithm and the result of the classification.

A pseudo-code example for the instance base learning methods is illustrated below. We consider the setup as following: X : *Training Data*, Y : *Class Labels of X*, x : *Unknown sample*

```

k Nearest Neighbours( $X, Y, x$ )
for  $i=1$  to  $m$  do
    Compute Distance  $d(X_i, x)$ 
end for
compute set  $I$  containing indices for the  $k$  smallest distances  $d(X_i, x)$ .
return majority label for  $Y_i$  where  $i \in I$ 

```

3.3 Clustering

Clustering techniques apply when the instances are to be divided into natural groups and there is no specified class to be predicted. These clusters presumably reflect a mechanism that causes some instances to bear a stronger resemblance to each other, than they do to the remaining ones. They are considered to be unknown and are inferred from the data, that is why Clustering is also known as Unsupervised Learning.

The groups that are identified may belong to one of the following different categories, based on the nature of the mechanisms that are thought to underlie in the particular clustering phenomenon:

- **Exclusive**, meaning any instance belongs to only one group
- **Overlapping**, as an instance may fall into several groups
- **Probabilistic**, because an instance may belong to each group with a certain probability
- **Hierarchical**, as a rough division of instances into groups at the top level and each group refined further, perhaps all the way down to individual instances.

However, because these mechanisms are rarely known, as the very existence of clusters is, after all, something that we're trying to discover, the characterisation of them is usually dictated by the clustering tools at our disposal.

A few real examples of the use of clustering involve dividing customers into homogeneous groups as a marketing focused procedure, a weather data collection and analysis for finding new insights into climatological and environmental trends, and bioinformatics' identification of groups of genes with similar patterns of expression, to determine which genes are responsible for specific hereditary diseases.

Due to the fact that the notion of cluster is not precisely defined, many clustering methods have been developed, each of which uses a different induction principle. They are divided into five main groups, as analyzed below:

- **Hierarchical Methods**, which construct the clusters by recursively partitioning the instances in either a top-down or bottom-up fashion. They can be sub-divided in Agglomerative hierarchical clustering and Divisive hierarchical clustering

- **Partitioning Methods**, which relocate instances by moving them from one cluster to another, starting from an initial partitioning. Types of partitioning methods involve Error Minimization Algorithms and Graph - Theoretic Clustering. The simplest algorithm, employing a squared error criterion is the K-means algorithm, which will be further analyzed as it was used on the current problem.
- **Density - Based Methods**, which assume that the points that belong to each cluster are drawn from a specific probability distribution
- **Model - Based Clustering Methods**, which attempt to optimize the fit between the given data and some mathematical models. The most frequently used methods in this category are Decision trees and Neural Networks
- **Grid - Based Methods**, which partition the space into a finite number of cells that form a grid structure, on which all of the operations for clustering are performed

3.3.1 Partitioning Methods-k-Means algorithm

Partitioning clustering methods, as mentioned, partition the data object set into clusters where every pair of object clusters is either distinct (hard clustering) or has some members in common (soft clustering).

The classic and most common technique is called k-Means. It can be also characterized as an Iterative Distance-Based Clustering method. Applying this algorithm, we first specify in advance the parameter k , which represents how many clusters are being sought. Then k points are chosen at random as cluster centers. We also define a maximum number of iterations for the algorithm to run, over which the process is terminated. All instances are assigned to their closest cluster center, based on the ordinary Euclidean distance metric. Next the centroid, or mean, of the instances in each cluster is calculated, and these centroids are considered to be new center values for their respective clusters. Finally, the whole process is repeated with the new cluster centers. Iteration continues until the same points are assigned to each cluster in consecutive rounds, at which stage the cluster centers have stabilized and will remain the same forever. If maximum number of iterations has been reached in the meantime, the algorithm is terminated.

A pseudocode of the above description is presented below, where the setup is the following:
 S : Instance set, k : Number of Clusters

```

k-Means( $S$ ,  $k$ )
  Initialize  $k$  cluster centers
  while termination condition is not satisfied do
    assign instances to the closest cluster center
    update cluster centers based on the assignment
  end while
  return clusters

```

This clustering method is simple and effective. It is easy to prove that choosing the cluster center to be the centroid minimizes the total squared distance from each of the cluster's points to its center. Once the iteration has stabilized, each point is assigned to its nearest cluster center, so the overall effect is to minimize the total squared distance from all points to their cluster centers. However, one should take into account that this minimum is a local one, there is no guarantee that it is the global minimum. To increase the chance of finding a global minimum, we often run the algorithm several times with different initial choices and choose the best final result.

3.4 Distances

[7] The notion of distance is the most important basis for Machine Learning, both Supervised and Unsupervised. In the first category, standard distances often do not lead to appropriate results, while in the second one, the calculation of means of objects of known groups is not always a valid method for the correct algorithms' application. By definition, the choice of the distance measure determines whether two objects naturally go together and, therefore, the right choice of the distance measure is one of the most decisive steps for the determination of learning properties. The distance should not only adequately represent the relevant scaling of the data, but also the study target, to obtain interpretable results. Some of the most widely used distance metrics, which were also implemented in the current thesis, are analyzed further below. [6]

3.4.1 Minkowski Metric

The Minkowski metric or L_q norm calculates the distance d between the two objects x and y by comparing the values of their n features. The Minkowski metric, as given in equation 9.1, can be applied to frequency, probability and binary values.

$$d(x, y) = L_q(x, y) = \sqrt[q]{\sum_{i=1}^n |x_i - y_i|^q} \quad (3.1)$$

The most important special case of the Minkowski metric is for $q=2$, the **Euclidean distance** or L_2 norm:

$$d(x, y) = L(x, y) = \sqrt{\sum_{i=1}^n |x_i - y_i|^2} \quad (3.2)$$

3.4.2 Cosine Distance

The cosine similarity (or Orchini similarity, angular similarity, normalized dot product) is a similarity on R^n , defined by

$$\cos(a) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \cdot \|\mathbf{y}\|} \quad (3.3)$$

where a is the angle between vectors x and y . In the binary case, it is called the Ochiai-Otsuka similarity. The cosine distance is defined as

$$d(x, y) = 1 - \cos a$$

3.4.3 Kullback-Leibler Divergence

The Kullback-Leibler divergence (KL) or relative entropy is a measure, calculated from information theory, which determines the inefficiency of assuming a model distribution, given the true distribution. It is generally used for x and y representing probability mass functions. The equation for its calculation is

$$d(x, y) = D(x \parallel y) = \sum_{i=1}^n x_i \star \log \frac{x_i}{y_i} \quad (3.4)$$

3.4.4 Kolmogorov-Smirnov Test

The Kolmogorov-Smirnov metric (or Kolmogorov metric, uniform metric) is a metric on probability space P , defined by

$$d(x, y) = \sup_{(x,y) \in R} |x - y| \quad (3.5)$$

considering again that x, y are different distribution functions. It is used in statistics as measure of goodness of fit.

3.4.5 Infinite Norm

The Infinite (or Uniform or Sup) Norm is the $L_{(\infty)}$ metric on the set $C_{[a,b]}$ of all real or complex continuous functions on a given segment $[a,b]$. For vectors x, y it is defined by

$$d(x, y) = \lim_{p \rightarrow \infty} \sqrt[p]{\sum_{i=1}^n |x_i - y_i|^p} \quad (3.6)$$

3.5 Evaluation

Evaluation is of crucial value in Data Mining, as the assessment of the application of our methods and the results of our techniques are the key to further conclusions, optimizations and progress. In this section, the concept is analyzed for both Classification and Clustering.

3.5.1 Classification Evaluation

The success rate of the classification learning is usually judged on the test data, for which the true classes are known, and is evaluated using different metrics, giving an objective measure of how well the concept has been learned by the data. Here, two different evaluation methods are being analyzed, Accuracy and Cohen's Kappa, which will be applied in Chapter 5 in the simulations' results.

Another major issue on the classifier's evaluation is its speed. A classifier that is 90% accurate may be preferred over one that is 95% accurate if it is 100 times faster in testing (and such differences in time-scales are not uncommon in neural networks, for example). This parameter is also examined in the experiments presented later in the current thesis.

Accuracy

The reliability of the classification result is represented by the proportion of correct classifications.

Usually, it is the accuracy on the unseen data, when the true classification is unknown, that is of practical importance. The generally accepted method for estimating this is to use the given data, in which we assume that all class memberships are known. Firstly, we use a substantial proportion (the training set) of the given data to train the procedure. This rule is then tested on the remaining data (the test set), and the results compared with the known classifications. The percentage of correctly predicted data in the test set is an unbiased estimate of the accuracy of the rule provided that the training set is randomly sampled from the given data.

There is a slight loss of efficiency here, as we do not use the full sample to train the decision rule, but with very large datasets this is not a major problem.

	1	2	Total
1	p ₁₁	p ₁₂	p _{1r}
2	p ₂₁	p ₂₂	p _{2r}
Total	p _{1c}	p _{2c}	1

Πίνακας 3.1: Confusion Matrix after Classification

Cohen's Kappa (Kappa coefficient)

[12] The Kappa statistic is a metric that compares an Observed Accuracy with an Expected Accuracy, and is used to evaluate a single classifier or classifiers amongst themselves. It takes into account random chance (agreement with a random classifier), which generally means it is less misleading than simply using accuracy as a metric. Computation of Observed Accuracy and Expected Accuracy is integral to comprehension of the kappa statistic, and is most easily illustrated through use of a confusion matrix.

For a two dimensional confusion matrix, as the previous one, the Kappa metric is equal to

$$k = \frac{p_0 - p_e}{1 - p_e} \quad (3.7)$$

where the Observed Accuracy is

$$p_0 = p_{11} + p_{22}$$

and the Expected Accuracy is

$$p_e = p_{1c}p_{1r} + p_{2c}p_{2r}$$

Kappa is always less than or equal to 1. A value of 1 implies perfect agreement and values less than 1 simply less than perfect agreement. In rare situations, Kappa can be negative. This is a sign that the two observers agreed less than would be expected just by chance.

The values and according interpretations of Kappa are summarized below:

- Poor agreement = Less than 0,20
- Fair agreement = 0,20 to 0,40
- Moderate agreement = 0,40 to 0,60
- Good agreement = 0,60 to 0,80
- Very good agreement = 0,80 to 1,00

3.5.2 Clustering Evaluation

Clustering evaluation demands an independent and reliable metric for the assessment and comparison of clustering experiments and results. In theory, the clustering researcher has acquired an intuition for the clustering evaluation, but in practice the mass of data, on the one hand, and the subtle details of data representation and clustering algorithms, on the other hand, make this random judgement impossible. An intuitive, introspective evaluation can only be plausible for small sets of objects, whilst large-scale experiments require an objective method. There is no absolute scheme for the desired assessment, but a variety of evaluation measures from diverse areas such as theoretical statistics, machine vision and web-page clustering that can be applied.

The theoretical definition of various clustering evaluation metrics, which were used during the thesis's experiments, is provided.

V-measure

[4] V-measure is an entropy-based metric, which measures how successfully the criteria of homogeneity and completeness have been satisfied. It is computed as the harmonic mean of distinct homogeneity and completeness scores.

A clustering result satisfies homogeneity if all of its clusters contain only data points, which are members of a single class. A clustering result satisfies completeness if all the data points, that are members of a given class, are elements of the same cluster. The homogeneity and completeness of a clustering solution run roughly in opposition: Increasing the homogeneity of a clustering solution often results in decreasing its completeness. More specifically:

- **Homogeneity:** In order to satisfy our homogeneity criteria, clustering must assign only those datapoints that are members of a single class to a single cluster. That is, the class distribution within each cluster should be skewed to a single class, having zero entropy. We determine how close a given clustering is to this ideal by examining the conditional entropy of the class distribution given, the proposed clustering.
- **Completeness:** Completeness is symmetrical to homogeneity. In order to satisfy the completeness criterion, a clustering must assign all of those datapoints that are members of a single class to a single cluster. To evaluate completeness, we examine the distribution of cluster assignments within each class. In a perfectly complete clustering solution, each of these distributions will be completely skewed to a single cluster.

So, V-Measure is given by the following equation

$$V_b = \frac{(1 + b) \star h \star c}{(b \star h) \star c}$$

The calculations of homogeneity, completeness and V-measure are completely independent of the number of classes, the number of clusters, the size of the data set and the clustering algorithm used. Thus, these measures can be applied to and compared across any clustering solution, regardless of the number of data points, the number of classes or the number of clusters

Rand - Index

The Rand - Index is a simple criterion used to compare an induced clustering structure (C_1) with a given clustering structure (C_2). Let

- **a** be the number of pairs of instances that are assigned to the same cluster in C_1 and in the same cluster in C_2
- **b** be the number of pairs of instances that are in the same cluster in C_1 , but not in the same cluster in C_2
- **c** be the number of pairs of instances that are in the same cluster in C_2 , but not in the same cluster in C_1
- **d** be the number of pairs of instances that are assigned to different clusters in C_1 and C_2 .

The quantities a and d can be interpreted as agreements, and b and c as disagreements. The Rand Index is defined as:

$$RAND = \frac{a + d}{a + b + c + d}$$

The Rand Index lies between 0 and 1. When the two partitions agree perfectly, Rand Index is 1.

Κεφάλαιο 4

Implementation

In this Chapter we present the implementation of all the concepts analyzed in Chapter 3 for our experiments, meaning all algorithms and their parameterization.

In Section 4.1, we present the two algorithms used for Classification, Nearest Centroid Classifier and k-Nearest Neighbours Classifier, and the important variables for their execution.

In Section 4.2, we focus on the Clustering procedure, giving a detailed description of the k-Means algorithm used for this purpose.

In Section 4.3, we describe the different distance calculation methods, for all simulations, and present all ways they were implemented and evaluated in both Classification and Clustering.

Finally, in Section 4.4, we explain in detail the way the results were evaluated and the metrics used for this purpose.

4.1 Classification

The objective of our approach was to classify the production chunk objects according to some of their attributes. Specifically, we wanted to categorize them based on two separate parameters, Product Cleaning Group, and Product Group, which are both attributes of each object. The Product Cleaning Group has 5 instances, and the Product Group has 12.

We also experimented with the classification algorithms for the Product Code attribute of each object, but the large number of separate production codes and therefore the small number of elements assigned to each class, resulted in classifier's overfitting. It was not possible to correctly train and test our data in a meaningful way, and the outcome of the process provided no more information than the existing about products' classes.

For the classification, two different algorithms were used, as analyzed previously: a custom Nearest Centroid algorithm and the known well documented k-Nearest Neighbours Algorithm.

Both algorithms are parameterized as following:

- **Classifying attribute:** both algorithms were implemented for two attributes, Product Cleaning Group and Product Group,
- **Split percentage:** We divided the initial product data into two groups, one training set, on which we calculated the separate classes' centers and applied the algorithms to produce our result, and one testing set, which we classified based on the training procedure's derived classes. The training and testing split percentages selected were 80% -20%, 65% -35% and 50% -50% respectively for both algorithms.
- **Distances:** We used the defined distances functions to calculate the distance between the transition matrices of each chunk object from either the class' centers or the nearest chunks, depending on the implementation.

4.1.1 Nearest Centroid Classifier

For our list of chunks with “Production” attribute only, we followed the steps mentioned below, to run the experiments of classification: Selecting the attribute on which we would classify the objects, which was either Product Cleaning Group or Product Group, we initially split the data based on the three training and test dataset percentages, as mentioned in parameterization before. We use the training data to calculate the centers of the algorithm. Focusing on the Transition Matrices of each object in this set, and examining the value of the classification attribute of the object, we found the average transition matrix of all the objects belonging to each instance. This matrix is considered to be the center of the class. Then, for each element in this class, from the initial test set, we calculated its distance from the different centers. Finally we assigned the test element to the class of which the center was nearest.

Depending on the experiment, different methods of distance algorithms were used.

Due to uneven sized classes, each experiment was executed for 30 iterations, meaning for 30 different sets of data, of the same split percentage. This enabled us to validate our results and to be sure of the integrity of our data.

For a pseudocode implementation of the above, please see [A.1.1](#).

4.1.2 K - Nearest Neighbors Classifier

K - Nearest Neighbors is a well-known classification algorithm. Although many implementations exist, we decided to implement our own, in order to have total control over the variables and its run-time execution. The steps of the execution are the following and the methodology is relatively close to the one of Nearest Centroid Classifier.

The first step was to choose the attribute for which we wanted to run the classifier, this can be either Product Cleaning Group or Product Group. Following, we split the data-set into training and test, for three different percentages, as in Nearest Centroid Classifier. For each iteration, 30 in total, we split the data-set into random sets, of specific percentage, in order to compare the integrity of the results. Having selected the distance method, we calculate the distance of each chunk object in the test set with all of the chunk objects in the training set. After that, we sorted the values for each test chunk in ascending order. The principal of K-NN algorithm is that each element of the test set is classified to the class that the majority of its k-Nearest Neighbors belongs to. Once we had created a sorted list of all the neighbors, it was easy to evaluate the classification for each value of k, within the desired limits.

At this point we have classified all the chunk objects of the test set for every k in our range, for a specific distance method. In Chapter 5 you can see the impact of k in the evaluation of k-NN. We repeated the process for all the distance methods that concerned us.

For a pseudocode implementation of the above, please see [A.1.2](#).

4.2 Clustering

The purpose of our clustering implementation is to assign our product chunks data into clusters based on some of their attributes. In this direction, we implemented a custom k - Means algorithm, analyzed as following:

4.2.1 K - Means Implementation

- **Clustering parameter k:** The number of clusters that the data is expected to be assigned to is initialized as the number of values in each of the attributes we performed the clustering on. Thus, for Product Cleaning Group, we have k=5 and for Product Group we have k=12.
- **Initial Centroids:** The centroid is represented as a transition matrix of a chunk class object. The experiments were held with a variety of initial centroids, both random and specific, so as to test the performance and accuracy of the algorithm. In the first scenario, each centroid was picked to

be a random transition table from one of the products of the input data, whilst in the second one, we initialized the centroids choosing objects with the clustering attribute in specific ranges or even of a certain value. More specifically, centroids were chosen to be the matrices of products either all in the same Production Cleaning Group (or Product Group, according to the clustering attribute), or all in different Production Cleaning Group/ Product Group, or, finally, all in random Production Cleaning Group/ Product Group. Each of these 3 separate experiments was executed 100 times to eliminate variation in our selections.

- **Average:** In each repetition of the algorithm, a new centroid of each cluster was calculated, based on its currently assigned members. In our implementation, this is represented by the average transition matrix of all transition matrices of the product chunks in this cluster.

For a pseudocode implementation of the above, please see [A.2.1](#).

4.3 Distances

The processes defined before, classification and clustering, and the algorithms implemented in both, require a distance calculation between the different elements. In our problem, considering the procedures are based on production messages, distance is calculated as the distance between the two dimensional transition matrices of two separate product chunks. Different distance algorithms were parameterized and applied in the experiments, so as to examine the performance and accuracy of each implementation. We used the following methods:

- Euclidean distance
- Cosine Distance
- Kullback - Leibler Divergence
- Kolmogorov - Smirnov Test
- Infinity Norm

Most of them, by definition, apply only on one-dimensional matrices (vectors), thus an important step was the correct transformation of the two dimensional matrices into vectors, in order to calculate the distances. Our approaches can be divided in three categories:

- **Average rows:** the distance algorithm is applied between the same rows of the matrices and the result was the average of all scores calculated.
- **Vector:** each row is appended to the first one, thus creating a $1 \times N$ vector, on which the distance method is applied.
- **Diagonal:** for $j > i$, the average of elements $[i, j]$ and $[j, i]$ is calculated and considered as one, thus creating an upper triangular matrix. Then, with the previous Flatten method, each row is appended to the first, for the distance method to be applied on.

The exact implementation of the previous methods can be found in the subsection [4.3.1](#) .
The total list of the implemented distance methods is the following:

- **Euclidean distance**
 - Euclidean distance between each corresponding cell
 - Average of the Euclidean distances between corresponding rows

- Average of the Euclidean distances between corresponding columns
- **Cosine distance**
 - Average of the cosine distance between corresponding rows
 - Cosine distance of the created vectors using the Vector transformation
 - Cosine distance of the created vectors using the Diagonal transformation
- **Kullback–Leibler Divergence**
 - Average of the KL Divergence between corresponding rows
 - KL Divergence of the created vectors using the Vector transformation
 - KL Divergence of the created vectors using the Diagonal transformation
- **Kolmogorov–Smirnov test**
 - Average of the KS test score between corresponding rows
 - KS test score of the created vectors using the Vector transformation
 - KS test score of the created vectors using the Diagonal transformation
- **Infinity Norm**
 - Infinity Norm of the difference of the two matrices

For the implementation of the above mentioned distance algorithms, the Python SciKit-Learn Library was used. This library provides many useful distance calculation functions but most of them are implemented for vectors.

4.3.1 Implementation of Distance Methods

Two of the functions that convert a $2D$ matrix into a *Vector*, can be found in a pseudocode implementation at [A.3.1](#).

The different implementations of the distance algorithms, as described in the beginning of this Section, can be found at [A.3.2](#).

4.3.2 Evaluation of Distance Methods

We were interested in evaluating the performance of all distance methods under different variation of the variables in question, for both classification and clustering cases. The reason we chose to run different evaluation tests on classification and clustering is that we wanted to understand the impact of the distance method in each case.

- Distance Methods Evaluation for Classification

We took the decision to only run the following experiments with the Nearest Centroid Classifier for the 80% train split. The classifier run for both Product Cleaning Group and Product Group for 30 iterations, as described in previous experiments, to ensure integrity of the results. For each iteration, we randomly split the data-set in 80%-20%, training-testing respectively, and calculated the centers of each instance, either 5 or 12 classes. After selecting a distance method, we assigned the chunk objects to their nearest center, based on the selected distance algorithm. A pseudocode implementation of the above algorithm can be found at [A.3.3](#). The results of this implementation can be found in Section [5.1.2](#).

- Distance Methods Evaluation for Clustering

In the clustering case, the only tested algorithm is K-Means, thus we evaluated all distance methods on K-Means algorithms. A pseudocode implementation of the above algorithm can be found at [A.3.4](#).

4.4 Outcome Evaluation

All of the above can only produce a usable result if they can be evaluated. In early experiments we used the Confusion Matrix of the outcome, which enabled us to have a more direct overview of the results. As the experiments evolved, we had to evaluate each outcome with a single score. The need for evaluation scores became more clear and we had to decide which scores to choose between the many evaluation methods that exist, for both classification and clustering. We used two different evaluation methods for each case. The main difference on evaluating Classification and Clustering results is that in the first case one knows the correct answer, because all data are labeled from the beginning. In the case of Clustering, however, one does not know the correct value of the data and has to evaluate the outcome based on different factors.

4.4.1 Implementation of Classification Evaluations

For the evaluation of the Classification Results we used the following methods:

- **Accuracy**

Accuracy is defined as the ration between the correctly classified items and the total number of classified items. In our case, the correctly classified chunk objects were the ones that the cluster's name was the same as the initial label of the respective attribute. The attributes of the classification were the Production Cleaning Group and the Production Group, as analyzed before.

$$\text{Accuracy} = \frac{\text{Number of Correctly Classified Elements}}{\text{Total Number of Classified Elements}}$$

- **Kappa**

The Kappa algorithm is part of the SciKit-Learn Laboratory Library and specifically, of the Metrics Module. SciKit is a well known Library for Python and has proved to be very helpful in our thesis project.

Source: [SciKit-Learn Laboratory -> Metrics -> Kappa](#)

4.4.2 Implementation of Clustering Evaluations

For the evaluation of the Clustering Results we used the following methods. Although in other clustering evaluation scenario one might not had the labels, in this case we had all our data labelled. Both algorithms require ground truth class labels, which are used as reference in order to evaluate the predicted.

- **V Measure**

Source: [SciKit-Learn -> Metrics -> V Measure Score](#)

- **Rand Index**

Source: [SciKit-Learn -> Metrics -> Adjusted Rand Score](#)

Κεφάλαιο 5

Results

In this chapter, we present the results from the application of classification and clustering algorithms on our data set. All algorithms, as described in Chapters 2 and 4, were executed for multiple parameters, thus providing different results and conclusions for each settings' combination.

In Section 5.1, we present the classification results, divided into three categories. At first, we compare both Nearest Centroid classifier and k-Nearest Neighbors classifier, for different training and test data splits, keeping the distance parameter fixed. Next on, we focus exclusively on the Nearest Centroid algorithm, for a specific data split and evaluate the classifier's performance for all different chunk distances. Finally, we examine the application of k-Nearest Neighbors classifier for different values of k and assess the outcome.

In Section 5.2, the clustering results are presented and evaluated. We used k-Means algorithm for clustering the chunk objects product data into different clusters, based again on two attributes, Product Cleaning Group and Product Group. The algorithm's input data are the chunk objects. The centroids and distances vary in each execution, so that we can evaluate the clustering result for different problem parameterizations, using two different metrics, V-Measure and Rand Index.

5.1 Classification

The classification was performed for two distinct attributes, Product Cleaning Group and Product Group.

The algorithms tested were:

- Nearest Centroid Classifier
- k-Nearest Neighbors Classifier

The distance algorithm was a variable in all cases.

For each run of every algorithm, a confusion matrix was produced. Its rows reflect the classifier's results, meaning the number of products per label that were classified as being in the column's label instance after running the algorithm. Its columns represent the ground truth, the actual values of each instance to be classified.

The sum of each column is the number of products classified as instances of this label, while the sum of each row represents the true number of this label's products in our initial sample. The sum of the diagonal of the matrix is the number of the correctly classified products.

The final confusion matrix of each run is further used to calculate the performance for each of the parameters considered, as mentioned above.

All data used were split into a training and a testing sample, which varies, so as to better examine the accuracy of the classifiers predicting the correct class for the test set, regarding the initial data they were trained on.

5.1.1 Baseline

A baseline result was required, in order to compare the classification results to, and evaluate the outcome of each method.

We used the ZeroR algorithm, which selects the class that has the most observations, and uses that class as the result for all predictions.

The baseline score for the input chunks was calculated by finding the percentage of each class's items in total given number of products and selecting the class with the bigger one as our accuracy metric.

The ZeroR Accuracy for the two attributes can be found below:

- ZeroR Product Cleaning Group Accuracy: *0.520*
- ZeroR Product Group Accuracy: *0.377*

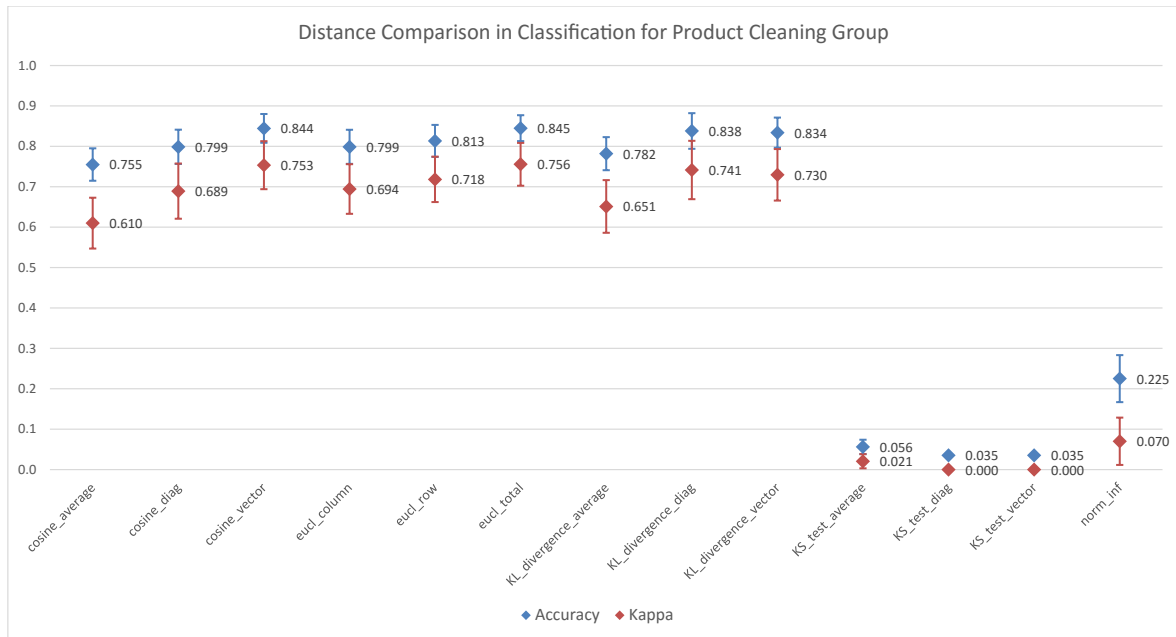
5.1.2 Distance Method Evaluation and Selection

Setup:

- *Algorithm:* Nearest Centroid Classifier
- *Attributes:*
 - Product Cleaning Group
 - Product Group
- *Split:* 80% training set, 20% test set
- *Distances:*
 - Euclidean Total
 - Euclidean Row
 - Euclidean Column
 - Cosine Average
 - Cosine Vector
 - Cosine Diagonal
 - KL - Divergence Average
 - KL - Divergence Vector
 - KL - Divergence Diagonal
 - KS - Test Average
 - KS - Test Vector
 - KS - Test Diagonal
 - Infinity Norm

In this experiment we assess the different Distance Methods, considering a fixed split of 80% train data and 20% test data, and using the Nearest Centroid classifier.

The results are visible in Figure [11.1](#) for Product Cleaning Group.



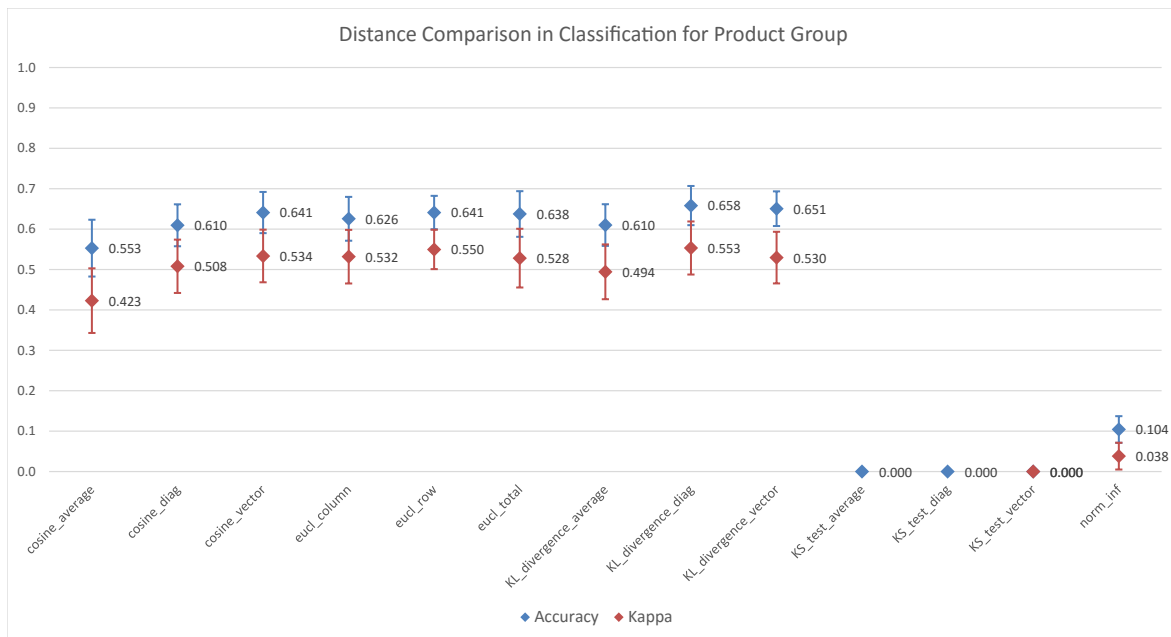
Σχήμα 5.1: Distance Comparison for Product Cleaning Group Classification

The best Accuracy and Kappa Coefficient scores appear when Euclidean Total distance is applied, and are followed by Cosine Vector Distance and KL - Divergence- Diagonal Scores, which are also very close. All other Euclidean, Cosine and KL metrics are in quite proximity with the above maximum scores.

Scores of KS - Test metrics and Infinity Norm metric appear to be significantly low, implying that they are inappropriate for distance calculation in the current problem.

The standard Deviation in Accuracy has a maximum value of 0,058, while in Kappa Coefficient it maximizes at 0,068, thus is considered negligible for our conclusions in both metrics.

For Product Group, results are displayed in Figure 11.2, with the KL - Divergence Diagonal and KL - Divergence Vector distances to be having the best Accuracy and Kappa Coefficient scores. Next in line follow the KL - Divergence Average and the Cosine Distance Average, with very close highest scores each.



Σχήμα 5.2: Distance Comparison for Product Group Classification

As in the previous experiment, the lowest scores appear in KS - tests and Uniform Norm results, which are concluded to be non- fitting for the current classification procedure.

Standard Deviation of both metrics has its maximum value at 0,056 and 0,799 respectively, so it is again insignificant to our results' accuracy.

5.1.3 Nearest Centroid and k-Nearest Neighbors Algorithms Comparison

Setup:

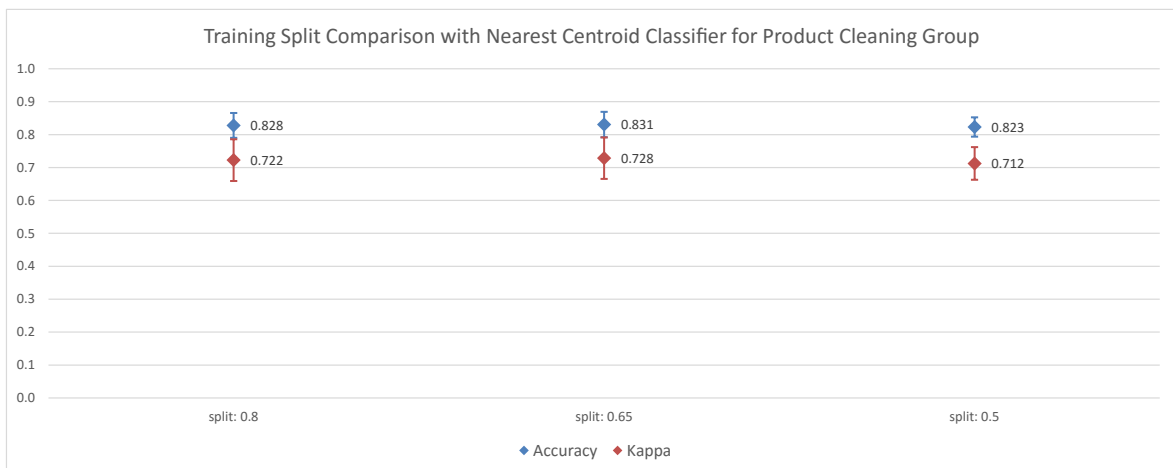
- *Algorithms:*
 - Nearest Centroid Classifier
 - k-Nearest Neighbors Classifier
- *Attributes:*
 - Product Cleaning Group
 - Product Group
- *Splits:*
 - 80% - 20%
 - 65% - 35%
 - 50% - 50%
- *Distance: Average of*
 - Euclidean Total
 - Cosine Vector
 - KL - Divergence Diagonal

In this section we compare the outcome of the application of the two algorithms, Nearest Centroid and k-Nearest Neighbours, on our data, focusing on the following parameters:

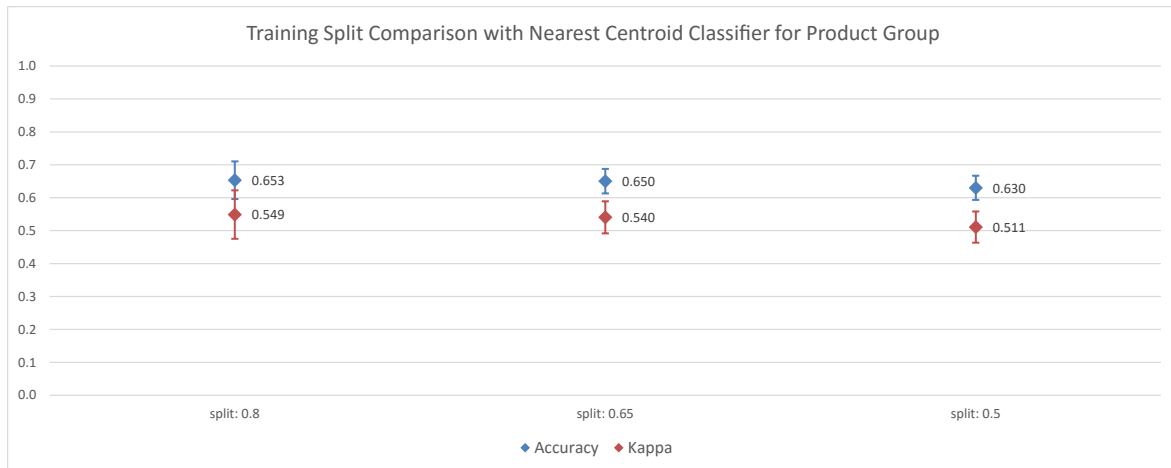
- Execution time
- Percentage of training- test data split
- Accuracy and Kappa Coefficient of classification result

The distance between the transition matrices of each chunk class object was kept fixed in this part of the process and was calculated as the average of the three different distances' results that appeared to have the maximum scores in the previous application : Euclidean Total, Cosine Distance- Vector and KL - Divergence Diagonal.

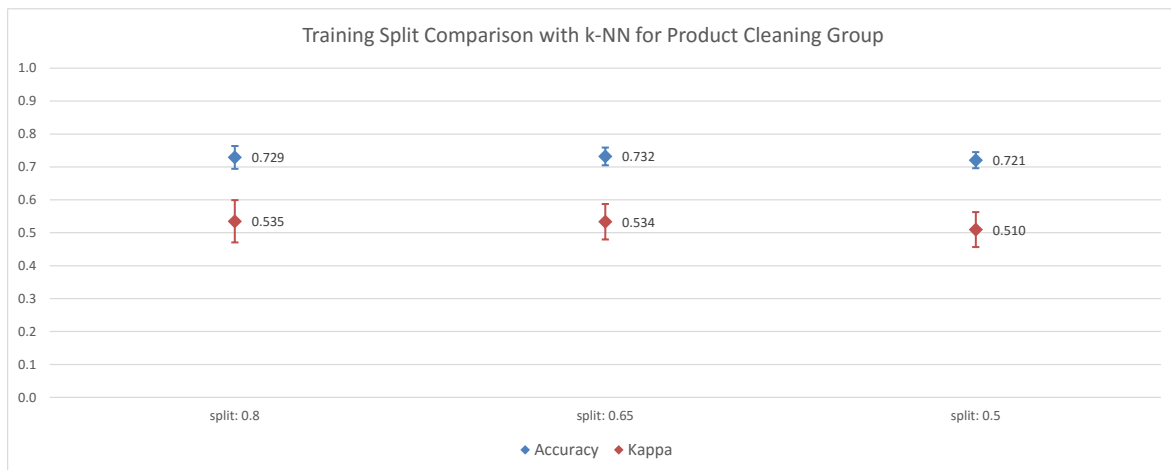
The results of the Nearest Centroid Algorithm's application are displayed in Figures 11.3 and 11.4, while the K -NN results in Figures 11.5 and 11.6.



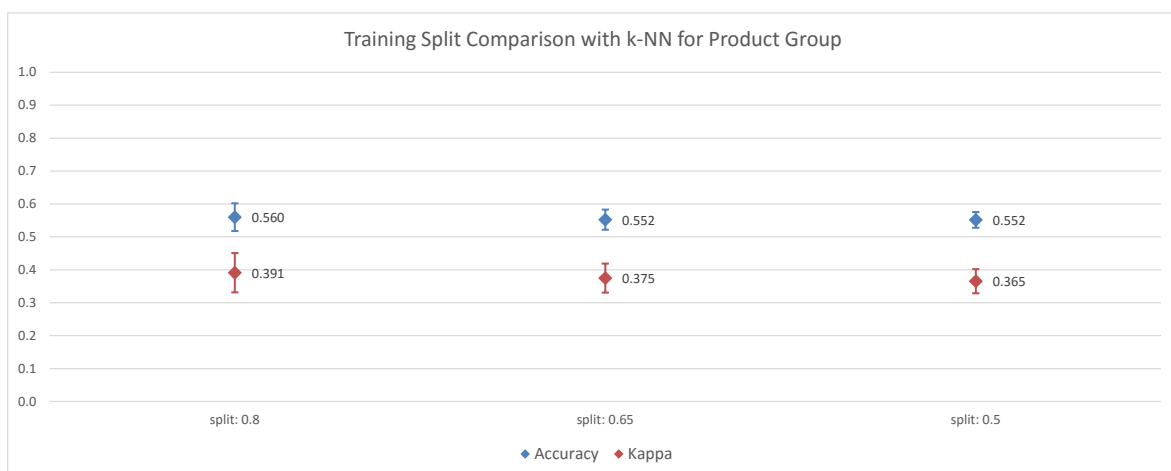
Σχήμα 5.3: Nearest Centroid Classifier for Product Cleaning Group



Σχήμα 5.4: Nearest Centroid Classifier for Product Group



Σχήμα 5.5: k-Nearest Neighbors Classifier for Product Cleaning Group



Σχήμα 5.6: k-Nearest Neighbors Classifier for Product Cleaning Group

It is obvious that, despite the three splits, the results of the classification are quite close, in both algorithms' charts. Regarding the Product Cleaning Group attribute, the Nearest Centroid algorithm gives maximum scores in the second split, with Accuracy in 83,1% and Kappa Coefficient in 72,8%,

while the same scores for the current attribute in k-Nearest Neighbors are 73,2% and 53,4%.

The conclusions are similar when examining the Product Group attribute charts. Here, differences are also very subtle between the distinct splits, with the first one to achieve a slightly bigger score than the others. For the Nearest Centroid algorithm, that maximum is a 65,3% in Accuracy and a 54,9% in Kappa Coefficient, whilst for the k -NN the respective scores are 56% and 39,1%.

From these results, it is obvious that the Nearest Centroid algorithm is superior to the k -NN, for the current problem and regarding the above parameters.

A strong argument for this is also each algorithm's runtime, as the first one required 20,15 seconds, while the second one run in 161,46 seconds, time amount at least 8 times greater than the other classifier's, which implies that the algorithm is not suitable for the examined classification.

5.1.4 k-Nearest Neighbors Algorithm for Different k

Setup:

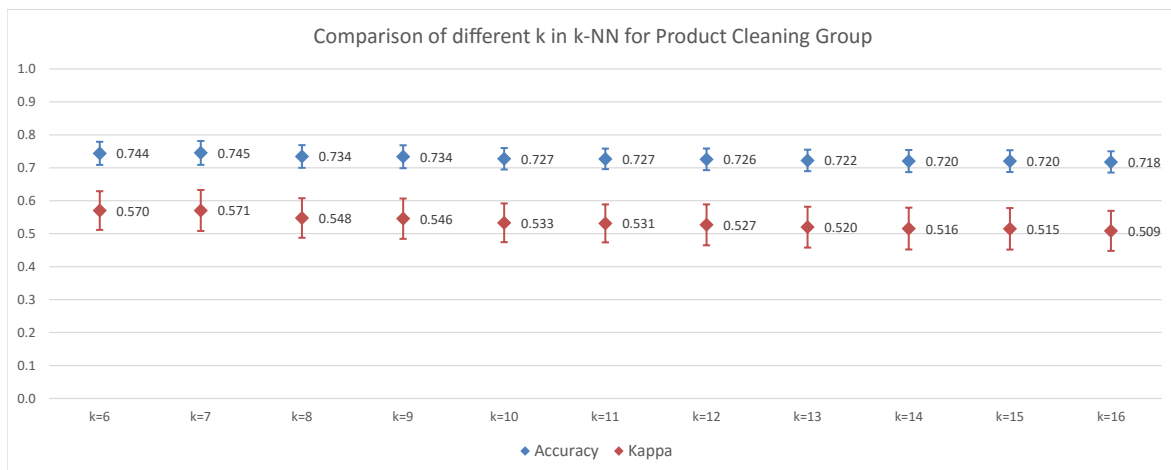
- *Algorithms:* k-Nearest Neighbors Classifier
- *Attributes:*
 - Product Cleaning Group
 - Product Group
- *Split:* 80% - 20%
- *Distance:* Average of
 - Euclidean Total
 - Cosine vector
 - KL - Divergence diagonal

The k -NN algorithm was applied for the above setup, so as to compare the effect of the selection of k in the classifier’s general accuracy and performance.

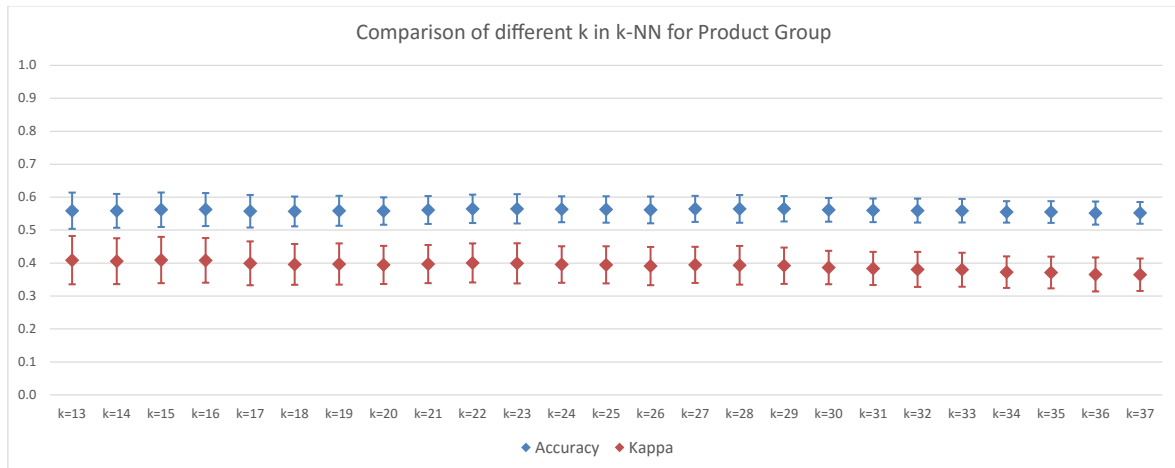
The algorithm was implemented for the two usually examined attributes, Product Cleaning Group and Product Group, and for a fixed split of 80% - 20% of training and testing data.

The distance was kept fixed as well, and equal to the average of the distances which had the maximum scores in the first experiment, Euclidean Total, Cosine Distance Vector and KL - Divergence Diagonal.

Results are presented in Figure 11.7 for Product Cleaning Group and Figure 11.8 for Product Group.



Σχήμα 5.7: K - NN for Different k in Product Cleaning Group



Σχήμα 5.8: K - NN for Different k in Product Group

Observing the charts, it is obvious that the selection of k has little impact on the Accuracy and Kappa Coefficient scores of the classification process, even for very different and distant values of k.

In Product Cleaning Group classification, the maximum score is achieved with k=7 for both Accuracy (74,5%) and Kappa Coefficient (57,1%) metrics, whilst classifying on Product Group attribute gives best scores for k=15, k=16, k=23, k=28 and k=29.

5.2 Clustering

5.2.1 Baseline

The baseline metric for the clustering procedure is calculated from the initial setup of the k - Means for each experiment, which provides us with two initial scores. The output of each run is compared to these results, which are depicted on all the following diagrams, so as to assess the performance of the experiments performed.

5.2.2 Distance Algorithms

Setup:

- *Algorithms:*
 - k - Means
 - Baseline
- *Attributes:*
 - Product Cleaning Group
 - Product Group
- *Initial Centroid Sets Type:*
 - All centroids of each set belonged to different clusters (Alldiff)
 - * Average of 20 sets
 - All centroids of each set belonged to the same cluster (Allsame)
 - * Average of 20 sets
- *Distances:*
 - Euclidean Total
 - Euclidean Row

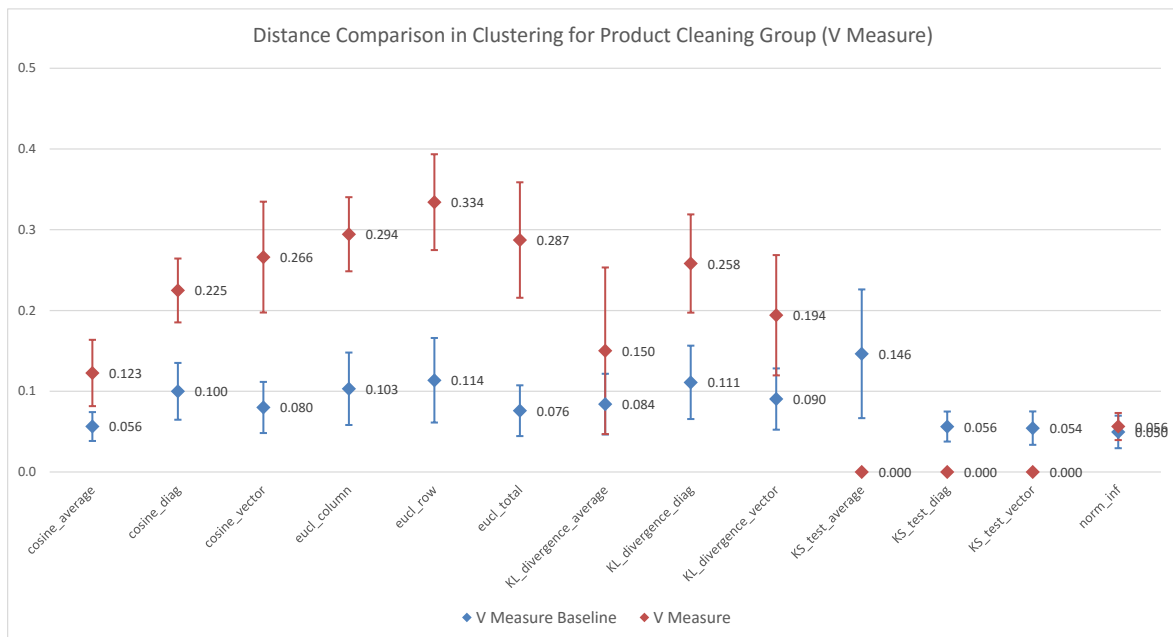
- Euclidean Column
- Cosine Average
- Cosine Vector
- Cosine Diagonal
- KL - Divergence Average
- KL - Divergence Vector
- KL - Divergence Diagonal
- KS - Test Average
- KS- Test Vector
- KS- Test Diagonal
- Infinity Norm

• *Evaluation Methods:*

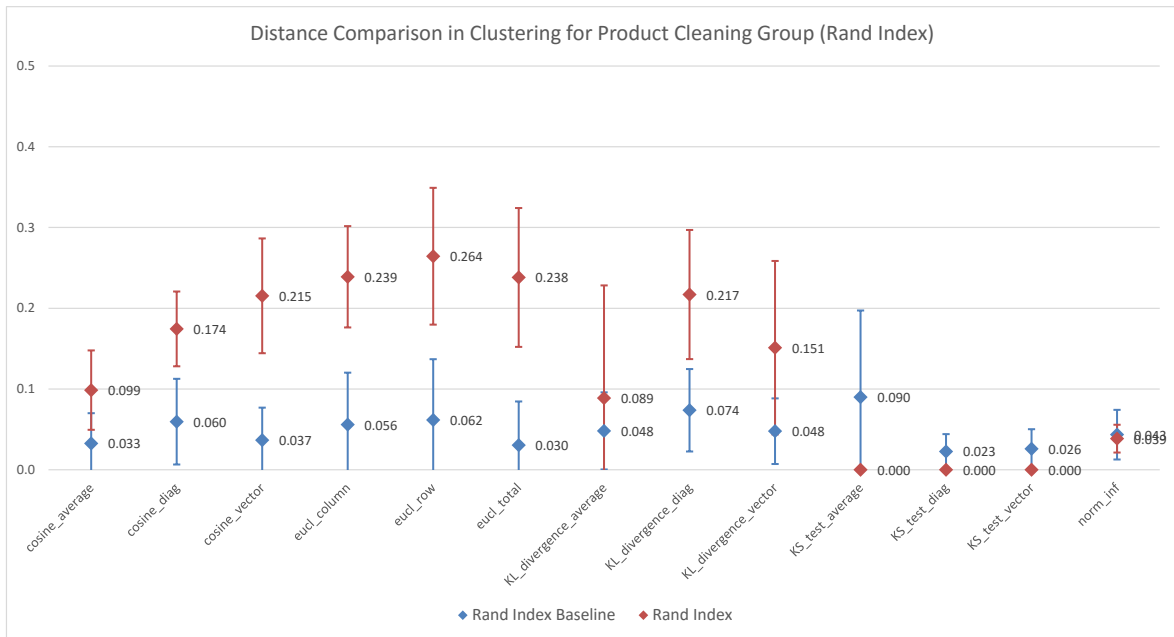
- V-Measure
- Rand Index

In Figures 11.9 and 11.10 the k - means algorithm was executed for 20 different sets of initial centroids, which were selected so that all of them, in each set, belonged to either a different or the same cluster, which is one of the 5 different Product Cleaning Groups. It is observed that the Euclidean distance methods perform the best, although the value variation is relatively big. The maximum score is achieved for the Euclidean Row distance metric and is a 33,4% V-Measure and a 26,4% Rand Index.

The low scoring percentages in all distances though, indicate that the clustering in general does not perform very well, yet its scores remain significantly above the baseline results.

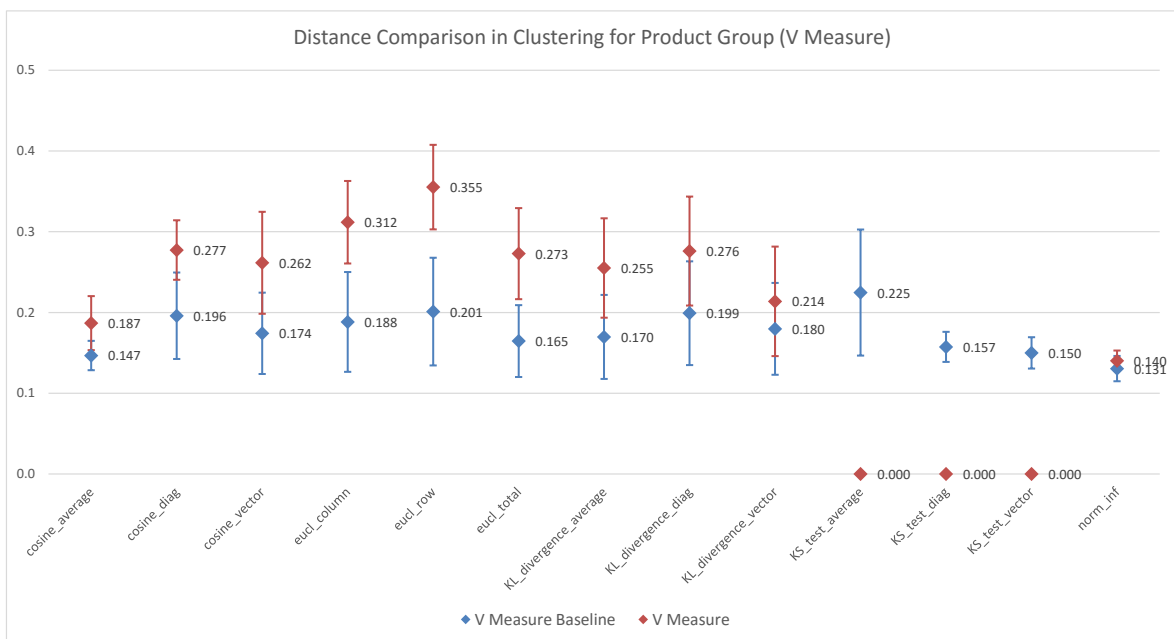


Σχήμα 5.9: k-means Clustering in Product Cleaning Group (V-Measure)

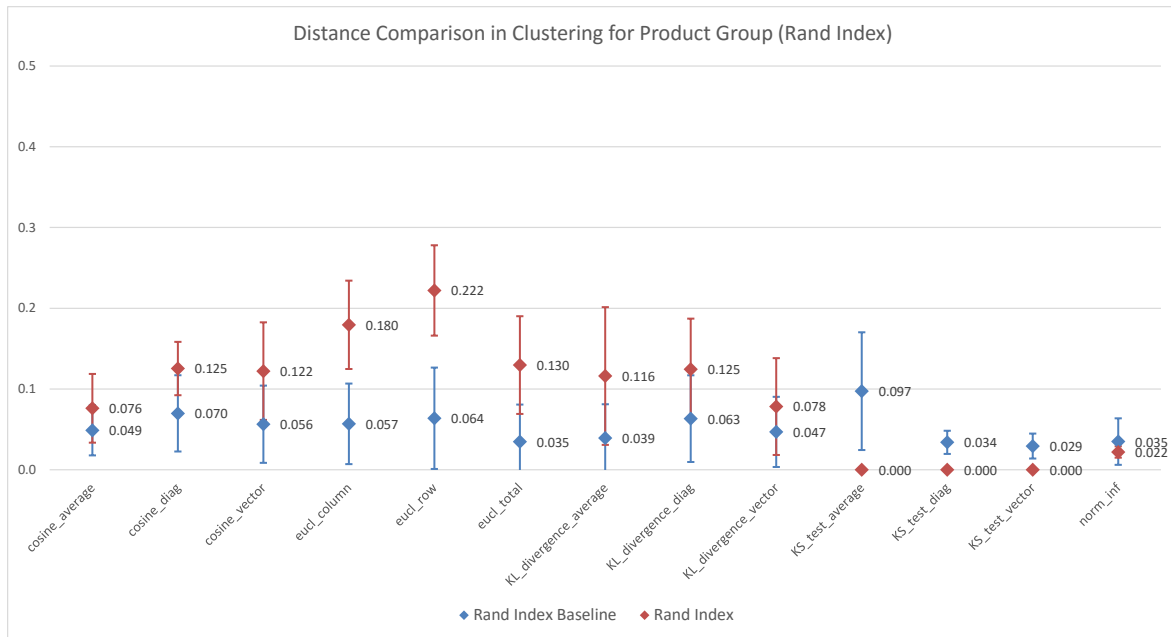


Σχήμα 5.10: k-means Clustering in Product Cleaning Group (Rand Index)

Figures 11.11 and 11.12 is of the same logic as the previous ones, but here the clustering was performed in the Product Group level. This means that we have 12 potential clusters instead of 5 we had in the previous two experiments.



Σχήμα 5.11: k-means Clustering in Product Group (V-Measure)



Σχήμα 5.12: k-means Clustering in Product Group (Rand Index)

Euclidean Distance Metrics have the dominant scores here again, with maximum accuracy for Euclidean Row, which scores a V-Measure of 35,5% and a Rand Index of 22,2%. KL - Divergence scores appear quite high as well.

In all cases, lowest scores come from the KS - Test and Infinity Norm metrics, and sometimes do not even exceed the baseline results.

5.2.3 Centroids

Setup:

- *Algorithms:*

- K - Means
- Baseline

- *Attributes:*

- Product Cleaning Group
- Product Group

- *Initial Centroid Sets Type:*

- All centroids of each set belonged to different clusters (Alldiff)
 - * Average of 100 sets
- All centroids of each set belonged to the same cluster (Allsame)
 - * Average of 100 sets
- All centroids of each set belonged to a random cluster (Allrand)
 - * Average of 100 sets

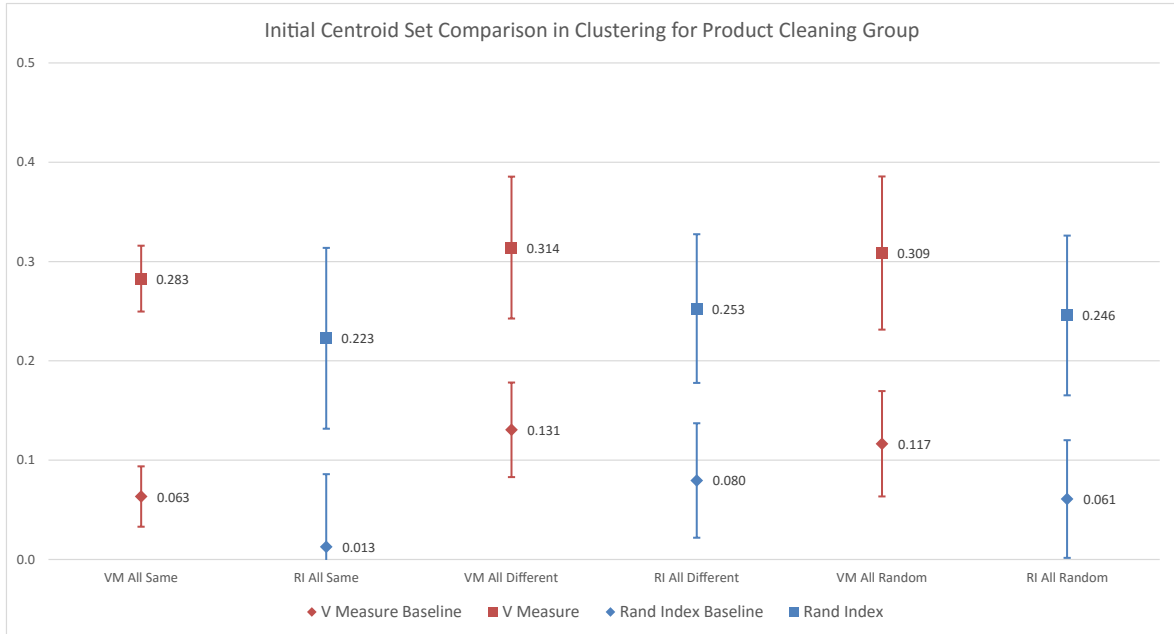
- *Distance: Average Of*

- Euclidean Total
- Euclidean Row
- Euclidean Column

- *Evaluation Methods:*

- V-Measure
- Rand Index

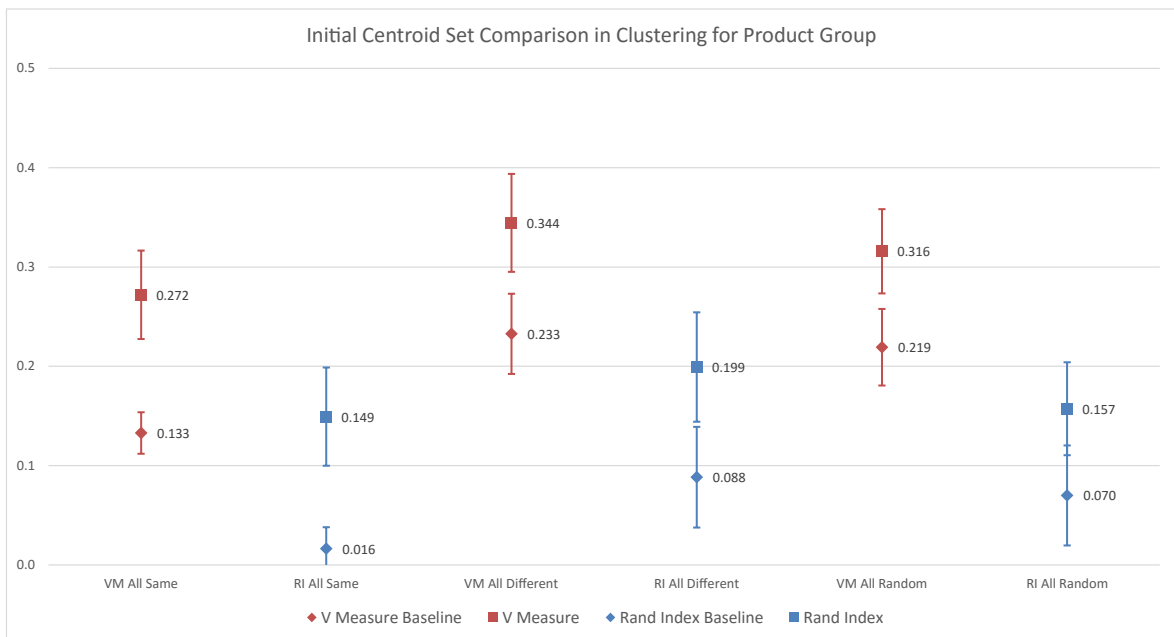
Figure 11.13 presents the results of the above setup for the scenario of clustering to 5 clusters, the Product Cleaning Groups.



Σχήμα 5.13: K means Clustering(Product Cleaning Group)

The best scores are visible when choosing initial centers that belong to different clusters, in which case there is a V-Measure of 31,4% and a Rand Index of 25,3%. On the contrary, the lowest score appears when the initial centers all belong to the same Product Cleaning Group.

Figure 11.14 follows the above methodology with the difference that the process is here targeted to 12 clusters, as we are examining the Product Group attribute.



Σχήμα 5.14: K means Clustering(Product Group)

As in the previous scenario, it can be observed that the result is best in the case of initial centroids belonging to different clusters and worst when they belong to the same.

It is quite obvious from the overall results that the clustering procedure has a poor general performance, as the scores do not exceed 35% for both attributes and all initial centers' combinations.

Κεφάλαιο 6

Epilogue

In this Chapter, we present the conclusions from the executed experiments and also our proposals for future work on the problem.

In Section 6.1 we present the conclusions derived from all experiments of Chapter 5, and summarize the performance of Classification and Clustering, evaluating the contribution of each parameter.

In Section 6.2, we describe possible topics that we did not have the opportunity to examine during the current thesis, and are worth looking into in the future.

6.1 Conclusions

In this section we summarize the observations from our experiments and discuss the final conclusions.

6.1.1 Classification

Regarding the efficiency of the distance methods, the Euclidean ones always gave top scores, along with the KL- Divergence and the Cosine distance. It was quite difficult to identify whether the average, diagonal or vector methods performed better, as for every metric, the optimal score was achieved for a different one, and there was no general pattern. The KS- Test and the Infinity Norm scores proved to be unsuitable for our data, and their results rarely overcame the baseline ones. We assume that this was due to the nature of our distributions, which came in the form of the probability transition matrices.

Comparing the classifiers, Nearest Centroid classifier and K-Nearest Neighbours classifier, the first one has a significantly better performance, considering both Accuracy and Kappa scores, for all classification attributes. This is an interesting observation. We conclude that regarding the classes as "centers" and assigning the data to the closest class outscored assigning the chunk to the closest neighbour, having examined the nearest ones. As mentioned in Chapter 5, the runtime for the k-NN was almost 8 times greater than that of the other classifier, which strongly validates the previous statement. This is due to the fact that k-NN has to calculate all distances from the k nearest class' centers, a task that can be very complex. The splitting in training and testing data had little effect on the resulting scores, and this is remarkable, if we take into account that the amount of data we used was few, and could show substantial differences for the different splitting percentages.

No significant variations were observed when running the k-NN for different k. This is a positive result, indicating that our data is insensitive to the selection of k, probably due to the form of our distributions.

6.1.2 Clustering

In all experiments, regardless of the parameterization of k-Means algorithm, the accuracy scores were quite low. So in general, the performance of k-Means on clustering the production chunks is considered to be poor.

As regard to the distances, the top three metrics were again Euclidean, Cosine and KL-Divergence, but with the Euclidean distances to rise significantly above the other two. Again, the KS-Test and

Infinity Norm scores are slightly over the baseline results, proving that they fail to correlate the data in an efficient way.

It was observed in the initialization of the algorithm with different types of centroids, that initial centroids belonging all to different clusters gave the maximum score, followed closely by centroids all belonging to random clusters. This is expected, as k-Means by definition has to assign all data to different clusters, so, starting off with distinct clusters as centroids, gave a more accurate result. On the contrary, when initial centroids were all in the same cluster, it would take many iterations to decentralize from the currently selected cluster, gradually calculate the others and assign the data correctly to them.

6.2 Future Work

Here are briefly presented the issues we did not have the time to investigate during the current thesis, and also issues that came up after our experiments and research, which could be further examined.

- **Distance Metrics**

More distance methods could be used to calculate the distances between the Product chunks. We could investigate methods especially targeted to two dimensional matrices, or vectorize the matrix in a different way (eg by row).

- **Clustering Algorithms** The poor performance of k - Means in Clustering could be a good motive to experiment with other clustering algorithms, which could prove to be more efficient in our distributions.

- **Classification using frequency vector and N-Dimensional transition matrix**

The Classification algorithms can be modified to receive different parameters. Instead of classifying the objects based on their messages' transition matrix, we could compute a simple frequency vector for each object, which depicts the percentage of appearance of each of our 45 different messages in this specific product. This vector is now the classifier's input, and all distances will be calculated between vectors, with the same methods as previously. In this way we evaluate each algorithm's performance, for all distances, when applied in the one dimensional frequency table, which obviously provides less information than the transitions one. The same thing could be applied for a three- (or N-) dimensional transition matrix, which examines more than one successive messages and finds the transition probabilities in their sequence. In this way, complexity of the information given to the classifier is increased, and accuracy of results and performance of the algorithm could be a lot different than the one we observed after our experiments.

- **Different Variable Selection**

In this work, we based our implementations on the Message Number sequence of each Production Chunk, as explained in the data Pre- Processing Chapter. In a different scenario, all variables characterizing our process (Temperature, Pressure, Agitation etc) could be taken into account in the algorithms' execution. This would require a different splitting of the initial data, and also a way to incorporate the extra parameters as attributes of the objects created. The input of each algorithm would need to be redefined, as the added variables are numerical, and so would the approach to calculate distances and centers, of either a class or a cluster in each execution. The experiments would have to run on a differently parameterized setup, to examine the contribution of the new parameters separately, keeping all others fixed or making suitable combinations to identify correlations.

- **Scoring of production process**

Production chunks refer to separate products being created. During the labelling process, we observed that production of the same product was not a solid fixed procedure, in terms of time duration, and order and kind of actions followed. This means that we had products all with the same attributes (Product Cleaning Group, Product Code, Product Group etc), but with different transition matrices, as the messages and their sequences in each one did not match.

In order to overcome this problem, a recording of the optimal procedure is needed, for all products, so as to compare the actual one to this and eliminate the variations. For each product, one could track the series of actions that are considered to be ideal, and create the chunk object for them. Then, the comparison of their transition matrices with the ones from the true products would give a measure of performance for each process, and allow us to score the Productions and determine if they were close to the ideal or not. This is a crucial first step for production optimization, and could boost the factory's overall performance, regarding quality of products, as well as minimization of production time.

- **Evaluation Metrics**

In evaluating the results, we used a few indicative metrics, commonly applied in classification and clustering assesment. There are though, many more which could be applied, providing different correlations and conclusions. Two examples on this are the Silhouette Coefficient, which measures classes' or clusters' consistency, and the Student's T- Test, which determines the probability that two sets of data are samples of the same distribution. In this way, evaluation of results could be advanced in different levels and provide us some interesting results.

Κεφάλαιο 7

Εισαγωγή

7.1 Ανάλυση Δεδομένων στην Παραγωγή

Τα συγχρονα συστήματα και διαδικασίες έχουν γίνει πολύ πολύπλοκα. Η ανάπτυξη της πληροφορικής και η ενσωμάτωση υπολογιστών υψηλής απόδοσης σε όλους τους κλάδους της βιομηχανικής διαδικασίας έχουν ως αποτέλεσμα τη συνεχή παραγωγή ενός τεράστιου όγκου δεδομένων. Μια παραγωγική διαδικασία παρακολουθείται από πολλούς αισθητήρες σε διάφορες στιγμές, έχοντας έτσι αυξημένη πολυπλοκότητα σε όλα της τα στάδια, και αποτελώντας μια πρόκληση για τη σημερινή επιστήμη. Τα τεχνολογικά επιτεύγματα έχουν βελτιώσει την παραγωγή από όλες τις απόψεις, έχουν όμως καταστήσει τις διαδικασίες πολύπλοκες και μη γραμμικές, οπότε είναι δύσκολο, ή αδύνατο να συλλέξουμε ακριβή ή άμεση πληροφορία μόνο από τον εξοπλισμό παρακολούθησης. Οι κλασικές μέθοδοι ανάλυσης χρειάζεται να αναπτυχθούν, ώστε να είναι αποδοτικές για τη διαδικασία, και νέες τεχνικές χρειάζεται να ενσωματωθούν σε παλαιότερες μεθόδους παραγωγής, ή ακόμα και να τις αντικαταστήσουν. [10]

Στην κατεύθυνση αυτή, οι βιομηχανίες υιοθετούν και αναπτύσσουν την ανάλυση των δεδομένων και στρατηγικές εκμετάλλευσής τους, με σκοπό να αυξήσουν την ποιότητα της παραγωγής, και την ασφάλεια της διαδικασίας και του ανθρώπου. Η προσέγγιση αυτή δεν είναι νέα, καθώς ιδέες business intelligence υλοποιούνται από τις αρχές της δεκαετίας του '70, και περιλαμβάνουν την εφαρμογή τεχνητών νευρωνικών δικτύων, και προβλεπτικού και προσαρμοστικού ελέγχου, σε μια προσπάθεια να βοηθήσουν τους ανθρώπινους χειριστές σε μια γραμμή παραγωγής, και, αν είναι δυνατόν, να τους εξαλείψουν. Εδώ έρχεται στο προσκήνιο η εξόρυξη δεδομένων, σαν ένας τρόπος να εξάγουμε έγκυρη, άγνωστη και κατανοητή πληροφορία από μια διαδικασία, και να την αξιοποιήσουμε στη λήψη σημαντικών επιχειρησιακών αποφάσεων. Όπως συνέβαινε παλαιότερα με τη χρήση σχεδίων και μαθηματικών εξισώσεων, για να παρακολουθήσουμε, να κατανοήσουμε και να βελτιστοποιήσουμε διαδικασίες, οι μέθοδοι ανάλυσης δεδομένων έχουν σήμερα παρόμοιο ρόλο. Η πολύπλευρη φύση του πεδίου, που περιλαμβάνει, μεταξύ άλλων, τη Μηχανική Μάθηση, την Επεξεργασία Εικόνας, και τη Στατιστική, στοχεύει στην εύρεση κανόνων και μοτίβων, που είναι μη ορατά στον άνθρωπο ή σε άλλες μεθόδους ανάλυσης.

Οι βιομηχανίες έχουν τη δυνατότητα να χρησιμοποιούν παλαιότερα δεδομένα, να αναγνωρίζουν συσχετίσεις μεταξύ ξεχωριστών διαδικασιών, όπως και τους παράγοντες που καθορίζουν την απόδοση μιας διαδικασίας, και έπειτα να βελτιώνουν τα στοιχεία που αποδεικνύεται ότι έχουν τη μεγαλύτερη επίδραση στην παραγωγή. Οι ίδιοι κανόνες εφαρμόζονται όταν εξετάζονται δεδομένα σε πραγματικό χρόνο, οπότε και χρειάζεται να αναπτυχθούν πιο πολύπλοκες τεχνικές, η άμεση ανάλυση τους όμως προσφέρει τη δυνατότητα βραχυπρόθεσμου προγραμματισμού και ενεργειών, αναγκαίων για τη βελτίωση και αποτελεσματικότητα της παραγωγής. Οι τεχνικές οπτικοποίησης είναι ένα σημαντικό εργαλείο για την αναγνώριση προτύπων, καθώς εφαρμόζονται ευρέως για καλύτερα αποτελέσματα, με τη χρήση γραφημάτων πιθανοτικών κατανομών και διαγραμμάτων clustering.

Οι ανάγκες και οι δυνατότητες αυτές αποτέλεσαν την αφετηρία της παρούσας διπλωματικής εργασίας, η οποία άρχισε σε μια προσπάθεια να αναγνωρίσουμε τα προβλήματα παρακολούθησης της παραγωγής και να αυτοματοποιήσουμε της διαδικασίες της ανάλυσης και της αξιοποίησης των δεδομένων.

7.2 Ορισμός Προβλήματος

Το περιβάλλον της Johnson & Johnson Hellas συνιστά ιδανικό πεδίο για να πειραματιστούμε και να υλοποιήσουμε όλα τα παραπάνω. Από τις αρχικές συζητήσεις, υπήρξε μεγάλο ενδιαφέρον και από τις δύο πλευρές για την εργασία. Ο κλάδος της ανάλυσης δεδομένων είναι τεράστιος, και οι εφαρμογές σε ένα τέτοιο χώρο είναι απεριόριστες.

Στο εργοστάσιο αυτό, κάθε γραμμή παραγωγής αποτελείται από πολλά μηχανήματα, το σημαντικότερο από τα οποία είναι το δοχείο ανάμιξης. Σε αντίθεση με άλλες μεθοδολογίες παραγωγής, το εργοστάσιο παράγει τα προϊόντα του σε παρτίδες. Αυτό σημαίνει ότι σε κάθε δοχείο παράγεται μια μεγάλη ποικιλία προϊόντων, ανάλογα με τη ζήτηση της αγοράς. Η ανάμιξη είναι ένα σημαντικό βήμα της παραγωγικής διαδικασίας, καθώς είναι το στάδιο στο οποίο δημιουργείται το προϊόν. Όπως θα εξηγηθεί και στο επόμενο κεφάλαιο λεπτομερώς, το δοχείο παρακολουθείται και ελέγχεται από ένα PLC, από το οποίο και μας δίδεται το αρχείο των δεδομένων.

Η πρώτη ανάγκη που μας μεταφέρθηκε ήταν αυτή της οπτικοποίησης των δεδομένων παραγωγής. Τα δεδομένα καταγράφονταν σε καθημερινή βάση, αλλά δε γινόταν κάποια περαιτέρω ανάλυση και αξιοποίησή τους. Η υλοποίηση μιας γραφικής αναπαράστασής τους μας επέτρεψε να τα κατανοήσουμε καλύτερα και διευκόλυνε την εφαρμογή τεχνικών Μηχανικής Μάθησης.

Η δεύτερη είχε αρκετά σκέλη, που περιελάμβαναν την αξιολόγηση της διαδικασίας ανάμιξης, και τη σύγκριση με την "ιδανική διαδικασία", με σκοπό τη βελτιστοποίηση της παραγωγής. Κάθε προϊόν έχει τα ίδια αυστηρά κριτήρια ποιότητας, τα οποία εντοπίζονται εύκολα μέσω της χημικής ανάλυσης του τελικού προϊόντος. Η διαδικασία που είναι δύσκολο να μετρηθεί και να αξιολογηθεί, είναι αυτή της εκτέλεσης της συνταγής κατά τη δημιουργία ενός συγκεκριμένου προϊόντος. Αυτό παρουσιάζει μεγάλο ενδιαφέρον, καθώς, αν και όλα τα προϊόντα που παράγονται πληρούν τα ποιοτικά κριτήρια, παρόλα αυτά δεν υπάρχουν δύο ακριβώς ίδιες χρονοσειρές ενεργειών. Παραθέτοντας ένα ενδεικτικό μέγεθος, κάθε παρτίδα ολοκληρώνεται σε περίπου 400 ενέργειες.

Αφετηρία της έρευνάς μας ήταν η εφαρμογή μεθόδων classification και clustering, με σκοπό να είμαστε σε θέση να αξιολογήσουμε και να συγκρίνουμε διαδικασίες. Θεωρούμε ότι αποτελεί ένα σημαντικό πρώτο βήμα στον πειραματισμό με χρονοσειρές, και στη μετατροπή και χρήση τους σε εφαρμογές Μηχανικής Μάθησης.

Η βασικές προκλήσεις, οι οποίες ανέκυψαν κατά την επεξεργασία των δεδομένων, ήταν οι ακόλουθες:

- **Καθαρισμός Δεδομένων**

Τα δεδομένα από το PLC ήταν σχεδόν πλήρη και αρκετά ακριβή, σε αντίθεση με τα δεδομένα που καταγράφονται από τους χειριστές σε κάθε δοχείο χωρίς αυτόματο τρόπο.

- **Επιλογή Μεταβλητών** Για τα δεδομένα μας, σε κάθε χρονική στιγμή είχαμε την τιμή 17 διαφορετικών μεταβλητών, που αφορούσαν στην κατάσταση του δοχείου, κάτι που αύξανε την πολυπλοκότητα συσχέτισής τους, και καθιστούσε δύσκολη την εξαγωγή ενός συμπεράσματος με την απλή παρατήρηση των τιμών τους.

- **Σύγκριση χρονοσειρών διαφορετικού μήκους** Ένα από τα πιο δύσκολα προβλήματα που αντιμετωπίσαμε ήταν η σύγκριση χρονοσειρών άνισου μήκους. Κάθε παρτίδα αποτελείται από ενέργειες που αντιπροσωπεύονται από κωδικούς (MsgNum). Καθώς δύο ίδια προϊόντα μπορούν να παραχθούν με διαφορετικές ενέργειες, όσον αφορά στη σειρά τους και στο συνολικό τους αριθμό, η σύγκριση τους δε μπορούσε να γίνει ένα προς ένα.

- **Υπολογισμός Αποστάσεων**

Οι περισσότερες υλοποιήσεις αλγορίθμων Μηχανικής Μάθησης χειρίζονται στοιχεία που συνιστούν ένα διάνυσμα δεδομένων. Χρειάστηκε να αποφασίσουμε ποια θα είναι η μορφή των

στοιχείων, των οποίων θα υπολογίσουμε την απόσταση, καθώς και τις μεθόδους για την εύρεση αυτής.

7.3 Προτεινόμενες Λύσεις και Υλοποίηση

Για να αντιμετωπίσουμε τις παραπάνω προκλήσεις, και να αξιοποιήσουμε βέλτιστα τα δεδομένα, ακολουθήσαμε τα εξής βήματα που συνιστούν τη λύση μας:

- **Προ-επεξεργασία Δεδομένων και Καθαρισμός:**

Συνδυάζοντας τα αρχικά δεδομένα από το PLC με αυτά από την ανθρώπινη καταγραφή της διαδικασίας, και, με βάση γνωστούς εμπειρικούς εσωτερικούς κανόνες, ομαδοποιήσαμε ξεχωριστές ενέργειες, με σκοπό τη δημιουργία αντικειμένων τα οποία αναπαριστούν μια διαδικασία που λαμβάνει χώρα. Αυτή ήταν είτε η παραγωγή προϊόντος, είτε ο καθαρισμός. Αποφασίσαμε να χρησιμοποιήσουμε μόνο μία μεταβλητή, τον κωδικό μηνύματος (MsgNum), καθώς μετά από πειράματα καταλήξαμε στο συμπέρασμα ότι ήταν η καλύτερη και πιο απλή προσέγγιση, που εξαλείφει κάθε πολυπλοκότητα λόγω του μεγάλου αριθμού μεταβλητών.

- **Visualization:**

Υλοποιήθηκε ένα εργαλείο οπτικοποίησης των δεδομένων, για την απεικόνιση όλων των διαφορετικών μεταβλητών, με σκοπό τη δημιουργία μιας γραφικής και διαδραστικής απόδοσης των αρχικών δεδομένων. Αυτό αποτέλεσε και μια προσπάθεια για την καλύτερη κατανόηση των δεδομένων, και την παρατήρηση πιθανών ορατών μοτίβων από τα αρχικά, από τις αρχικά δύσκολες και δύσκολες στη διαχείριση χρονοσειρές.

- **Σύγκριση χρονοσειρών διαφορετικού μήκους:**

Εξαλείψαμε την εξάρτηση των ενεργειών από την παράμετρο του χρόνου. Αυτό συνέβη λόγω της επαναληψιμότητάς τους, και του γεγονότος ότι πολλές από τις ενέργειες πραγματοποιούνται σε συνεχείς επαναλήψεις, χωρίς να επηρεάζουν το αποτέλεσμα. Η αρχική μας προσέγγιση ήταν να δημιουργήσουμε ένα κανονικοποιημένο διάνυσμα συχνοτήτων για όλα τα μηνύματα σε μια παρτίδα. Αποφασίσαμε όμως να ενσωματώσουμε περισσότερες πληροφορίες από τα αρχικά δεδομένα. Η τελικά μας προσέγγιση ήταν να κατασκευάσουμε έναν κανονικοποιημένο πίνακα μεταβάσεων για κάθε παρτίδα. Κάθε κελί του πίνακα αντιπροσωπεύει την πιθανότητα μετάβασης από ένα μήνυμα σε ένα άλλο. Με απλούστερα λόγια, το κελί[A][B] είναι η πιθανότητα εκτέλεσης της ενέργειας B μετά την ενέργεια A. Ο τρόπος αυτός αποδείχθηκε ιδιαίτερα χρήσιμος, και μας βοήθησε στο κομμάτι της Μηχανικής Μάθησης.

- **Υπολογισμός Αποστάσεων:**

Σχεδιάστηκε και υλοποιήθηκε ένας τρόπος υπολογισμού αποστάσεων μεταξύ διαφορετικών διδιάστατων αντικειμένων (πινάκων), με σκοπό τη σύγκριση των αντικειμένων που δημιουργήσαμε και την εύρεση της σχετικής τους απόστασης. Η πρώτη προσέγγιση ήταν η μετατροπή του πίνακα σε διάνυσμα, και η δεύτερη ήταν η δημιουργία των μεθόδων υπολογισμού αποστάσεων και η χρήση τους στους αλγορίθμους. Η υλοποίηση αυτή περιείχε πολλές γνωστές μεθόδους (Euclidean, Cosine, KL-Divergence, KS-Test, Infinity Norm) και αξιοποιήθηκε περαιτέρω στο κομμάτι της εφαρμογής των αλγορίθμων Μηχανικής Μάθησης. Κατά την εκτέλεση των πειραμάτων μας, παρατηρήσαμε σημαντικές διαφορές στα αποτελέσματα, ανάλογα με τη μέθοδο απόστασης που είχαμε επιλέξει.

- **Classification:**

Πειραματιστήκαμε με δύο αλγορίθμους Classification στα αντικείμενα που δημιουργήθηκαν μετά την προ-επεξεργασία των δεδομένων, τον Nearest Centroid Classifier, και τον k-Nearest Neighbours Classifier. Κάθε προϊόν έχει διαφορετικά γνωρίσματα (Product Group, Product

Cleaning Group), ήδη γνωστά από τη διαδικασία απόδοσης ετικέτας, επομένως το classification πραγματοποιήθηκε με βάση αυτά τα χαρακτηριστικά. Η απόσταση των αντικειμένων ήταν μία από της παραμέτρους κάθε εκτέλεσης. Για τον αλγόριθμο k-NN, η τιμή του k ήταν μια παράμετρος που επίσης εξετάστηκε. Τα δεδομένα εκπαιδεύτηκαν, και επιλέχθηκε και ένα δείγμα testing, ώστε να προσδιορίσουμε την αποδοτικότητα της διαδικασίας, η οποία μετρήθηκε με δύο μεθόδους, τις Accuracy και Kappa Coefficient. Τα πειράματα έδωσαν ενδιαφέροντα αποτελέσματα, με σκορ αποδοτικότητας μέχρι και 85%.

- **Clustering:**

Προσαρμόσαμε τον αλγόριθμο k-Means για Clustering, και, βασιζόμενοι και πάλι στα διαφορετικά γνωρίσματα των αντικειμένων- προϊόντων, τα χωρίσαμε σε διαφορετικά clusters. Η παραμετροποίηση του αλγορίθμου εστίασε σε δύο κατευθύνσεις, την επιλογή των αρχικών κέντρων, και την κατάλληλη επιλογή της μεθόδου υπολογισμού απόστασης ανάμεσα στα αντικείμενα. Δύο ήταν μέθοδοι αξιολόγησης των πειραμάτων, οι V-Measure και Rand Index, και η απόδοση αυτών προέκυψε αρκετά χαμηλή, με μέγιστο σκορ γύρω στο 36%.

7.4 Δομή Διπλωματικής Εργασίας

Στο κεφάλαιο 8 ορίζουμε το πρόβλημα στο οποίο θα εργαστούμε, και όλες τις διαφορετικές παραμέτρους του, περιγράφουμε όλα τα βήματα για την επεξεργασία και κατανόηση των δεδομένων, και παρουσιάζουμε το εργαλείο visualization που δημιουργήσαμε για τις ανάγκες της εταιρείας από τα δεδομένα της παραγωγής.

Στο κεφάλαιο 9 επικεντρωνόμαστε στο θεωρητικό υπόβαθρο που απαιτείται για τις εφαρμογές Εξόρυξης Δεδομένων και Μηχανικής Μάθησης. Ορίζουμε τις παραπάνω έννοιες, και αναλύουμε αυτές του Classification και του Clustering, που χρησιμοποιήθηκαν στην παρούσα εργασία. Παρουσιάζουμε επίσης τις έννοιες της Απόστασης και της Αξιολόγησης στη Μηχανική Μάθηση, και τις μεθόδους και τεχνικές που χρησιμοποιήθηκαν για αυτούς τους υπολογισμούς.

Στο κεφάλαιο 10 παρουσιάζουμε την υλοποίηση της λύσης μας. Οι διαδικασίες των Classification και Clustering παραμετροποιούνται και εξηγούνται αναλυτικά, ενώ εφαρμόζονται όλοι οι τρόποι υπολογισμού της απόστασης και του σκορ απόδοσης στα πειράματα που εκτελούνται.

Στο κεφάλαιο 11 παραθέτουμε τα αποτελέσματα και τα διαγράμματα των πειραμάτων των Classification και Clustering, για τις διαφορετικές αρχικοποιήσεις των αλγορίθμων και τις υλοποιήσεις τους.

Τέλος, στο κεφάλαιο 12, συνοψίζουμε τα συμπεράσματά μας από τα αποτελέσματα, και αναφερόμαστε στα θέματα στα οποία δεν είχαμε την ευκαιρία να εργαστούμε, και που μπορούν να ληφθούν υπόψη σε μελλοντική εργασία.

Κεφάλαιο 8

Περιγραφή του Προβλήματος και των Δεδομένων

Σε αυτό το Κεφάλαιο, περιγράφουμε όλα τα χαρακτηριστικά του προβλήματος της βιομηχανίας που αντιμετωπίζουμε στην παρούσα διπλωματική εργασία, και τα αρχικά βήματα που κάναμε για να αξιοποιήσουμε τα δεδομένα παραγωγής που μας δόθηκαν, με ένα αποτελεσματικό τρόπο.

Στην Ενότητα 8.1, αρχικά περιγράφουμε την διαδικασία ανάμειξης, στην οποία και επικεντρωνόμαστε, ορίζοντας τα πιο σημαντικά της σημεία. Στη συνέχεια, παρουσιάζουμε όλες τις παραμέτρους του αρχικού συνόλου δεδομένων και ορίζουμε συνοπτικά τα βήματα που ακολουθήθηκαν στην διαδικασία και τις βασικές προκλήσεις που αντιμετωπίσαμε.

Στην Ενότητα 8.2, περιγράφουμε τις τεχνικές προ-επεξεργασίας που αξιοποιήσαμε, οι οποίες είναι χωρισμένες σε τρεις κατηγορίες, Flattening, Labelling και Object Creation.

Τέλος στην Ενότητα 8.3, αναλύουμε τον τρόπο με τον οποίο αναπαραστήσαμε γραφικά τα δεδομένα μας και παρουσιάζουμε το εργαλείο που παραμετροποιήσαμε για τους σκοπούς της εταιρίας.

8.1 Περιγραφή Προβλήματος

Σε αυτή την Ενότητα, παρουσιάζονται οι βασικές παράμετροι και μεταβλητές.

8.1.1 Περιγραφή Θέματος

Το θέμα το οποίο μελετήσαμε είναι ένα υπαρκτό πρόβλημα. Τα δεδομένα, προέρχονται από τη Johnson & Johnson Hellas, και συγκεκριμένα, τα δεδομένα που χρησιμοποιήθηκαν στη Μηχανική Μάθηση, είναι από ένα συγκεκριμένο Δοχείο το οποίο βρίσκεται σε γραμμή παραγωγής παρτίδων, σε ένα εργοστάσιο στην Ελλάδα. Σε κάθε Δοχείο μπορεί να πραγματοποιηθεί ένας συγκεκριμένος αριθμός ενεργειών, μέσω των οποίων μπορεί να παραχθεί μια μεγάλη ποικιλία προϊόντων. Παρακάτω μπορείτε να βρείτε περισσότερες πληροφορίες για το κάθε στοιχείο.

- **Δοχείο**

Αυτό το δοχείο έχει κάποιους μηχανισμούς αυτοματοποίησης αλλά οι περισσότερες διεργασίες εκτελούνται από τον χειριστή. Δεν θεωρούμε χρήσιμο το να αναλύσουμε ακριβώς ποιες είναι αυτόματες και ποιες εκτελούνται από τον χειριστή. Η επεξήγηση του αρχικού συνόλου δεδομένων θα γίνει στην επόμενη υπο-ενότητα.

- **Ανάμειξη/Παραγωγή**

Η βασική λειτουργία του δοχείου είναι η ανάμιξη διαφορετικών υλών με σκοπό να παράξουμε κάποιο συγκεκριμένο προϊόν. Τα προϊόντα ποικίλλουν, από λοσιόν μέχρι έλαια και κρέμες. Κάποια από τα βασικά υλικά παρέχονται από ένα δίκτυο του εργοστασίου και άλλα από μια καταπακτή από τον χειριστή του δοχείου.

- **Πιθανές Ενέργειες**

Οι πιθανές ενέργειες που μπορούν να πραγματοποιηθούν είναι: Ανάδευση, Θέρμανση, Ψύξη, Εισαγωγή Ρευστών, Δημιουργία Κενού. Πολλά από τα παραπάνω μπορούν να γίνουν ταυτόχρονα, για παράδειγμα, θέρμανση και ανάδευση.

- **Καταγεγραμμένες Τιμές**

Όπως θα περιγραφεί στην επόμενη υποενότητα, οι ενέργειες που εκτελούνται στο δοχείο, αποθηκεύονται από ένα PLC ως μια χρονοσειρά από μηνύματα και τις κατάλληλες μετρήσεις και χρονοσημάνσεις. Υπάρχει ένα συγκεκριμένο χαρακτηριστικό σχετικά με τη μορφή του CSV στην έξοδο, το οποίο είναι το ακόλουθο:

Οι τιμές καταγράφονται μόνο κατά τις χρονικές στιγμές που συμβαίνει κάποια ενέργεια. Για κάθε μεταβλητή έχουμε δύο τιμές, μια είναι η πραγματική τιμή που καταγράφεται και η άλλη είναι η τιμή που έχει οριστεί από τον χειριστή και προσπαθεί να φθάσει το δοχείο μέσω των εντολών του PLC. Για παράδειγμα, για την θερμοκρασία έχουμε την μετρούμενη τιμή 50, 87, και τον στόχο που θέλει να επιτύχει το δοχείο, που μπορεί να είναι η θερμοκρασία 75 βαθμοί Κελσίου. Αυτές οι τιμές αποθηκεύονται ως δύο διαφορετικές γραμμές στο αρχικό CSV.

8.1.2 Περιγραφή του αρχικού Συνόλου Δεδομένων

Αυτό είναι ένα πολύ σημαντικό μέρος της εργασίας μας, το οποίο θα βοηθήσει τον αναγνώστη να κατανοήσει καλύτερα την διαδικασία προ-επεξεργασίας των δεδομένων. Τα δεδομένα που λάβαμε, είναι η απευθείας έξοδος του PLC που ελέγχει το δοχείο που εξετάζουμε. Η μορφοποίηση του αρχείου είναι CSV και τα δεδομένα αντιπροσωπεύουν ένα έτος εργασιών.

- **Πληροφορίες του αρχικού Συνόλου Δεδομένων**

- Μορφοποίηση: CSV
- Χρονική Περίοδος: 16/02/2015 – 13/02/2016
- Σύνολο Γραμμών Δεδομένων: 132005
- Μέγεθος Αρχείου: 25Mb
- Ιδιότητες (Στήλες):
 - * StateAfter
 - * MsgNumber
 - * Temp
 - * Pressure
 - * Agitation1
 - * Agitation2
 - * Homogen
 - * Pump Power
 - * Raw Meterials
 - * TimeString

- **Screenshot του αρχικού CSV**

	A	B	C	D	E	F	G	H	I	J
1	StateAfter	MsgNumber	Temp	Pressure	Agitation1	Agitation2	Homogen	Pump Power	Raw Materials	TimeString
2	1	552	78.62	0	14.84028	24.75579	0	0.8101852		16.02.2015 09:41:08
3	1	564	78.62	0	14.84028	24.75579	0	0.8101852		16.02.2015 09:41:08
4	1	553	75	-500	15	30	750	60		16.02.2015 09:41:08
5	1	565	75	-500	15	30	750	60		16.02.2015 09:41:08
6	1	552	78.89	1.15741	14.84028	24.75579	0	0.8391203		16.02.2015 09:43:39
7	1	564	78.89	1.15741	14.84028	24.75579	0	0.8391203		16.02.2015 09:43:39
8	0	553	75	-500	15	30	750	60		16.02.2015 09:43:39
9	0	565	75	-500	15	30	750	60		16.02.2015 09:43:39
10	1	522	50.715	1.15741	6.534722	25.62182	0	0.8101852	21.20999	16.02.2015 10:25:43
11	1	523	75	-500	15	31	750	60	500	16.02.2015 10:25:43
12	0	536	77.1	0	14.83333	24.77691	0	0.8101852		16.02.2015 10:25:43
13	0	562	77.1	0	14.83333	24.77691	0	0.8101852		16.02.2015 10:25:43
14	1	537	75	-500	15	30	750	60		16.02.2015 10:25:43
15	1	563	75	-500	15	30	750	60		16.02.2015 10:25:43

Σχήμα 8.1: Αρχική μορφή του αρχείου δεδομένων CSV

• Επεξήγηση των σημαντικών μεταβλητών των δεδομένων

- *StateAfter*: Αυτή η μεταβλητή μπορεί να είναι είτε 0 είτε 1. Αυτό σημαίνει ότι αν είναι 1 η ενέργεια εκκινεί, ενώ αν είναι 0 η ενέργεια σταματάει εκείνη την χρονική στιγμή.
- *MsgNumber*: Μέσω του αρχείου απεικόνισης μπορούμε να δούμε την εξήγηση του συγκεκριμένου μηνύματος, δηλαδή ποια ενέργεια έγινε στο δοχείο τη χρονική στιγμή αυτή. Παρακάτω μπορείτε να δείτε μερικά παραδείγματα τέτοιων μηνυμάτων.

	A	B
1	500	ΑΝΑΚΥΚΛΟΦΟΡΙΑ ΑΠΌ ΠΑΝΩ
2	501	ΑΝΑΚΥΚΛΟΦΟΡΙΑ ΑΠΌ ΠΑΝΩ (Set Points)
3	504	ΑΔΕΙΑΣΜΑ ΣΕ TNT
4	505	ΑΔΕΙΑΣΜΑ ΣΕ TNT (Set Points)
5	522	ΕΙΣΑΓΩΓΗ ΑΠΙΟΝΙΣΜΕΝΟΥ ΝΕΡΟΥ
6	523	ΕΙΣΑΓΩΓΗ ΑΠΙΟΝΙΣΜΕΝΟΥ ΝΕΡΟΥ (Set points)

Σχήμα 8.2: Παραδείγματα MsgNumbers

- *TimeString*: Αυτή η μεταβλητή μας δίνει την χρονική στιγμή που έγινε η συγκεκριμένη ενέργεια με ανάλυση δευτερολέπτου.

• Διαχωρισμός

Η συνεχής χρονοσειρά μηνυμάτων ενεργειών δεν μας έδωσε καμία πληροφορία σχετικά με την έναρξη ή την λήξη των παρτίδων. Ο διαχωρισμός τους ήταν ένα πολύ σημαντικό κομμάτι, το οποίο μας επέτρεψε να εφαρμόσουμε τις τεχνικές Μηχανικής Μάθησης, όπως θα εξηγηθεί στο Κεφάλαιο της Υλοποίησης.

Ο διαχωρισμός των αρχικών δεδομένων σε "chunks" έγινε αξιοποιώντας συγκεκριμένους κανόνες που μας τους παρείχαν ανώτερα στελέχη του τομέα διαχείρισης παραγωγής.

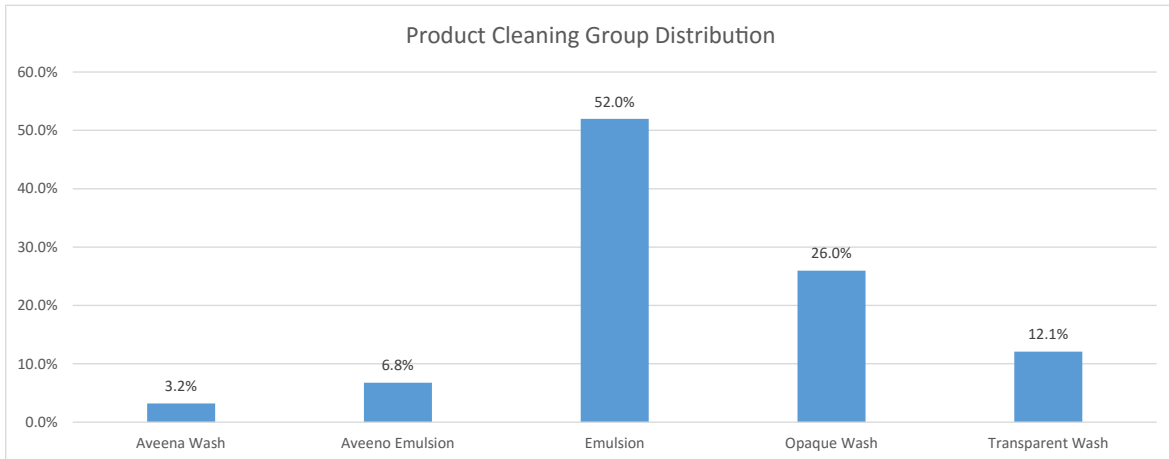
• Τιτλοφόρηση

Ένα ακόμα σημαντικό βήμα ήταν η διαδικασία τιτλοφόρησης των δεδομένων, δηλαδή το να θέσουμε τις κατάλληλες ετικέτες στα κατάλληλα στοιχεία. Στα αρχικά δεδομένα δεν υπήρχε πληροφορία σχετικά με τα προϊόντα που παρήχθησαν ή τις κατηγορίες που αυτά ανήκουν. Με σκοπό να μπορέσουμε να θέσουμε τις κατάλληλες ετικέτες, μας δόθηκαν δύο αρχεία.

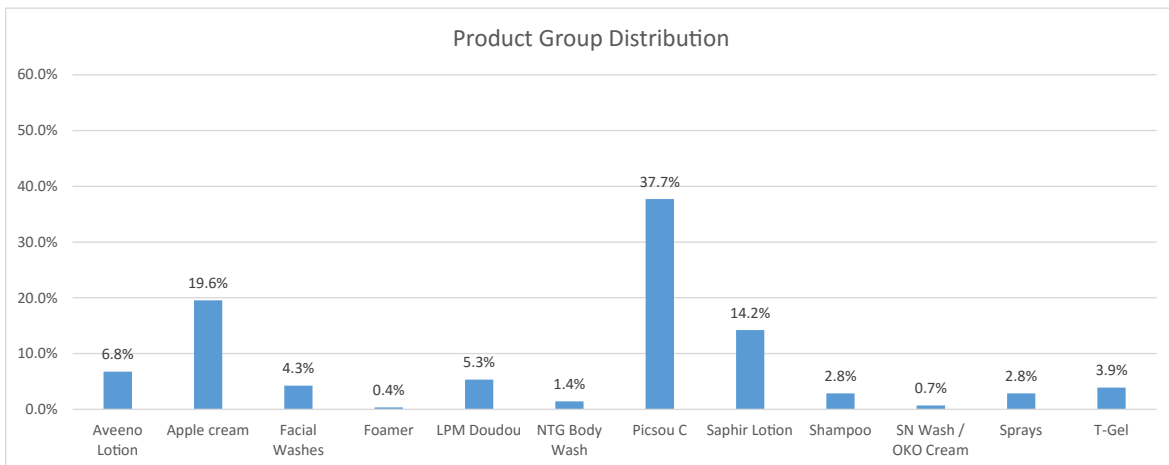
Το πρώτο ήταν ένα αρχείο καταγραφής από τους χειριστές, όπου ο καθένας σημείωνε τον κωδικό του προϊόντος που φτιάχτηκε στην εκάστοτε παρτίδα και την ώρα λήξης της. Το δεύτερο ήταν ένα αρχείο το οποίο περιείχε τους κωδικούς των προϊόντων και διάφορα χαρακτηριστικά τους, όπως το όνομα του προϊόντος, και τις συγκεκριμένες κατηγορίες που ανήκει, αναλόγως με την ομαδοποίηση, είτε σε επίπεδο Product Cleaning Group είτε σε επίπεδο Product Group.

Αυτές οι δύο κατηγοριοποιήσεις θα παίξουν πολύ σημαντικό ρόλο τόσο στο classification όσο και στο clustering των δεδομένων μας.

- Το Product Cleaning Group αποτελείται από 5 γκρουπ και όλα τα προϊόντα εντάσσονται σε κάποιο από αυτά. Η ακριβής κατανομή μπορεί να βρεθεί στην Εικόνα 8.3
- Το Product Group είναι μια διαφορετική κατηγοριοποίηση και αποτελείται από 12 γκρουπ, και, όπως και στην προηγούμενη, έτσι και σε αυτή, όλα τα προϊόντα ανήκουν σε κάποια κατηγορία. Η ακριβής κατανομή μπορεί να βρεθεί στην Εικόνα 8.4



Σχήμα 8.3: Κατανομή του Product Cleaning Group



Σχήμα 8.4: Κατανομή του Product Group

8.1.3 Ορισμός των υποπροβλημάτων

Δύο από τα πιο σημαντικά προβλήματα που αντιμετωπίσαμε, και ξεπεράσαμε, είναι τα ακόλουθα.

- **Επιλογή Μεταβλητών**

Στα αρχικά δεδομένα που λάβαμε είχαμε αρκετές μεταβλητές με τις οποίες μπορούσαμε να δουλέψουμε. Ύστερα από κάποιους αρχικούς ελέγχους, αποφασίσαμε να προχωρήσουμε αξιοποιώντας μόνο την μεταβλητή MsgNum. Στην Ενότητα 12.2 αναφέρουμε ότι πρέπει να δοθεί ιδιαίτερη έμφαση σε μελλοντική δουλειά για να διερευνηθούν τα αποτελέσματα των αλγορίθμων αν ενσωματωθούν και άλλες μεταβλητές των αρχικών δεδομένων.

- **Σύγκριση ακολουθιών δεδομένων άνισου μήκους**

Όπως καταλαβαίνει κανείς, κάθε παρτίδα μπορεί να παραχθεί με διαφορετικούς συνδυασμούς δράσεων. Αυτό οφείλεται στο γεγονός ότι κάθε χειριστής μπορεί να κάνει μικρές αλλαγές στη διαδικασία. Κάθε διεργασία αποτελείται από περίπου 400 μηνύματα, οι οδηγίες δεν είναι τόσο λεπτομερείς και τα υλικά που χρησιμοποιούνται έχουν κάποια όρια ανοχής.

Αυτό το χαρακτηριστικό της διαδικασίας καθιστά κάθε παρτίδα μοναδική, όσον αφορά στο μήκος της χρονοσειράς και στην ακριβή σειρά των εντολών. Όπως θα γίνει σαφές στο κεφάλαιο 10, τα στοιχεία που θέλουμε να συγκρίνουμε θα πρέπει να έχουν το ίδιο μέγεθος, προκειμένου να γίνει το classification ή το clustering. Η ιδέα που υλοποιήσαμε είναι η κατασκευή ενός Πίνακα Μεταβάσεων για κάθε batch. Αυτός ο Πίνακας έχει πάντα μέγεθος $M \times M$, όπου είναι ο αριθμός των μοναδικών κωδικών μηνυμάτων, το οποίο είναι το ίδιο για κάθε παρτίδα. Στην περίπτωση μας, $M = 45$.

8.2 Προεπεξεργασία των Δεδομένων

Σε αυτή την Ενότητα όλες οι μέθοδοι που χρησιμοποιήθηκαν για την προεπεξεργασία των δεδομένων θα αναλυθούν. Ο ρόλος αυτού του τμήματος της εργασίας ήταν να γίνει καθαρισμός του αρχικού συνόλου δεδομένων, ώστε να είναι πιο εύκολο να γίνει η απεικόνιση των δεδομένων, και να εφαρμοστούν οι αλγόριθμοι Classification και Clustering.

8.2.1 Flattening - Καθαρισμός

Όπως εξηγήθηκε στην υποενότητα 8.1.1 τα δεδομένα έρχονται σε ζεύγη μετρήσεων και Set-Point. Για τολόγο αυτό, έχουμε αναπτύξει μια διαδικασία flattening, προκειμένου να δημιουργούμε μία γραμμή για κάθε ενέργεια, η οποία να περιέχει και όλες τις μετρήσεις και όλα τα Set-Points. Παράλληλα, τα δεδομένα καθαρίστηκαν από τις δοκιμαστικές τιμές που είχαν εισαχθεί από τους χρήστες. Αυτό έγινε με βασικό σκοπό να βοηθήσει τη διαδικασία οπτικοποίησης.

8.2.2 Τιτλοφόρηση

Όπως αναφέρθηκε προηγουμένως, η απόδοση ετικετών ήταν ένα σημαντικό και απαραίτητο βήμα, διότι, μέσα από τις ετικέτες που δώσαμε στο δεδομένα, ήμασταν σε θέση να αξιολογήσουν την απόδοση των αλγορίθμων.

Τα αρχεία που χρησιμοποιήθηκαν σε αυτό το κομμάτι ήταν τρία: το αρχείο καταγραφής από τους χειριστές, το αρχείο απεικόνισης μεταξύ των κωδικών των προϊόντων και των χαρακτηριστικών του προϊόντος, καθώς και η αρχική χρονοσειρά μηνυμάτων δράσεων που περιέχει όλα τα αρχικά στοιχεία, διαχωρισμένα σε κομμάτια.

Το πρώτο βήμα ήταν να ενσωματώσουμε τα χαρακτηριστικά των προϊόντων για τις παρτίδες που βρέθηκαν στο αρχείο καταγραφής. Αυτό ήταν σχετικά εύκολο, γιατί είχαμε το μοναδικό κωδικό του προϊόντος από το αρχείο καταγραφής, ο οποίος ήταν αυτό που χρειαζόταν προκειμένου να ανακτηθούν όλες οι απαραίτητες πληροφορίες από το αρχείο απεικόνισης. Το μόνο πρόβλημα που συναντήσαμε ήταν η σπάνια περίπτωση ορισμένων λαθών στους κωδικούς προϊόντος, όπως αυτοί ήταν καταγεγραμμένοι στο αρχείο καταγραφής. Τα χαρακτηριστικά που ανακτήθηκαν ήταν το Production Cleaning Group και το Production Group.

Το δεύτερο βήμα ήταν να συνενώσουμε τα προϊόντα, όπως μας δόθηκαν στο αρχείο καταγραφής με τις παρτίδες από το αρχικό σύνολο δεδομένων. Ενώ η διαδικασία φαινόταν εύκολη, αποδείχθηκε ότι ήταν αρκετά απαιτητική. Ο τρόπος που έγινε αυτό ήταν να βρεθεί ο χρόνος λήξης μιας συγκεκριμένης παρτίδας στο αρχείο καταγραφής και, στη συνέχεια, να εντοπισθεί από το αρχικό σύνολο δεδομένων, το κομμάτι που είχε την πλησιέστερη ώρα λήξης. Τρία ήταν τα κύρια προβλήματα σε αυτή τη διαδικασία:

- Ο λανθασμένος τρόπος με τον οποίο εισέρχονται οι χρονοσημάνσεις στο αρχείο καταγραφής. Αυτή είναι μια διαδικασία που εκτελείται με το χέρι, το οποίο σημαίνει ότι υπήρχαν πολλά λάθη και μη χρησιμοποιήσιμες πληροφορίες.
- Επειδή η διαδικασία διαχωρισμού του αρχικού συνόλου δεδομένων έγινε από μη-τέλειους κανόνες, μερικές παρτίδες συγχωνεύθηκαν και άλλες διαχωρίστηκαν σε δύο.
- Η ασυμφωνία μεταξύ των δύο αρχείων. Προκειμένου να γίνει η αντιστοιχία, υλοποιήσαμε έναν αλγόριθμο, ο οποίος απέδιδε τις ετικέτες στο πλησιέστερο κομμάτι παραγωγής. Επειδή βρέθηκαν πολύ μεγάλες διαφορές, αποφασίσαμε να χρησιμοποιήσουμε, για το μέρος της Μηχανικής Μάθησης, μόνο τα κομμάτια εκείνα που είχαν χρονική διαφορά μικρότερη από 7 , μεταξύ του χρόνου λήξης που βρέθηκε στο αρχείο καταγραφής και εκείνου που διαπιστώθηκε από τον διαχωρισμό του αρχικού συνόλου δεδομένων.

8.2.3 Δημιουργία Αντικειμένων για τη Μηχανική Μάθηση

Το τελευταίο βήμα της προ-επεξεργασίας των δεδομένων ήταν να κατασκευασθεί ένα νέο σύνολο δεδομένων, που περιείχε τα προηγούμεως επισημασμένα κομμάτια, σε μια αξιοποιήσιμη μορφή. Αποφασίσαμε να δημιουργήσουμε την **Κλάση Chunk** και να μετατρέψουμε την κάθε παρτίδα σε ένα στιγμιότυπο αυτής της κλάσης. Όταν δημιουργήθηκαν όλα, τα αποθηκεύσαμε σε ένα αρχείο JSON, για ευκολότερη πρόσβαση και για να διασφαλισθεί η ακεραιότητα της διαδικασίας μας.

• Περιγραφή της Κλάσης Chunk

```

1 class Chunk(object):
2
3     def __init__(self, **entries):
4         super(Chunk, self).__init__()
5
6         #The start time of this chunk as found on the initial data-set
7         self.start_time = None
8         #The end time of this chunk as found on the initial data-set
9         self.end_time = None
10        #The chunk_type can be either "Production" or "Cleaning"
11        self.chunk_type = None
12
13        #The product code as found from the log file.
14        self.pr_code = None
15        #The end time as found from the log file
16        self.pr_logged_end_time = None
17
18        #The following attributes were found from the referance file through the
19        production code
20        self.pr_name = None
21        self.pr_group = None
22        self.pr_cl_group = None
23
24        #This is used for both classification and clustering
25        #This attribute is set to the name of the class/cluster it belongs
26        #and is being evaluated at the end of each experiment
27        self.cluster = None
28
29        #If this chunk represents the center of a cluster,
30        #this variable is set to its name
31        self.name = None

```

```

32     #This is the Transition Matrix of this chunk
33     self.TM = None
34
35     #This function is used by create_TM in order to initiate the Transition
36     #Matrix
37     def init_TM(self, num_msgs):
38         self.TM = [[0 for _ in range(num_msgs)] for _ in range(num_msgs)]
39
40     #This function creates the Transition Matrix of the chunk
41     #Given the message sequence and a message dictionary
42     #1)initialises the Transitions Matrix
43     #2)for each transition from message A to message B increases by 1 the
44     #respective cell
45     #3)normalises the matri by row
46     def create_TM(self, msg_sequence, msgs_dict):
47         self.init_TM(len(msgs_dict))
48         cur_msg = msg_sequence[0]
49
50         for i in range(len(msg_sequence)-1):
51             past_msg = cur_msg
52             cur_msg = msg_sequence[i+1]
53             self.TM[self.msgs_dict[past_msg]][self.msgs_dict[cur_msg]] += 1
54
55         for row_id, row in enumerate(self.TM, 0):
56             row_sum = sum(row)*1.0
57             self.TM[row_id] = [0.0 if row_sum == 0.0 else cell/row_sum for cell
58                               in row]

```

Listing 8.1: Υλοποίηση την Κλάσης Chunk

- **Δημιουργία της λίστας Chunk και εξαγωγή σε JSON**

Κάθε κομμάτι που βρέθηκε στην αρχική διαχωρισμένη λίστα δεδομένων, μετατράπηκε σε ένα αντικείμενο της κλάσης Chunk. Τα στοιχεία αυτά χωρίζονται σε δύο βασικά είδη, "Cleaning Chunks" και "Production Chunks". Όπως μπορεί κανείς να δει στην υλοποίηση της κλάσης Chunk, ορισμένα χαρακτηριστικά συμπληρώθηκαν κατά τη δημιουργία του κάθε αντικειμένου, και άλλα έμειναν κενά.

Ένα από τα πιο σημαντικά χαρακτηριστικά των Chunks είναι ο Πίνακας Μεταβάσεων, ο οποίος δημιουργήθηκε από την αλληλουχία των μηνυμάτων ενεργειών σε κάθε παρτίδα.

Όταν όλες οι παρτίδες επισημάνθηκαν και μετατράπηκαν σε αντικείμενα της κλάσης Chunk, αποθηκεύσαμε την λίστα αυτή, προκειμένου να έχουμε ένα σταθερό σύνολο δεδομένων για περαιτέρω πειραματισμό και να μην απαιτείται εκ νέου εκτέλεση του αλγόριθμου. Η έξοδος ήταν ένα αρχείο JSON, το οποίο αποδείχθηκε πολύ χρήσιμο στη συνέχεια.

8.3 Οπτικοποίηση των Δεδομένων

Στην αρχική συνομιλία με την Johnson & Johnson, είχε επικοινωνηθεί με σαφήνεια ότι μια κρίσιμη ανάγκη για εκείνους, ήταν η δυνατότητα να απεικονιστούν τα δεδομένα από την παραγωγική διαδικασία. Αυτή η εργασία ήταν πολύ δύσκολη και κάποιες φορές αδύνατη μέσω του αρχικού σετ δεδομένων, λόγω του γεγονότος ότι κάθε παρτίδα μπορεί να αποτελείται από έως και 800 σειρές δεδομένων. Για το σκοπό αυτό, δημιουργήσαμε ένα εργαλείο οπτικοποίησης, για να απεικονίσουμε γραφικά όλες τις μεταβλητές που παρέχονται. Αυτό το εργαλείο είναι ήδη σε χρήση και έχει αποδειχθεί πολύ χρήσιμο στην παρακολούθηση της παραγωγικής διαδικασίας.

8.3.1 Υποδομές και Μεθοδολογία

Η οπτικοποίηση των δεδομένων σε διαγράμματα έγινε με τη χρήση μιας βιβλιοθήκης javascript. Συνδυάσαμε διάφορα διαγράμματα για να δημιουργήσουμε και να προσαρμόσουμε το συγκεκριμένο εργαλείο.

Τα βήματα που ο χρήστης έχει να κάνει είναι η τοποθέτηση του αρχείου CSV σε κατάλληλο φάκελο και η έναρξη της εφαρμογής προ-επεξεργασίας, από τη αρχική σελίδα του εργαλείου, όπως φαίνεται στα επόμενα σχήματα 8.5. Η ιστοσελίδα αυτή δημιουργήθηκε με χρήση PHP και HTML. Από αυτό το βήμα και έπειτα, η διαδικασία είναι αυτοματοποιημένη και ο χρήστης πρέπει να περιμένει μόνο μερικά δευτερόλεπτα προτού να είναι σε θέση να δει το διάγραμμα.

Η προ-επεξεργασία αρχίζει με την ανάγνωση του αρχείου εισόδου CSV, το αρχείο καθαρίζεται από τις τιμές που λείπουν και μετασχηματίζεται σε μια πιο κατάλληλη μορφή για την οπτικοποίηση, όπως εξηγείται στο σημείο 8.2.1. Το προηγούμενο μέρος υλοποιήθηκε μόνο με χρήση Python. Στη συνέχεια, το αρχείο αποθηκεύεται σε έναν πίνακα της βάσης δεδομένων. Χρησιμοποιήθηκε το WAMP, το οποίο μας επέτρεψε να δημιουργήσουμε μια βάση δεδομένων MySQL.

Μόλις είναι όλα έτοιμα και αποθηκευμένα στη βάση δεδομένων, ο χρήστης μπορεί να μεταβεί στη σελίδα με το διάγραμμα. Αυτή η σελίδα είναι το τελικό αποτέλεσμα της οπτικοποίησης μας. Μπορούμε να δούμε περισσότερες λεπτομέρειες και εικόνες του εργαλείου στην επόμενη υποενότητα.

8.3.2 Επεξήγηση του Διαγράμματος

- **Περιγραφή**

Στόχος μας ήταν να απεικονισθούν όλες οι μεταβλητές σε ένα γράφημα, έτσι ώστε ο χρήστης να είναι σε θέση να δει όλες τις συσχετίσεις και τις εξαρτήσεις μεταξύ τους, να μπορεί να επιλέξει ποιες να απεικονισθούν και ποιες να κρύψει. Το διάγραμμα παρέχει τη δυνατότητα εμφάνισης ή απόκρυψης κάθε ξεχωριστής μεταβλητής, κάνοντας κλικ στο όνομά της στο κάτω μέρος της οθόνης, και την οποία μεταβλητή στη συνέχεια εμφανίζει ή αποκρύπτει τόσο στο γράφημα όσο και στον κάθετο άξονα. Επιπλέον, ο χρήστης μπορεί να επικεντρώσει το γράφημα στο συγκεκριμένο χρονικό πλαίσιο που τον ενδιαφέρει, είτε πληκτρολογώντας τις ακριβείς ημερομηνίες στα επάνω δεξιά κουτιά, ή επιλέγοντας μία από τις δυνατότητες στην πάνω αριστερή γωνία, είτε, τέλος, με την κίνηση της μπάρας στο κάτω μέρος του χάρτη. Προσθέσαμε επίσης μια επιπλέον μεταβλητή, σε μορφή σημαιών, η οποία τοποθετείται στο κάτω μέρος του διαγράμματος και, με το πέρασμα του δείκτη του ποντικιού, εμφανίζει το μήνυμα της ενέργειας που συνέβη εκείνη τη χρονική στιγμή. Μερικά στιγμιότυπα από το εργαλείο παρουσιάζονται παρακάτω.

- **Στιγμιότυπα**

Παρακάτω μπορείτε να βρείτε τα προαναφερθέντα στιγμιότυπα, τα οποία μπορούν να βοηθήσουν στην κατανόηση της λειτουργίας του εργαλείου.

Visualisation Tool

NTUA Thesis Project in cooperation with Johnson & Johnson Hellas

Developed by:

Haris Michailidis haris.michailidis@gmail.com

Isidora Tourni isidora.tourni@gmail.com

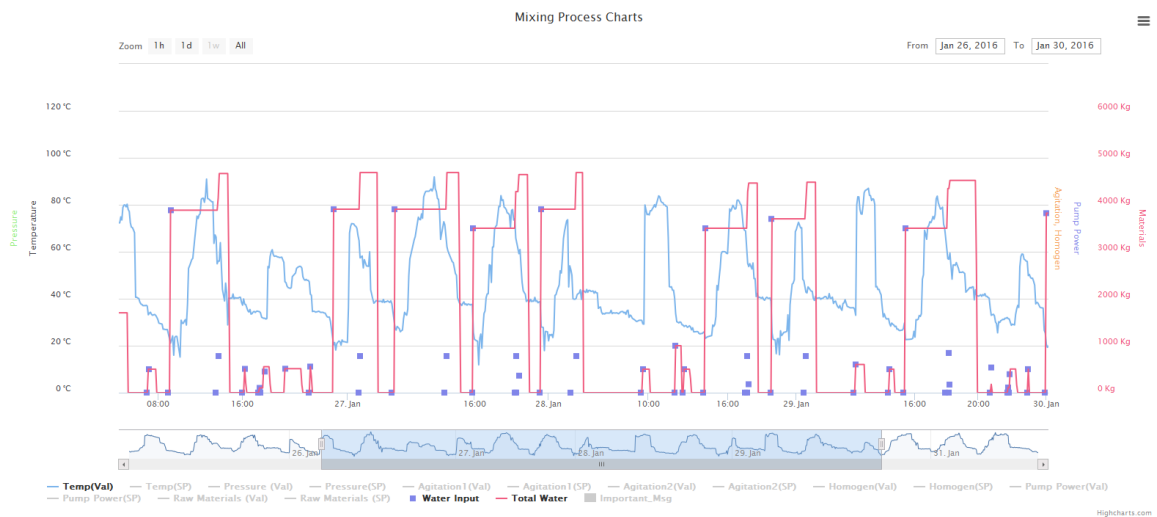
Instructions

- 1) Copy the proper CSV to the VT_input folder
- 2) Rename it to "input.csv"
- 3) Click on the "Run Main Script" button
- 4) Wait for ~1min
- 5) The Log and an extra button will appear
- 6) Click on the "Go to Chart" button

Run Main Script

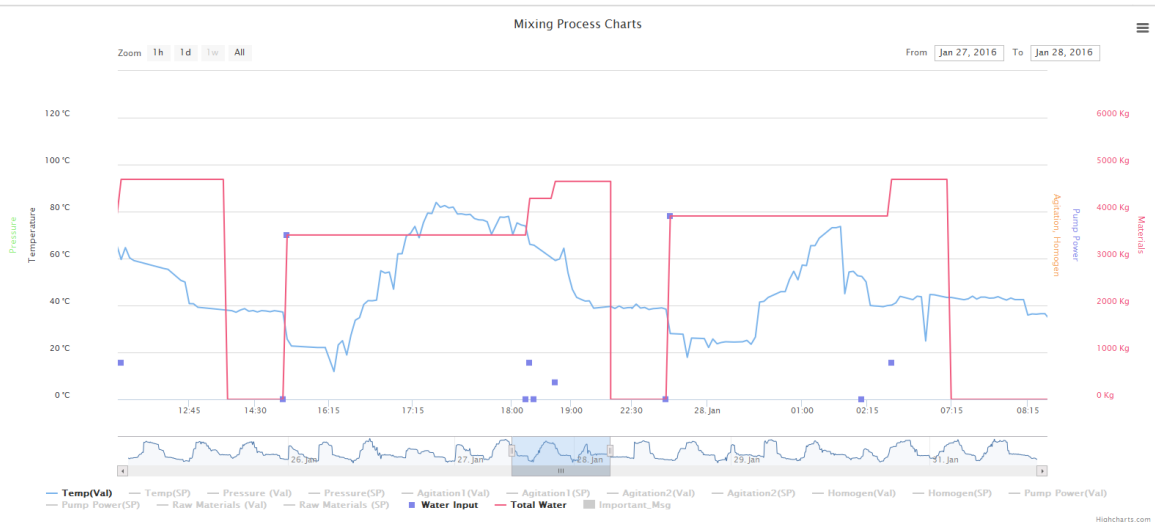
Σχήμα 8.5: Εργαλείο Οπτικοποίησης - Landing Page

Αυτό το στιγμιότυπο δείχνει την αρχική σελίδα του εργαλείου, όπου ο χρήστης μπορεί να διαβάσει τις οδηγίες. Όταν έχει ολοκληρωθεί η τοποθέτηση του αρχείου στο σωστό φάκελο, κάνει κλικ στο κουμπί "Run Main Script".



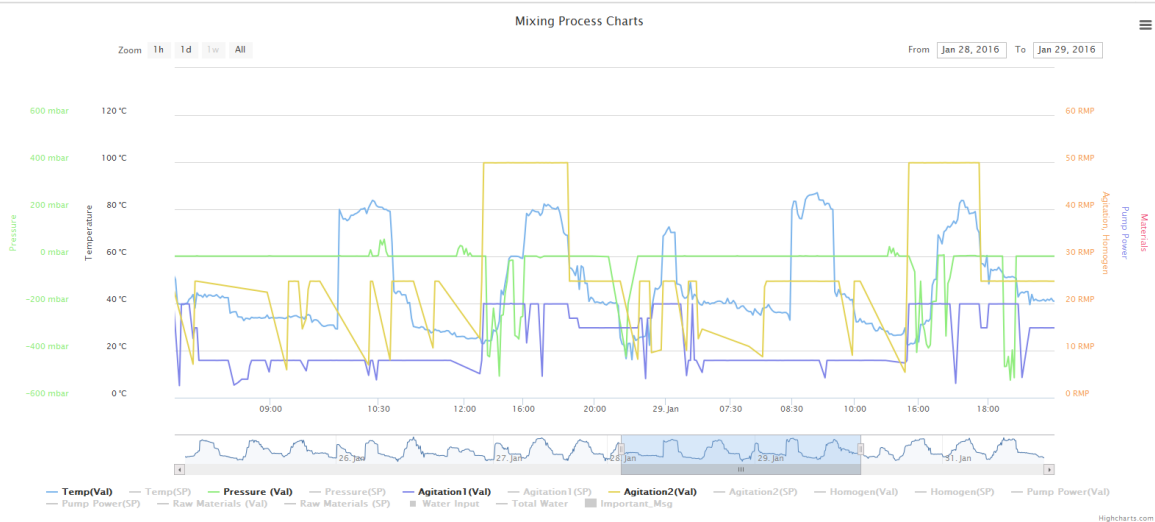
Σχήμα 8.6: Εργαλείο Οπτικοποίησης - Temperature, Water - 4 days

Σε αυτή την οθόνη, ο χρήστης ενδιαφέρεται για την παρακολούθηση της θερμοκρασίας και το συνολικό νερό που υπάρχει στο δοχείο. Το χρονικό παράθυρο είναι σχετικά μεγάλο, με επισκόπηση 4 ημερών.



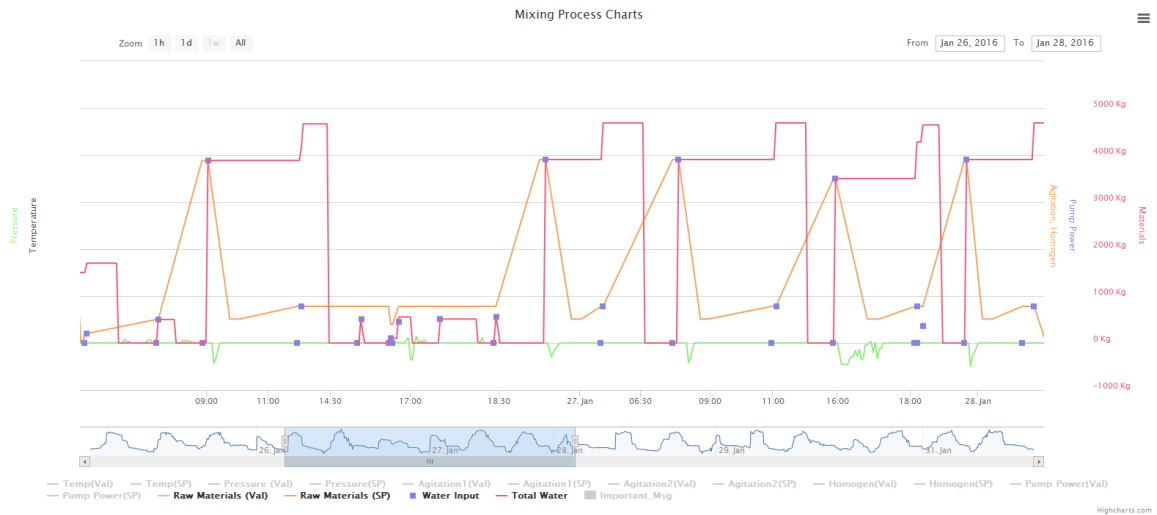
Σχήμα 8.7: Εργαλείο Οπτικοποίησης - Temperature, Water - 1 day

Αυτή είναι η ίδια οθόνη όπως στο προηγούμενο σχήμα, με εξαίρεση το χρονικό παράθυρο που, σε αυτή την περίπτωση, είναι πιο μικρό και εμφανίζει μόνο μία ημέρα της δραστηριότητας.



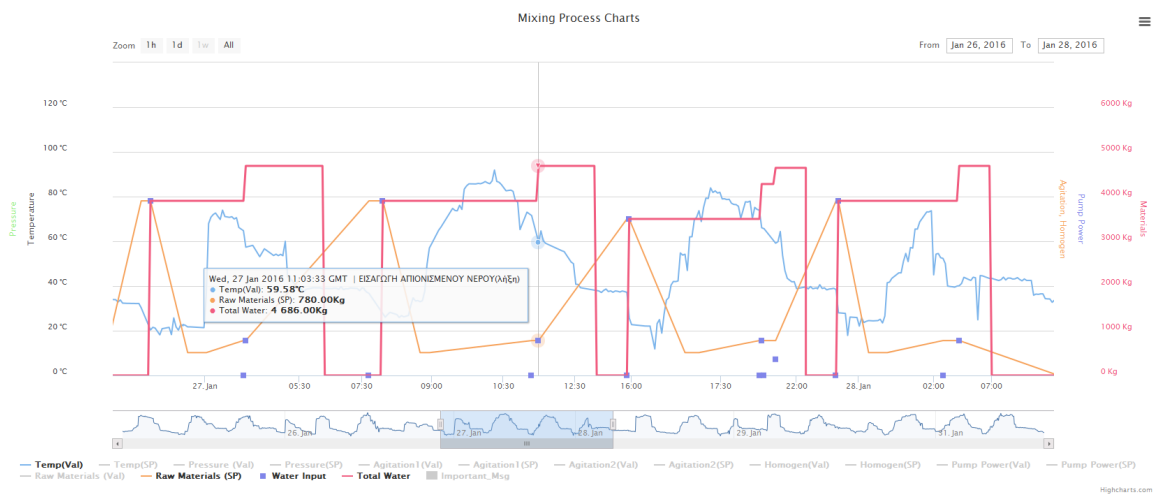
Σχήμα 8.8: Εργαλείο Οπτικοποίησης - Four Variables Correlation

Σε αυτή την οθόνη, ο χρήστης έχει ενεργοποιήσει τέσσερις ξεχωριστές γραφικές παραστάσεις προς προβολή, θερμοκρασία, πίεση ανάδευση 1 και ανάδευση 2. Μέσω αυτού του γραφήματος, ο χρήστης μπορεί να κατανοήσει καλύτερα τη σχέση μεταξύ των μεταβλητών κατά τη διάρκεια της διαδικασίας ανάμιξης.



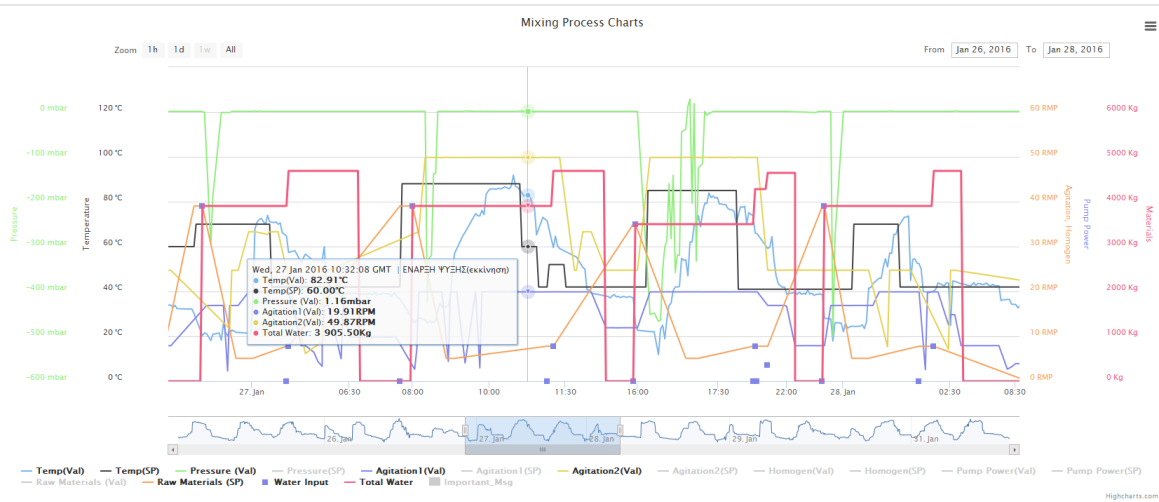
Σχήμα 8.9: Εργαλείο Οπτικοποίησης - Raw Materials Addition

Ένα από τα πιο σημαντικά στοιχεία της διαδικασίας ανάμιξης είναι η προσθήκη νερού. Σε αυτό το γράφημα μπορούμε να δούμε με μοβ χρώμα την απόλυτη τιμή του εισαχθέντος νερού σε κάθε στιγμή. Με την κόκκινη γραμμή παρατηρούμε το άθροισμα όλου του νερού που βρίσκεται στο δοχείο, μέχρι τη στιγμή που αποστραγγίζεται από αυτό, οπότε η κόκκινη γραμμή πηγαίνει στο μηδέν.



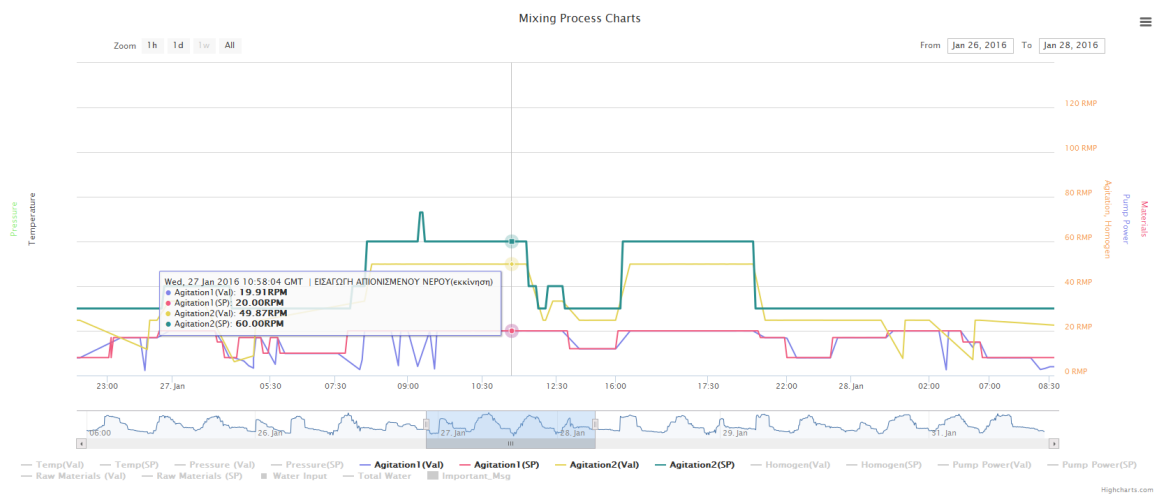
Σχήμα 8.10: Εργαλείο Οπτικοποίησης - Action Message Presentation

Προκειμένου ο χρήστης να είναι σε θέση να εξετάσει την διαδικασία λεπτομερώς, υπάρχει ένα κουτί, το οποίο περιλαμβάνει όλες τις τιμές από τις τρέχουσες μεταβλητές του γραφήματος σε οποιαδήποτε χρονική στιγμή. Ο χρήστης μπορεί να μετακινεί το ποντίκι πάνω στο γράφημα και να βλέπει τις ακριβείς τιμές της κάθε μεταβλητής, μαζί με το μήνυμα ενέργειας εκείνης της στιγμής.



Σχήμα 8.11: Εργαλείο Οπτικοποίησης - Variable Values Presentation

Όπως παρουσιάστηκε στο προηγούμενο στιγμιότυπο, στο πλαίσιο που εμφανίζεται στην άκρη των μετρήσεων, όλες οι ενεργές μεταβλητές μπορούν να εμφανισθούν, επιτρέποντας στο χρήστη να έχει μια πλήρη εικόνα οποιαδήποτε στιγμή της διαδικασίας.



Σχήμα 8.12: Εργαλείο Οπτικοποίησης - Values and Set Points

Σε αυτό το στιγμιότυπο, παρουσιάζεται η έννοια των set-points. Το set-point είναι μια τιμή, που καθορίζεται από το χρήστη, η οποία λειτουργεί ως οδηγός για τη διαδικασία. Όπως μπορεί κανείς να δει, οι τιμές SP είναι πάντα μπροστά από τις μετρούμενες.

Κεφάλαιο 9

Μηχανική Μάθηση

Σε αυτό το κεφάλαιο, θα αναλύσουμε το θεωρητικό υπόβαθρο όλων των τεχνικών Μηχανικής Μάθησης που υλοποιούνται, και τις παραμέτρους τους, και θα ορίσουμε τις έννοιες που είναι απαραίτητες για την κατανόηση των μεθόδων εφαρμογής.

Στην ενότητα 9.1, εισάγουμε τον αναγνώστη στους τομείς του Data Mining και της Μηχανικής Μάθησης, στις διάφορες κατηγορίες από τις οποίες αποτελούνται, και αυτές στις οποίες θα επικεντρωθεί η παρούσα εργασία.

Στην ενότητα 9.2, ορίζουμε περαιτέρω το Classification και υποκατηγορίες του, και αναλύουμε τις Instance-Based τεχνικές Μάθησης, και, ειδικότερα, τον k- Nearest Neighbours αλγόριθμο, το οποίο εφαρμόζουμε αργότερα στο σύνολο των δεδομένων μας.

Στην ενότητα 9.3, περιγράφουμε την έννοια του Clustering, απαριθμώντας τις ομάδες στις οποίες χωρίζεται, και αφιερώνοντας το υπόλοιπο της ενότητας στον k- Means αλγόριθμο, ως τη μέθοδο Clustering που θα εφαρμοστεί περαιτέρω.

Στην ενότητα 9.4, εξηγούμε το ρόλο της απόστασης όσον αφορά στις προηγουμένως καθορισμένες μεθόδους μάθησης, και αναλύουμε κάποιες σημαντικές μεθόδους μετρήσεις απόστασης.

Τέλος, στην ενότητα 9.5, εξετάζουμε κάποιες προσεγγίσεις και τεχνικές αξιολόγησης των αλγορίθμων που χρησιμοποιούνται στη Μηχανική Μάθηση, και θεωρητικά περιγράφουμε αυτές που θα χρησιμοποιηθούν περαιτέρω.

9.1 Εξόρυξη δεδομένων και Μηχανική Μάθηση

[11] [3] Η διαδικασία ανακάλυψης μοτίβων σε δεδομένα ορίζεται ευρέως ως Εξόρυξη Δεδομένων (Data Mining). Πρόκειται για μια αυτόματη ή ημιαυτόματη διαδικασία, που εφαρμόζεται σε σημαντικές ποσότητες δεδομένων, προκειμένου να ανακαλυφθούν συνδυασμοί και να αποκτήσουμε αξιόλογες πληροφορίες. Αυτή η εύρεση μοτίβων είναι εξαιρετικά χρήσιμη, καθώς μας επιτρέπει να πραγματοποιούμε με τετριμμένες προβλέψεις για τα νέα δεδομένα. [2]

[8] Όταν τα μοτίβα, τα οποία αναγνωρίζονται, αναπαρίστανται σε δομές, που μπορούν να εξεταστούν, να αναλυθούν και να χρησιμοποιηθούν για την ενημέρωση των μελλοντικών αποφάσεων, οι δομές αυτές ονομάζονται διαρθρωτικές. Μηχανική μάθηση ονομάζεται μια συλλογή από τεχνικές για την εύρεση και την περιγραφή αυτών των προτύπων διαρθρωτικών δομών, και αποτελεί ένα εργαλείο για την εξήγηση των δεδομένων και την πραγματοποίηση προβλέψεων από αυτά.

Τα δεδομένα λάμβάνουν στη συνέχεια τη μορφή ενός συνόλου παραδειγμάτων ή καταστάσεων, που χαρακτηρίζονται ως instances, και η έξοδος της διαδικασίας λαμβάνει τη μορφή πρόβλεψης και συνόλων κανόνων σχετικά με τη νέα παραδείγματα, υπό δεδομένες συνθήκες. Έτσι, η μάθηση μπορεί να θεωρηθεί ότι έχει δύο ξεχωριστούς ορισμούς: η δημιουργία γνώσης και η ικανότητα να χρησιμοποιηθεί για περαιτέρω σκοπούς.

Σε εφαρμογές εξόρυξης δεδομένων, οι διαδικασίες μάθησης μπορούν να ταξινομηθούν σε τέσσερις διαφορετικούς τομείς:

- **Classification Learning**, όπου το σύστημα εκμάθησης τροφοδοτείται με ένα σύνολο γνωστών παραδειγμάτων, από τα οποία αναμένεται να μάθει έναν τρόπο ταξινόμησης αυτών

- **Association Learning**, όπου αναζητείται κάθε συσχέτιση μεταξύ χαρακτηριστικών, όχι μόνο αυτών που προβλέπουν μια συγκεκριμένη class
- **Clustering**, όπου πραγματοποιείται η ομαδοποίηση παραδειγμάτων που ανήκουν μαζί
- **Numeric Prediction**, όπου το αποτέλεσμα που ζητείται να προβλεφθεί δεν είναι μια διακριτή class, αλλά μια αριθμητική ποσότητα

Στην παρούσα εργασία, δύο από τις παραπάνω μεθόδους θα εξεταστούν λεπτομερώς και θα εφαρμοστούν στα δεδομένα μας, τα Classification και Clustering, τα οποία αποτελούν πιο ειδικούς ορισμούς για τις δύο μεγαλύτερες κατηγορίες Μηχανικής Μάθησης: Επιβλεπόμενη και Μη Επιβλεπόμενη Μάθηση.

9.2 Classification

Το Classification καλύπτει ένα ευρύ φάσμα της ανθρώπινης δραστηριότητας. Στην ευρύτερη έννοιά του, ο όρος περιλαμβάνει οποιαδήποτε απόφαση ή πρόβλεψη έγινε με βάση τις διαθέσιμες πληροφορίες, και η διαδικασία classification ορίζεται ως μέθοδος για την κατ'επανάληψη πραγματοποίηση των εν λόγω αποφάσεων σε επερχόμενες καταστάσεις. Πιο συγκεκριμένα, το πρόβλημα αφορά σε μια διαδικασία που θα εφαρμόζεται σε μία συνεχή αλληλουχία από δεδομένα, κάθENA από τα οποία χρειάζεται να ανατεθεί σε μια προκαθορισμένη class, αναλόγως με χαρακτηριστικά ή ιδιότητες. [1]

Το Classification ονομάζεται και Επιβλεπόμενη Μάθηση, επειδή, κατά μία έννοια, το σύστημα λειτουργεί υπό εποπτία, καθώς διαθέτουμε και το πραγματικό αποτέλεσμα για κάθENA από τα παραδείγματα training. Αυτό το αποτέλεσμα ορίζεται ως class.

Μερικά από τα πιο σημαντικά προβλήματα της επιστήμης, της βιομηχανίας και του εμπορίου, που απαιτούν περίπλοκα και συχνά πολλά δεδομένα, μπορούν να θεωρηθούν ως Classification προβλήματα, όπως η προκαταρκτική διάγνωση της νόσου του ασθενούς, ενώ αναμένουμε τα αποτελέσματα, προκειμένου να επιλέγει η άμεση θεραπεία, ή η ταξινόμηση των ατόμων βάσει των οικονομικών και προσωπικών τους πληροφοριών.

[5] Ένας μεγάλος αριθμός από τεχνικές Classification έχει αναπτυχθεί. Μπορούμε ονομαστικά να τις χωρίσουμε στους ακόλουθους τομείς:

- **Logical and Symbolic techniques**, όπως τα Decision Trees και τα Learning Rulesets
- **Perception-based techniques**, που αναλύονται σε Single Layer Perceptrons, Multilayered Perceptrons και Radial Basis Function (RBF) Networks
- **Statistics**, που περιλαμβάνουν Naive Bayes Classifiers και τα Bayesian Networks
- **Instance - Based Learning**
- **Support Vector Machines**

Στην επόμενη ενότητα θα επικεντρωθούμε στις Instance- Based τεχνικές, μερικές από τις οποίες εφαρμόστηκαν στην παρούσα μελέτη.

9.2.1 Instance - Based Learning

[9] Οι Instance - Based Learning αλγόριθμοι μάθησης καθυστερούν τη διαδικασία γενίκευσης, έως ότου γίνει η διαδικασία classification. Μόλις ένα σύνολο instances έχει απομνημονευθεί, για ένα νέο instance η μνήμη αναζητά το υπάρχον εκείνο που του μοιάζει περισσότερο. Με άλλα λόγια, τα γνωστά instances αποθηκεύονται και τα νέα, των οποίων η class είναι άγνωστη, συσχετίζονται με αυτά. Επομένως, όλη η πραγματική κατηγοριοποίηση γίνεται όταν χρειάζεται να ταξινομήσουμε ένα

νέο instance, και όχι όταν επεξεργαζόμαστε τα training δεδομένα, και έτσι οι αλγόριθμοι απαιτούν λιγότερο χρόνο κατά τον υπολογισμό παρά κατά το classification.

Η παραπάνω διαδικασία ονομάζεται k-Nearest Neighbours Classification. Η απόλυτη θέση των instances μέσα σε αυτό το χώρο δεν είναι τόσο σημαντική, σε αντίθεση με τη σχετική απόσταση μεταξύ τους. Χρησιμοποιώντας μια κατάλληλη μέθοδο υπολογισμού απόστασης, η οποία ελαχιστοποιεί ιδανικά η απόσταση μεταξύ δύο ομοίως κατηγοριοποιημένων instances, μεγιστοποιώντας παράλληλα την απόσταση μεταξύ των instances των διαφορετικών classes, το κοντινότερο instance χρησιμοποιείται για να αναθέσουμε στο νέο σε αυτή την class. Μερικές φορές, περισσότερα από ένα του πλησιέστερου γειτονικού στοιχείου εξετάζονται, και η πλειοψηφούσα class ανατίθεται στο νέο instance. Αυτή είναι η μέθοδος k-Nearest Neighbours. Η επιλογή του k επηρεάζει την απόδοση του αλγορίθμου και το αποτέλεσμα της ταξινόμησης.

Ένα παράδειγμα ψευδοκώδικα για τον αλγόριθμο αυτό παρουσιάζεται ακολούθως: X : Training Data, Y : Class Labels of X , x : Unknown sample

```
k Nearest Neighbours( $X, Y, x$ )  
for  $i=1$  to  $m$  do  
    Compute Distance  $d(X_i, x)$   
end for  
compute set  $I$  containing indices for the  $k$  smallest distances  $d(X_i, x)$ .  
return majority label for  $Y_i$  where  $i \in I$ 
```

9.3 Clustering

Οι τεχνικές Clustering εφαρμόζονται όταν τα instances χρειάζεται να διαιρεθούν σε ομάδες χωρίς να υπάρχει κάποια class για να προβλεφθεί. Αυτές οι ομάδες (clusters) στηρίζονται στο ότι κάποια instances μοιάζουν περισσότερο με κάποια απ'όσο με κάποια άλλα. Θεωρούνται άγνωστες και προκύπτουν από τα δεδομένα αυτά καθαυτά, γι αυτό και το Clustering είναι γνωστό σαν Μη Επιβλεπόμενη Μάθηση.

Τα clusters που αναγνωρίζονται μπορούν να ανήκουν σε μία από τις ακόλουθες κατηγορίες, με βάση τη φύση των μηχανισμών που θεωρούμε ότι περιγράφουν το εκάστοτε φαινόμενο ομαδοποίησης:

- **Exclusive**, δηλαδή κάθε instance ανήκει μόνο σε ένα cluster
- **Overlapping**, όπου κάθε instance μπορεί να ανήκει σε πολλά clusters
- **Probabilistic**, γιατί ένα instance μπορεί να ανήκει σε κάθε cluster με μια συγκεκριμένη πιθανότητα
- **Hierarchical**, δηλαδή ένας αρχικός διαχωρισμός των instances σε ομάδες στο κορυφαίο επίπεδο, και επανεξέταση του κάθε cluster, μέχρι τα μεμονωμένα instances

Ωστόσο, επειδή οι μηχανισμοί αυτοί είναι σπάνια γνωστοί, ο χαρακτηρισμός τους συνήθως υπαγορεύεται από τα εργαλεία clustering που έχουμε στη διάθεσή μας.

Μερικά πραγματικά παραδείγματα της χρήσης της ομαδοποίησης είναι τη διαίρεση των πελατών σε ομοιογενείς ομάδες ως τεχνική μάρκετινγκ, η συλλογή δεδομένων του καιρού και η ανάλυσή τους για την εύρεση νέων θεωριών στα κλιματολογικά και περιβαλλοντικά πεδία, και η ταυτοποίηση, μέσω της βιοπληροφορικής, ομάδων γονιδίων με σκοπό να καθοριστεί ποια γονίδια είναι υπεύθυνα για συγκεκριμένες κληρονομικές ασθένειες.

Λόγω του γεγονότος ότι η έννοια της συστάδας δεν είναι επακριβώς καθορισμένη, πολλές μέθοδοι έχουν αναπτυχθεί στην ομαδοποίηση, καθεμία από τις οποίες χρησιμοποιεί μια διαφορετική αρχή επαγωγής. Διαιρούνται σε πέντε βασικές ομάδες, όπως αναλύεται παρακάτω:

- **Hierarchical Methods**, οι οποίες κατασκευάζουν τα clusters διαιρώντας αναδρομικά τα δεδομένα με τρόπο top-down ή bottom-up. Μπορούν να διαχωρισθούν περαιτέρω σε Agglomerative hierarchical clustering και Divisive hierarchical clustering
- **Partitioning Methods**, οι οποίες μεταφέρουν instances από ένα cluster σε άλλο, μετά από μια αρχική διαμέρισή τους. Περιλαμβάνουν τα Error Minimization Algorithms και Graph - Theoretic Clustering. Ο απλούστερος αλγόριθμος, ο οποίος αξιοποιεί τη μέθοδο του squared error, είναι ο K-means algorithm, ο οποίος θα αναλυθεί στη συνέχεια καθώς χρησιμοποιήθηκε στο παρόν πρόβλημα .
- **Density - Based Methods**, οι οποίες υποθέτουν ότι τα στοιχεία που ανήκουν στο ίδιο cluster προκύπτουν από μια συγκεκριμένη πιθανοτική κατανομή.
- **Model - Based Clustering Methods**, που επιχειρούν να βελτιστοποιήσουν το ταίριασμα μεταξύ των δεδομένων και κάποιων μαθηματικών μοντέλων. Οι πιο συχνά χρησιμοποιούμενες μεθόδους στην κατηγορία αυτή είναι τα Decision trees και τα Neural Networks.
- **Grid - Based Methods**, που διαμερίζουν το χώρο σε ένα πεπερασμένο αριθμό κελιών πίνακα, στον οποίο εφαρμόζονται και όλες οι τεχνικές Clustering.

9.3.1 Partitioning Methods- K - Means algorithm

Οι Partitioning clustering τεχνικές, όπως αναφέρθηκε, διαμερίζουν τα δεδομένα σε ομάδες όπου κάθε ζεύγος αντικειμένων του cluster είναι είτε διακριτό (hard clustering) είτε έχει κάποια κοινά χαρακτηριστικά (soft clustering).

Η κλασική και πιο κοινή τεχνική ονομάζεται k-Means. Εφαρμόζοντας αυτό τον αλγόριθμο, αρχικά ορίζουμε την παράμετρο k, η οποία αναπαριστά τον αριθμό των clusters που αναζητούμε. Έπειτα, k σημεία επιλέγονται τυχαία σαν κέντρα των clusters. Ορίζουμε επίσης ένα μέγιστο αριθμό επαναλήψεων, οπότε και η διαδικασία τερματίζεται αν αυτός ξεπεραστεί. Κάθε instance τοποθετείται στο κοντινότερο κέντρο κάποιου cluster, με βάση την Euclidean μέθοδο υπολογισμού απόστασης. Ύστερα, υπολογίζεται το κέντρο, ή ο "μέσος" κάθε cluster, και θεωρείται το νέο του κέντρο. Τέλος, η παραπάνω διαδικασία επαναλαμβάνεται για τα νέα αυτά κέντρα. Οι επαναλήψεις συνεχίζονται μέχρι τα ίδια σημεία να κατηγοριοποιηθούν σε κάποιο cluster διαδοχικά, οπότε και τα κέντρα των clusters σταθεροποιούνται και παραμένουν τα ίδια για πάντα. Αν ο μέγιστος αριθμός επαναλήψεων έχει ξεπεραστεί, εν τω μεταξύ, ο αλγόριθμος τερματίζεται.

Ο ψευδοκώδικας για την παραπάνω διαδικασία παρουσιάζεται στη συνέχεια: S : Instance set, k : Number of Clusters

k - Means(S, k)

Initialize k cluster centers

while termination condition is not satisfied **do**

 assign instances to the closest cluster center

 update cluster centers based on the assignment

end while

return clusters

Αυτή η μέθοδος clustering είναι ιδιαίτερα απλή και αποτελεσματική. Είναι εύκολο να αποδείξουμε ότι η επιλογή του κέντρου κάθε φορά ελαχιστοποιεί την συνολική απόσταση κάθε μέλους του cluster από αυτό. Χρειάζεται βέβαια να λάβουμε υπόψη μας ότι το ελάχιστο που υπολογίζεται είναι τοπικό και όχι ολικό. Για να αυξήσουμε την πιθανότητα να εντοπίσουμε ένα ολικό ελάχιστο, συχνά εκτελούμε τον αλγόριθμο πολλαπλές φορές, με διαφορετικές αρχικές επιλογές, ώστε να καταλήξουμε στα βέλτιστα αποτελέσματα.

9.4 Αποστάσεις

[7] Η έννοια της απόστασης είναι η πιο σημαντική βάση για τη Μηχανική Μάθηση, τόσο την Επιβλεπόμενη όσο και τη Μη Επιβλεπόμενη. Για την πρώτη κατηγορία, οι συνήθεις αποστάσεις πολλές φορές οδηγούν σε ακατάλληλα αποτελέσματα, ενώ στη δεύτερη ο υπολογισμός των μέσων των αντικειμένων γνωστών ομάδων δεν είναι πάντα έγκυρη μέθοδος για τη σωστή εφαρμογή του αλγορίθμου. Εξ ορισμού, η επιλογή της μεθόδου απόστασης καθορίζει κατά πόσο δύο αντικείμενα ομαδοποιούνται μεταξύ τους, επομένως, η ορθότητά της είναι από τα πλέον αποφασιστικά βήματα για την αποδοτικότητα της Μάθησης. Η απόσταση χρειάζεται όχι μόνο να μπορεί αναπαριστά τη μορφή των δεδομένων, αλλά και να μελετά το στόχο, ώστε να λάβουμε κατανοητά αποτελέσματα. Μερικές από τις πλέον χρησιμοποιούμενες μεθόδους υπολογισμού της απόστασης, οι οποίες υλοποιήθηκαν και από εμάς, αναλύονται ακολούθως. [6]

9.4.1 Minkowski Metric

Η απόσταση Minkowski ή η L_q νόρμα υπολογίζει την απόσταση d μεταξύ δύο αντικειμένων x και y συγκρίνοντας τις τιμές n χαρακτηριστικών τους. Δίνεται από την εξίσωση 9.1, και εφαρμόζεται σε τιμές συχνότητας, πιθανοτήτων και δυαδικές.

$$d(x, y) = L_q(x, y) = \sqrt[q]{\sum_{i=1}^n |x_i - y_i|^q} \quad (9.1)$$

Η πιο σημαντική ειδική περίπτωση αυτής προκύπτει για $q=2$ και είναι η **Euclidean distance** ή L_2 νόρμα:

$$d(x, y) = L(x, y) = \sqrt{\sum_{i=1}^n |x_i - y_i|^2} \quad (9.2)$$

9.4.2 Cosine Distance

Η cosine similarity ορίζεται στο R^n , από την

$$\cos(a) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \cdot \|\mathbf{y}\|} \quad (9.3)$$

όπου a η γωνία μεταξύ των διανυσμάτων x και y . Η cosine distance ορίζεται ως

$$d(x, y) = 1 - \cos a$$

9.4.3 Kullback- Leibler Divergence

Η Kullback-Leibler divergence (KL) αποφασίζει την μη συμφωνία μιας μοντελοποιημένης κατανομής δεδομένης της πραγματικής κατανομής. Γενικά χρησιμοποιείται για δύο συναρτήσεις πιθανοτικών κατανομών, x και y , και δίνεται από την εξίσωση

$$d(x, y) = D(x \parallel y) = \sum_{i=1}^n x_i \star \log \frac{x_i}{y_i} \quad (9.4)$$

9.4.4 Kolmogorov- Smirnov Test

Η Kolmogorov Smirnov μέθοδος ορίζεται στο χώρο πιθανοτήτων P , από την

$$d(x, y) = \sup_{(x,y) \in R} |x - y| \quad (9.5)$$

θεωρώντας και πάλι ότι οι x και y είναι διαφορετικές συναρτήσεις κατανομών. Χρησιμοποιείται στη στατιστική ανάλυση σαν μέθοδος μέτρησης της "ποιότητας της αντιστοιχίας".

9.4.5 Infinite Norm

Η Infinite Norm είναι η $L_{(\infty)}$ νόρμα στο $C_{[a,b]}$ όλων των πραγματικών ή φανταστικών συνεχών συναρτήσεων σε ένα δεδομένο διάστημα $[a,b]$. Για διανύσματα x, y ορίζεται από

$$d(x, y) = \lim_{p \rightarrow \infty} \sqrt[p]{\sum_{i=1}^n |x_i - y_i|^p} \quad (9.6)$$

9.5 Αξιολόγηση

Η αξιολόγηση είναι μια σημαντική έννοια στο Data Mining, καθώς τα αποτελέσματα της εφαρμογής των μεθόδων μας είναι το κλειδί για περαιτέρω συμπεράσματα, βελτιστοποίηση και γενικότερη πρόοδο. Στην ενότητα αυτή, θα την αναλύσουμε τόσο για Classification όσο και για Clustering.

9.5.1 Αξιολόγηση Classification

Το ποσοστό επιτυχίας του classification συνήθως κρίνεται από τα δεδομένα test, για τα οποία οι πραγματικές ετικέτες είναι γνωστές, και αξιολογείται με διάφορες μεθόδους, δίνοντάς μας ένα αντικειμενικό μέγεθος του κατά πόσο η έννοια αφομοιώθηκε από τα δεδομένα. Εδώ, αναλύονται δύο διαφορετικές μέθοδοι, οι Accuracy και Cohen's Kappa, οι οποίες και εφαρμόζονται στο κεφάλαιο 11 στα αποτελέσματα των πειραμάτων.

Ένα σημαντικό επίσης θέμα, αναφορικά με την αξιολόγηση του αλγορίθμου, είναι η ταχύτητά του. Ένας αλγόριθμος classification, ο οποίος είναι 90% ακριβής, μπορεί να προτιμηθεί έναντι ενός 95% ακριβή, αν ο πρώτος είναι 100 φορές πιο γρήγορος κατά την εκτέλεση του testing (και τέτοια αποτελέσματα είναι συνηθισμένα στα νευρωνικά δίκτυα, για παράδειγμα). Η παράμετρος αυτή εξετάζεται και στην εργασία μας.

Accuracy

Η αξιοπιστία του αποτελέσματος Classification αντιπροσωπεύεται από το ποσοστό των σωστών ταξινομήσεων.

Συνήθως είναι η ακρίβεια των άγνωστων δεδομένων, όταν οι πραγματικές ετικέτες είναι άγνωστες, η οποία είναι πρακτικής σημασίας. Η πλέον αποδεκτή μέθοδος για αυτό τον υπολογισμό είναι να χρησιμοποιήσουμε τα δεδομένα, θεωρώντας ότι όλες οι αρχικές ετικέτες είναι γνωστές. Πρώτα, εκπαιδεύουμε τον αλγόριθμο με βάση το training set, δηλαδή ένα κομμάτι των αρχικών δεδομένων. Έπειτα τον εφαρμόζουμε και στα εναπομείναντα δεδομένα, και τα αποτελέσματα συγκρίνονται με τις γνωστές ετικέτες. Το ποσοστό των σωστά κατηγοριοποιημένων στο test set είναι μια εκτίμηση της ακρίβειας, γνωρίζοντας ότι τα training δεδομένα είναι δείγμα των αρχικών.

Υπάρχει μια μικρή περίπτωση ανακρίβειας εδώ, καθώς δε χρησιμοποιούμε το πλήρες δείγμα για να εκπαιδεύσουμε τον αλγόριθμο, αλλά με μεγάλες ομάδες δεδομένων αυτό δεν αποτελεί σημαντικό πρόβλημα.

Cohen's Kappa (Kappa coefficient)

[12] Ο δείκτης αυτός συγκρίνει την παρατηρούμενη ακρίβεια με την αναμενόμενη, και χρησιμοποιείται για την αξιολόγηση ενός αλγορίθμου classification ή για τη σύγκριση μεταξύ πολλών αλγορίθμων. Λαμβάνει υπόψη του την τυχαιότητα συμφωνίας με έναν τυχαίο classifier, κάτι που

	1	2	Total
1	p_{11}	p_{12}	p_{1r}
2	p_{21}	p_{22}	p_{2r}
Total	p_{1c}	p_{2c}	1

Πίνακας 9.1: Confusion Matrix after Classification

σημαίνει ότι είναι λιγότερο παραπλανητικός από την απλή ακρίβεια. Ο υπολογισμός της παρατηρούμενης ακρίβειας και της αναμενόμενης ακρίβειας είναι αναγκαίος για την μέθοδο, και εύκολα γίνεται κατανοητός με τη χρήση ενός confusion πίνακα.

Για ένα διδιάστατο πίνακα, όπως ο προηγούμενος, η Kappa τιμή είναι ίση με

$$k = \frac{p_0 - p_e}{1 - p_e} \quad (9.7)$$

όπου η παρατηρούμενη ακρίβεια είναι

$$p_0 = p_{11} + p_{22}$$

και η αναμενόμενη ακρίβεια είναι

$$p_e = p_{1c}p_{1r} + p_{2c}p_{2r}$$

Το Kappa είναι πάντα μικρότερο ή ίσο με 1. Η τιμή 1 σημαίνει τέλεια συμφωνία, και τιμές λιγότερο από 1 κάτι λιγότερο από τέλεια συμφωνία. Σε σπάνιες περιπτώσεις, το Kappa μπορεί να είναι αρνητικό. Αυτό είναι ένα δείγμα ότι οι δύο παρατηρητές συμφώνησαν λιγότερο από ό,τι θα αναμενόταν κατά τύχη.

Οι τιμές και σύμφωνες ερμηνείες του Kappa συνοψίζονται παρακάτω:

- Poor agreement = Λιγότερο από 0,20
- Fair agreement = 0,20 έως 0,40
- Moderate agreement = 0,40 έως 0,60
- Good agreement = 0,60 έως 0,80
- Very good agreement = 0,80 έως 1,00

9.5.2 Αξιολόγηση Clustering

Η αξιολόγηση του Clustering απαιτεί μια ανεξάρτητη και αξιόπιστη μέθοδο για τη σύγκριση των πειραμάτων και των αποτελεσμάτων. Στη θεωρία, ο ερευνητής έχει αποκτήσει μια διαίσθηση για το αποτέλεσμα, αλλά στην πράξη ο όγκος των δεδομένων, από τη μία πλευρά, και οι λεπτομέρειες της αναπαράστασης δεδομένων και αλγορίθμων clustering, από την άλλη, καθιστούν αυτή την τυχαία κρίση αβάσιμη. Μια διαισθητική, αξιολόγηση μπορεί να είναι μόνο εύλογη για μικρά σύνολα αντικειμένων, ενώ μεγάλης κλίμακας πειράματα απαιτούν μια αντικειμενική μέθοδο. Δεν υπάρχει μία απόλυτη οδός για την επιθυμητή εκτίμηση, αλλά μια ποικιλία τεχνικών αξιολόγησης που προέρχονται από διάφορα πεδία, όπως τη στατιστική και την Τεχνητή Νοημοσύνη.

Ο θεωρητικός ορισμός των διαφόρων μεθόδων αξιολόγησης, οι οποίες χρησιμοποιήθηκαν κατά τη διάρκεια των πειραμάτων της εργασίας, παρουσιάζεται στη συνέχεια.

V-measure

[4] Η V-measure είναι μια μέθοδος, η οποία μετρά πόσο ικανοποιούνται τα κριτήρια της ομοιογένειας και της πληρότητας. Υπολογίζεται ως ο αρμονικός μέσος των σκορ των δύο αυτών μεγεθών.

Το αποτέλεσμα του Clustering διαθέτει ομοιογένεια αν όλα τα clusters του περιέχουν μόνο δεδομένα, τα οποία είναι μέλη της ίδιας class. Επίσης, διαθέτει πληρότητα και αν όλα τα δεδομένα, τα οποία είναι μέλη μιας συγκεκριμένης class, είναι στοιχεία του ίδιου cluster. Η ομοιογένεια και η πληρότητα μιας λύσης είναι έννοιες αντίθετες: Η αύξηση της ομοιογένειας του cluster συχνά οδηγεί σε μείωση της πληρότητας του. Πιο συγκεκριμένα:

$$V_b = \frac{(1 + b) * h * c}{(b * h) * c}$$

Οι υπολογισμοί της ομοιογένειας, πληρότητας και V-measure είναι εντελώς ανεξάρτητοι του αριθμού των classes, του αριθμού των clusters, του μεγέθους του συνόλου δεδομένων και το αλγορίθμου. Έτσι, οι μέθοδοι αυτοί μπορούν να εφαρμοστούν σε οποιαδήποτε λύση clustering.

Rand Index

Ο Rand Index είναι ένα απλό κριτήριο που χρησιμοποιείται για να συγκρίνουμε μια επαγόμενη δομή clustering (C_1) με μια δεδομένη δομή clustering (C_2). Θεωρούμε

- **a** τον αριθμό των ζευγών των instances που τίθενται στο ίδιο cluster του C_1 και στο ίδιο cluster του C_2
- **b** τον αριθμό των ζευγών των instances που τίθενται στο ίδιο cluster του C_1 , αλλά όχι στο ίδιο cluster του C_2
- **c** τον αριθμό των ζευγών των instances που τίθενται στο ίδιο cluster του C_2 , αλλά όχι στο ίδιο cluster του C_1
- **d** τον αριθμό των ζευγών των instances που τίθενται σε διαφορετικά clusters του C_1 και C_2 .

Το ποσότητες a και d μπορούν να ερμηνευθούν σα μεταβλητές "συμφωνίας", και οι b και c ως "διαφωνίας". Ο Rand Index ορίζεται ως:

$$RAND = \frac{a + d}{a + b + c + d}$$

Ο Rand Index βρίσκεται μεταξύ 0 και 1. Όταν οι δύο κατανομές συμφωνούν πλήρως, τότε ισούται με 1.

Κεφάλαιο 10

Υλοποίηση

Σε αυτό το κεφάλαιο παρουσιάζουμε τις υλοποιήσεις όλων των εννοιών που αναλύθηκαν στο Κεφάλαιο 9 σχετικά με τα πειράματά μας, δηλαδή όλους τους αλγόριθμους και τις παραμέτρους τους.

Στην Ενότητα 10.1, παρουσιάζουμε τους δύο αλγορίθμους που χρησιμοποιήθηκαν για Classification, τον Nearest Centroid Classifier και τον k-Nearest Neighbours Classifier, και τις σημαντικότερες μεταβλητές για την εκτέλεσή τους.

Στην Ενότητα 10.2, επικεντρωνόμαστε στην διαδικασία Clustering, δίνοντας μια αναλυτική περιγραφή του αλγορίθμου k-Means που χρησιμοποιήθηκε για αυτόν τον σκοπό.

Στην Ενότητα 10.3, περιγράφουμε τις διαφορετικές μεθόδους υπολογισμού απόστασης για όλα τα πειράματα, και παρουσιάζουμε όλους τους τρόπους που υλοποιήθηκαν και αξιολογήθηκαν και στο Classification και στο Clustering.

Τέλος, στην Ενότητα 10.4, εξηγούμε με λεπτομέρεια τον τρόπο με τον οποίο αξιολογήθηκαν τα αποτελέσματα και τις μεθόδους που χρησιμοποιήθηκαν για αυτόν τον σκοπό.

10.1 Classification

Σκοπός της προσέγγισής μας ήταν να κατηγοριοποιήσουμε αντικείμενα chunk σύμφωνα με κάποιες από τις ιδιότητές τους. Συγκεκριμένα, θέλαμε να τα κατηγοριοποιήσουμε σύμφωνα με δύο διαφορετικές παραμέτρους, τα Product Cleaning Group, και Product Group. Και οι δύο είναι ιδιότητες κάθε αντικειμένου. Το Product Cleaning Group έχει 5 δυνατές περιπτώσεις και το Product Group έχει 12.

Επιπρόσθετα πειραματιστήκαμε, στους αλγορίθμους Classification, και με την μεταβλητή Product Code, όμως ο μεγάλος αριθμός διαφορετικών κωδικών προϊόντων και άρα ο μικρός αριθμός στοιχείων που αντιστοιχούσε σε κάθε κλάση, αποτέλεσαν στο overfitting του αλγορίθμου μας. Δεν ήταν δυνατό να εκπαιδύσουμε σωστά τον αλγόριθμο πάνω σε αυτά τα δεδομένα, ούτε να μετρήσουμε τα αποτελέσματα, και το αποτέλεσμα της διαδικασίας δεν μας έδωσε περισσότερες πληροφορίες από τις υπάρχουσες, σχετικά με την κατηγορία του προϊόντος.

Για το Classification, δύο διαφορετικοί αλγόριθμοι χρησιμοποιήθηκαν, όπως αναλύσαμε προηγουμένως: μια υλοποίηση του Nearest Centroid και ο γνωστός, καλά τεκμηριωμένος αλγόριθμος k-Nearest Neighbours.

Και οι δύο αλγόριθμοι είναι παραμετροποιημένοι ακολούθως:

- **Ιδιότητα Classifying:** Υλοποιήθηκαν για δύο ιδιότητες, Product Cleaning Group και Product Group,
- **Ποσοστό διαχωρισμού:** Διαχωρίσαμε τα αρχικά δεδομένα παραγωγής σε δύο σύνολα, ένα για εκπαίδευση, μέσω του οποίου υπολογίσαμε τα διαφορετικά κέντρα των κλάσεων και εφαρμόσαμε τους αλγορίθμους μας για να παράξουμε τα αποτελέσματα, και στο δοκιμαστικό μέρος, στο οποίο δοκιμάσαμε τους αλγορίθμους μας σύμφωνα με τα μοντέλα που προέκυψαν από το μέρος της εκπαίδευσης. Τα ποσοστά εκπαίδευσης και ελέγχου είναι 80% -20%, 65% -35% and 50% -50% αντίστοιχα και για τους δύο αλγορίθμους.

- **Αποστάσεις:** Χρησιμοποιήσαμε τις συναρτήσεις απόστασης, που έχουμε ορίσει, για να υπολογίσουμε την απόσταση μεταξύ των πινάκων μετάβασης του εκάστοτε αντικειμένου chunk, είτε από το κέντρο της κλάσης είτε από το κοντινότερο chunk, αναλόγως με την υλοποίηση.

10.1.1 Nearest Centroid Classifier

Για την λίστα με τα chunks που είχαν μόνο την ιδιότητα “Production”, ακολουθήσαμε τα παρακάτω βήματα, για να τρέξουμε τα παραδείγματα classification: Έχοντας επιλέξει την ιδιότητα, με βάση την οποία θέλουμε να κατηγοριοποιήσουμε τα αντικείμενα, η οποία ήταν είτε Product Cleaning Group είτε Product Group, αρχικά διαχωρίσαμε τα δεδομένα με βάση τρία διαφορετικά ποσοστά εκπαίδευσης-ελέγχου, όπως αναφέραμε στο κομμάτι της παραμετροποίησης προηγουμένως. Χρησιμοποιήσαμε τα δεδομένα εκπαίδευσης για να υπολογίσουμε τα κέντρα του αλγορίθμου. Από τον Πίνακα Μεταβάσεων του κάθε αντικειμένου του σετ, και μελετώντας την τιμή της μεταβλητής κατηγοριοποίησης του αντικειμένου, βρήκαμε τον μέσο πίνακα μεταβάσεων όλων των αντικειμένων που ανήκουν στην εκάστοτε κλάση αυτής της μεταβλητής. Αυτός ο πίνακας θεωρείται το κέντρο της κλάσης. Στη συνέχεια, για κάθε στοιχείο της κλάσης, από το αρχικό σετ ελέγχου, υπολογίσαμε την απόστασή του από τα διαφορετικά κέντρα. Τέλος αναθέσαμε αυτό το στοιχείο στην κλάση, της οποίας το κέντρο ήταν το κοντινότερο.

Αναλόγως με το πείραμα, χρησιμοποιήθηκαν διαφορετικές μέθοδοι υπολογισμού της απόστασης.

Λόγω του άνισου μεγέθους των κλάσεων, κάθε πείραμα εκτελέστηκε για 30 επαναλήψεις, δηλαδή, για 30 διαφορετικά σύνολα δεδομένων, του ίδιου ποσοστού διαχωρισμού. Αυτό μας επέτρεψε να επιβεβαιώσουμε τα αποτελέσματά μας και να είμαστε σίγουροι για την ακεραιότητα των συμπερασμάτων μας.

Την υλοποίηση του παραπάνω, σε ψευδοκώδικα, μπορείτε να την βρείτε εδώ [A.1.1](#).

10.1.2 K - Nearest Neighbors Classifier

Ο K - Nearest Neighbors είναι ένας γνωστός αλγόριθμος classification. Παρόλο που πολλές υλοποιήσεις υπάρχουν, αποφασίσαμε να υλοποιήσουμε τον δικό μας, προκειμένου να έχουμε τον απόλυτο έλεγχο σχετικά με τις μεταβλητές του και τον τρόπο εκτέλεσής του. Τα βήματα της εκτέλεσης είναι τα ακόλουθα και η μεθοδολογία είναι σχετικά κοντινή με αυτή του Nearest Centroid Classifier.

Το πρώτο βήμα ήταν να επιλέξουμε το χαρακτηριστικό για το οποίο θέλαμε να τρέξουμε τον classifier, αυτό μπορεί είτε να είναι το Product Cleaning Group είτε το Product Group. Ακολουθώντας, διαχωρίσαμε τα αρχικά μας δεδομένα σε δύο κατηγορίες, εκπαίδευσης και ελέγχου, για τρία διαφορετικά ποσοστά, όπως και στον Nearest Centroid Classifier. Για κάθε επανάληψη, διαχωρίσαμε τα δεδομένα σε τυχαία σετ, συγκεκριμένου ποσοστού διαχωρισμού, με σκοπό να επαληθεύσουμε την ακεραιότητα των αποτελεσμάτων μας. Έχοντας επιλέξει την μέθοδο υπολογισμού απόστασης, υπολογίσαμε την απόσταση κάθε αντικειμένου chunk του σετ ελέγχου με όλα τα αντικείμενα chunk του σετ εκπαίδευσης. Στη συνέχεια, ταξινομήσαμε τις τιμές για κάθε chunk σε αύξουσα σειρά. Η βασική αρχή του αλγορίθμου k-NN είναι ότι κάθε αντικείμενο του σετ ελέγχου κατηγοριοποιείται στην κλάση που ανήκει η πλειοψηφία των k κοντινότερων γειτόνων του. Μόλις είχε δημιουργηθεί η ταξινομημένη λίστα με όλους τους γείτονες, ήταν εύκολο να υπολογίσουμε την κλάση για κάθε τιμή του k, εντός των επιθυμητών ορίων.

Σε αυτό το σημείο είχαμε κατηγοριοποιήσει όλα τα αντικείμενα του σετ ελέγχου για κάθε k στην επιθυμητή εμβέλεια, για μια συγκεκριμένη μέθοδο υπολογισμού απόστασης. Στο Κεφάλαιο 11 μπορείτε να δείτε την επίδραση του k στην ακρίβεια του k-NN. Επαναλάβαμε το πείραμα για όλους τους αλγορίθμους υπολογισμού απόστασης που μας ενδιέφεραν.

Την υλοποίηση του παραπάνω, σε ψευδοκώδικα, μπορείτε να την βρείτε εδώ [A.1.2](#).

10.2 Clustering

Ο σκοπός της Clustering υλοποίησής μας ήταν να αναθέσουμε τα chunks προϊόντων σε clusters, με βάση κάποιες από τις ιδιότητές τους. Σε αυτή την κατεύθυνση, υλοποιήσαμε μια έκδοση του αλγορίθμου k-Means, ο οποίος αναλύεται στη συνέχεια:

10.2.1 Υλοποίηση του k-Means

- **Παράμετρος k του Clustering:** Ο αριθμός των clusters στον οποίο αναμένουμε τα δεδομένα μας να κατηγοριοποιηθούν, αρχικοποιείται με τον αριθμό των διαφορετικών τιμών που μπορεί να πάρει η ιδιότητα με βάση την οποία τα διακρίνουμε. Εκ τούτου, για το Product Cleaning Group, έχουμε $k=5$ και για το Product Group έχουμε $k=12$.
- **Αρχικά Κέντρα:** Το κέντρο αναπαρίσταται ως ένας Πίνακας Μεταβάσεων ενός αντικειμένου κλάσης chunk. Τα πειράματα διεξήχθησαν σε μια σειρά από αρχικά κέντρα, και τυχαία και συγκεκριμένα, με σκοπό να ελέγξουμε την απόδοση και την ακρίβεια του αλγορίθμου. Στο πρώτο σενάριο, κάθε κέντρο επιλέχθηκε τυχαία από τα δεδομένα εισαγωγής, ενώ στο δεύτερο, αρχικοποιήσαμε τα κέντρα επιλέγοντας σε κάθε σετ όλα να ακολουθούν έναν συγκεκριμένο κανόνα. Πιο συγκεκριμένα, επιλέχθηκαν είτε να ανήκουν όλα, του ίδιου σετ, στην ίδια κλάση, είτε το καθένα να ανήκει σε μια ξεχωριστή κλάση. Κάθε μία από τις 3 διαφορετικές περιπτώσεις εκτελέστηκε 100 φορές για να εξαλείψουμε οποιαδήποτε τυχαιότητα στις επιλογές μας.
- **Μέσος Όρος:** Σε κάθε μία από τις επαναλήψεις του αλγορίθμου, ένα νέο κέντρο του κάθε cluster υπολογιζόταν. Στην υλοποίησή μας, αυτός ο υπολογισμός γινόταν μέσω του Πίνακα Μεταβάσεων, δηλαδή, τον μέσο όρο όλων των Πινάκων Μεταβάσεων όλων των chunk παραγωγής του εκάστοτε cluster.

Την υλοποίηση του παραπάνω, σε ψευδοκώδικα, μπορείτε να την βρείτε εδώ [A.2.1](#).

10.3 Αποστάσεις

Οι προηγούμενες διεργασίες, Classification και Clustering, και οι αλγόριθμοι που υλοποιήθηκαν και στις δύο, απαιτούν τον υπολογισμό αποστάσεων μεταξύ των διαφόρων στοιχείων (chunks). Στο δικό μας πρόβλημα, επειδή οι διαφορετικές παρτίδες αναπαρίστανται μέσω των Πινάκων Μεταβάσεων, η απόσταση πρέπει να υπολογιστεί ως απόσταση μεταξύ αυτών των διδιάστατων πινάκων. Διαφορετικοί αλγόριθμοι παραμετροποιήθηκαν και εφαρμόστηκαν στα πειράματά μας, με σκοπό να μελετήσουμε την απόδοση και την ευστοχία της κάθε υλοποίησης. Χρησιμοποιήσαμε τις ακόλουθες μεθόδους:

- Euclidean distance
- Cosine Distance
- Kullback - Leibler Divergence
- Kolmogorov - Smirnov Test
- Infinity Norm

Οι περισσότερες από αυτές, εξ ορισμού, εφαρμόζονται μόνο σε μονοδιάστατους πίνακες, δηλαδή σε διανύσματα, γι' αυτό και ένα σημαντικό βήμα ήταν να μετασχηματίσουμε σωστά τους διδιάστατους πίνακες σε διανύσματα, με σκοπό να υπολογίσουμε τις αποστάσεις τους. Οι προσεγγίσεις μας χωρίζονται σε τρεις κατηγορίες:

- **Average rows:** Σε αυτή την περίπτωση, ο αλγόριθμος εφαρμόζεται μεταξύ των αντίστοιχων σειρών και στο τέλος επιστρέφουμε τον μέσο όρο των μετρήσεων.

- **Vector:** κάθε γραμμή προσαρτάται στην πρώτη και έτσι δημιουργείται το επιθυμητό διάνυσμα $1 \times N$.
- **Diagonal:** για κάθε $j > i$, ο μέσος όρος των στοιχείων $[i, j]$ και $[j, i]$ υπολογίζεται, έτσι δημιουργείται ένας άνω τριγωνικός πίνακας. Στην συνέχεια ακολουθούμε την προηγούμενη μέθοδο και προσαρτούμε την ημιγραμμή στην προηγούμενη.

Την υλοποίηση του παραπάνω, σε ψευδοκώδικα, μπορείτε να την βρείτε εδώ [10.3.1](#) .

Η συνολική λίστα των υλοποιήσεων των μεθόδων υπολογισμού απόστασης φαίνεται στη συνέχεια:

- **Euclidean distance**
 - Euclidean distance between each corresponding cell
 - Average of the Euclidean distances between corresponding rows
 - Average of the Euclidean distances between corresponding columns
- **Cosine distance**
 - Average of the cosine distance between corresponding rows
 - Cosine distance of the created vectors using the Vector transformation
 - Cosine distance of the created vectors using the Diagonal transformation
- **Kullback–Leibler Divergence**
 - Average of the KL Divergence between corresponding rows
 - KL Divergence of the created vectors using the Vector transformation
 - KL Divergence of the created vectors using the Diagonal transformation
- **Kolmogorov–Smirnov test**
 - Average of the KS test score between corresponding rows
 - KS test score of the created vectors using the Vector transformation
 - KS test score of the created vectors using the Diagonal transformation
- **Infinity Norm**
 - Infinity Norm of the difference of the two matrices

Για την εφαρμογή των παραπάνω αλγορίθμων απόστασης, η βιβλιοθήκη Python SciKit-Learn χρησιμοποιήθηκε. Αυτή η βιβλιοθήκη παρέχει πολλές χρήσιμες συναρτήσεις υπολογισμού αποστάσεων, αλλά οι περισσότερες αφορούν υπολογισμό μεταξύ διανυσμάτων.

10.3.1 Υλοποίηση Μεθόδων Αποστάσεων

Δύο από τις συναρτήσεις που μετασχηματίζουν έναν $2D$ πίνακα σε ένα , μπορούν να βρεθούν σε υλοποίηση ψευδοκώδικα στο [A.3.1](#).

Οι διαφορετικές υλοποιήσεις των αλγορίθμων απόστασης, όπως περιγράφηκαν στην αρχή της Ενότητας, μπορούν να βρεθούν στο [A.3.2](#).

10.3.2 Αξιολόγηση των Μεθόδων Αποστάσεων

Μας ενδιέφερε να αξιολογήσουμε την απόδοση όλων των μεθόδων συναρτήσει των διαφορετικών παραμέτρων υπό αξιολόγηση, τόσο για το classification όσο και για το clustering. Ο λόγος που θέλαμε να το τρέξουμε για όλα αυτά τα σενάρια ήταν για να καταλάβουμε το αντίκτυπο που έχουν σε κάθε περίπτωση.

- Αξιολόγηση των Μεθόδων Αποστάσεων στο Classification

Αποφασίσαμε να τρέξουμε το ακόλουθο πείραμα με τον Nearest Centroid Classifier για τον διαχωρισμό 80%. Ο αλγόριθμος έτρεξε τόσο για το Product Cleaning Group όσο και για το Product Group για 30 επαναλήψεις, όπως περιγράφηκε προηγουμένως. Για κάθε επανάληψη, επιλέξαμε τυχαία έναν διαχωρισμό, πάντα με σταθερά τα ποσοστά 80%-20%, εκπαίδευσης και ελέγχου αντίστοιχα, και υπολογίσαμε τα κέντρα του κάθε διαφορετικού γκρουπ. Αφού επιλέξαμε μια μέθοδο, θέσαμε το κάθε αντικείμενο στο κοντινότερο κέντρο, σύμφωνα με τον επιλεγμένο αλγόριθμο.

Την υλοποίηση του παραπάνω, σε ψευδοκώδικα, μπορείτε να την βρείτε εδώ [A.3.3](#). Τα αποτελέσματα μπορούν να βρεθούν στην Ενότητα [11.1.2](#).

- Αξιολόγηση των Μεθόδων Αποστάσεων στο Clustering

Στην περίπτωση του clustering, ο μόνος αλγόριθμος που δοκιμάσαμε ήταν ο K - Means. Την υλοποίηση του παραπάνω, σε ψευδοκώδικα, μπορείτε να την βρείτε εδώ [A.3.4](#).

10.4 Αξιολόγηση των Αποτελεσμάτων

Όλα τα προηγούμενα μπορούν να φανούν χρήσιμα μόνο αν το αποτέλεσμά τους μπορεί να αξιολογηθεί. Στα αρχικά πειράματα χρησιμοποιήσαμε τον Confusion Matrix για να μπορούμε να έχουμε μια πιο άμεση εποπτεία του αποτελέσματος. Καθώς τα πειράματα εξελισσότουσαν, δημιουργήθηκε η ανάγκη να μπορούμε να αξιολογούμε τα αποτελέσματα με ένα μοναδικό σκορ, και έπρεπε να επιλέξουμε τον κατάλληλο τρόπο να γίνει αυτή η αξιολόγηση. Η βασική διαφορά στον τρόπο αξιολόγησης μεταξύ Classification και Clustering είναι ότι στην πρώτη περίπτωση γνωρίζουμε από πριν την σωστής ετικέτες για τα δεδομένα μας, ενώ στην δεύτερη όχι.

10.4.1 Υλοποίηση της Αξιολόγησης του Classification

Για την αξιολόγηση των αποτελεσμάτων του Classification αξιοποιήσαμε τις παρακάτω μεθόδους:

- **Ακρίβεια (Accuracy)**

Η Ακρίβεια είναι ορισμένη ως ο λόγος των σωστά κατηγοριοποιημένων αντικειμένων προς όλα τα αντικείμενα που δόθηκαν στην είσοδο του αλγορίθμου.

$$\text{Accuracy} = \frac{\text{Number of Correctly Classified Elements}}{\text{Total Number of Classified Elements}}$$

- **Kappa Coefficient**

Ο αλγόριθμος Kappa είναι μέρος του SciKit-Learn Laboratory Library και συγκεκριμένα, του Metrics Module. Το SciKit είναι μια γνωστή βιβλιοθήκη για Python και μας φάνηκε πολύ χρήσιμη κατά τη διάρκεια αυτής της εργασίας.

Πηγή: [SciKit-Learn Laboratory -> Metrics -> Kappa](#)

10.4.2 Υλοποίηση της Αξιολόγησης του Clustering

Για την αξιολόγηση των αποτελεσμάτων του Clustering χρησιμοποιήθηκαν οι παρακάτω μέθοδοι. Παρόλο που σε άλλες περιπτώσεις μπορεί να μην είχαμε τις σωστές ετικέτες, σε αυτή τις είχαμε και τις αξιοποιήσαμε ως είσοδο για την καλύτερη αξιολόγηση των αποτελεσμάτων. Και οι δύο τρόποι αξιολόγησης χρειάζονται κάποιες ετικέτες ως ground truth για να μπορέσουν να μας δώσουν αποτελέσματα.

- **V-Measure**

Source: [SciKit-Learn -> Metrics -> V Measure Score](#)

- **Rand Index**

Source: [SciKit-Learn -> Metrics -> Adjusted Rand Score](#)

Κεφάλαιο 11

Αποτελέσματα

Στο Κεφάλαιο αυτό παρουσιάζουμε τα αποτελέσματα από την εφαρμογή των αλγορίθμων Classification και Clustering στα δεδομένα μας. Όλοι οι αλγόριθμοι, όπως περιγράφηκαν στα κεφάλαια 8 και 10, εκτελέστηκαν για πολλαπλές παραμέτρους, δίνοντας μας διαφορετικά αποτελέσματα και συμπεράσματα για κάθε συνδυασμό.

Στο 11.1 παρουσιάζουμε τα αποτελέσματα του Classification, σε τρεις διαφορετικές κατηγορίες. Πρώτα, συγκρίνουμε τους αλγορίθμους Nearest Centroid και k-Nearest Neighbours, για διαφορετικά ποσοστά training, testing, κρατώντας σταθερή τη μέθοδο υπολογισμού της απόστασης. Έπειτα, επικεντρώναστε αποκλειστικά στον πρώτο, για συγκεκριμένο ποσοστό train-test, ίσο με 80% - 20%, και αξιολογούμε την απόδοση του αλγορίθμου για όλες τις διαφορετικές μεθόδους εύρεσης απόστασης μεταξύ των αντικειμένων. Τέλος, εξετάζουμε την εφαρμογή του k-NN για διαφορετικές τιμές του k και το αποτέλεσμα αυτής.

Στο 11.2 παρουσιάζονται και αξιολογούνται τα αποτελέσματα του Clustering. Χρησιμοποιήσαμε τον αλγόριθμο k-Means για το διαχωρισμό των αντικειμένων σε clusters, βασιζόμενοι ξανά σε δύο γνωρίσματά τους, το Product Cleaning Group και το Product Group. Τα δεδομένα εισόδου του αλγορίθμου είναι τα αντικείμενα chunk. Τα κέντρα και οι αποστάσεις αλλάζουν σε κάθε εκτέλεση, έτσι ώστε να μπορούμε να αξιολογήσουμε το αποτέλεσμα για διαφορετικές παραμετροποιήσεις του προβλήματος, με τη χρήση δύο μεθόδων, των V-Measure και Rand Index.

11.1 Classification

Το Classification πραγματοποιήθηκε για δύο ξεχωριστά γνωρίσματα, τα Product Cleaning Group και Product Group.

Οι αλγόριθμοι που εξετάστηκαν ήταν οι:

- Nearest Centroid Classifier
- k-Nearest Neighbors Classifier

Ο υπολογισμός της απόστασης ήταν διαφορετικός σε κάθε περίπτωση.

Όλα τα δεδομένα χωρίστηκαν σε training και test, το ποσοστό των οποίων αλλάζει, ώστε να εξετάσουμε καλύτερα την ακρίβεια των αλγορίθμων και την πρόβλεψη της σωστής κλάσης.

11.1.1 Baseline

Χρησιμοποιήσαμε ένα baseline αποτέλεσμα, με σκοπό να συγκρίνουμε τα αποτελέσματα του classification με αυτό, και να αξιολογήσουμε το αποτέλεσμα της κάθε μεθόδου.

Ο αλγόριθμος που υλοποιήσαμε ήταν ο zeroR, ο οποίος επιλέγει την κλάση που περιλαμβάνει τις περισσότερες παρατηρήσεις, και τη χρησιμοποιεί σαν το αποτέλεσμα για όλες τις προβλέψεις.

Το σκορ για τα δεδομένα εισόδου υπολογίστηκε με την εύρεση του ποσοστού των αντικειμένων κάθε κλάσης στο συνολικό αριθμό προϊόντων, και με επιλογή της κλάσης που μας δίνει το μέγιστο σκορ.

Η ακρίβεια zeroR για τα δύο γνωρίσματα φαίνεται ακολούθως:

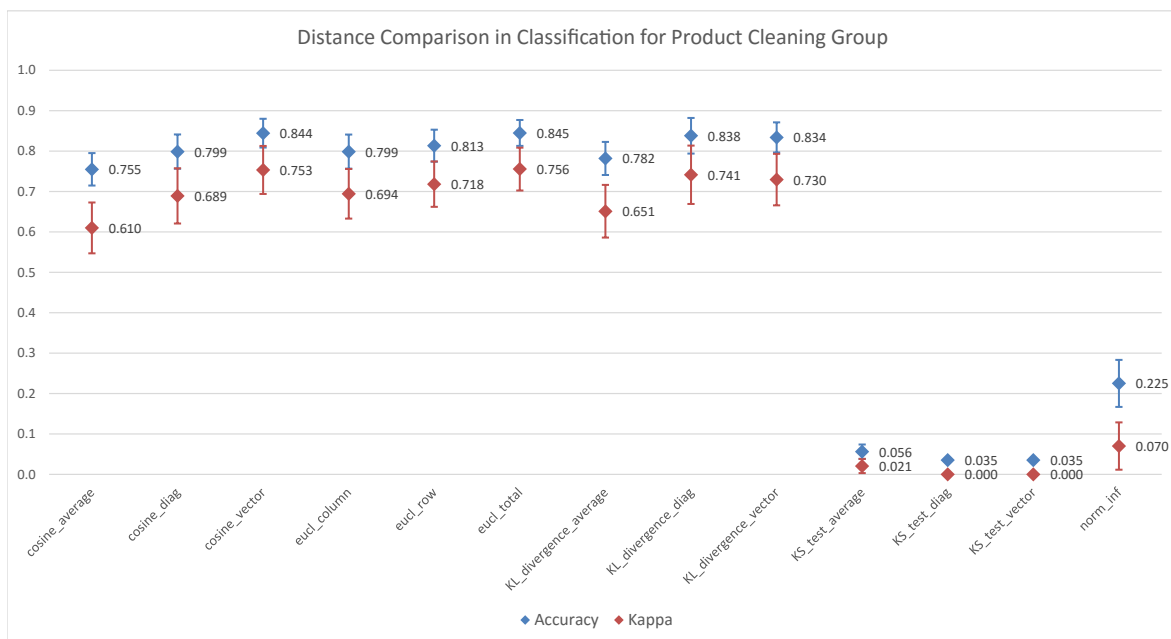
- ZeroR Product Cleaning Group Accuracy: 0.520
- ZeroR Product Group Accuracy: 0.377

11.1.2 Αξιολόγηση Μεθόδου Απόστασης

Setup:

- *Algorithm*: Nearest Centroid Classifier
- *Attributes*:
 - Product Cleaning Group
 - Product Group
- *Split*: 80% training set, 20% test set
- *Distances*:
 - Euclidean Total
 - Euclidean Row
 - Euclidean Column
 - Cosine Average
 - Cosine Vector
 - Cosine Diagonal
 - KL - Divergence Average
 - KL - Divergence Vector
 - KL - Divergence Diagonal
 - KS - Test Average
 - KS - Test Vector
 - KS - Test Diagonal
 - Infinity Norm

Στο πείραμα αυτό εξετάζουμε τις διαφορετικές μεθόδους υπολογισμού απόστασης, θεωρώντας σταθερό το ποσοστό 80% train και 20% test, και χρησιμοποιώντας τον αλγόριθμο Nearest Centroid. Τα αποτελέσματα φαίνονται στο 11.1 για το Product Cleaning Group.



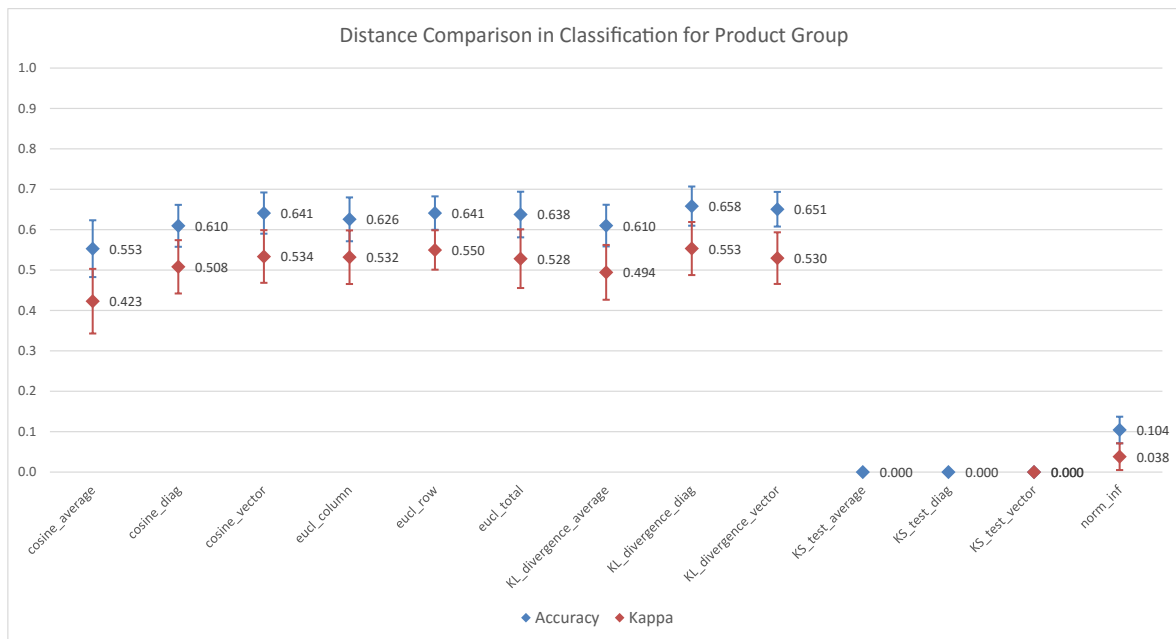
Σχήμα 11.1: Σύγκριση Μεθόδων Απόστασης για Classification στο Product Cleaning Group

Τα μέγιστα σκορ Accuracy και Kappa Coefficient εμφανίζονται στη Euclidean Total απόσταση, και ακολουθούνται από αυτά των Cosine Vector Distance και KL - Divergence- Diagonal. Όλα τα άλλα Euclidean, Cosine και KL σκορ είναι σε αρκετά κοντινή απόσταση με τα βέλτιστα.

Τα σκορ των KS και Infinity Norm εμφανίζονται σημαντικά χαμηλά, υποδεικνύοντας ότι είναι ακατάλληλα για τον υπολογισμό της απόστασης στο παρόν πρόβλημα.

Η τυπική απόκλιση του Accuracy έχει μέγιστη τιμή 0,058, ενώ στο Kappa Coefficient μεγιστοποιείται στο 0,068, επομένως θεωρούμε ότι μπορούμε να την αγνοήσουμε στα συμπεράσματά μας και για τις δύο μεθόδους.

Για το Product Group, τα αποτελέσματα φαίνονται στο 11.2, με τα KL - Divergence Diagonal και KL - Divergence Vector να εμφανίζουν τα βέλτιστα σκορ. Αμέσως μετά ακολουθούν τα KL - Divergence Average και Cosine Average.



Σχήμα 11.2: Σύγκριση Μεθόδων Απόστασης για Classification στο Product Group

Όπως και στο προηγούμενο πείραμα, τα χαμηλότερα σκορ εμφανίζονται στα KS και Infinity Norm, τα οποία ξανά θεωρούμε ότι δε ταιριάζουν στη συγκεκριμένη διαδικασία κατηγοριοποίησης.

Η τυπική απόκλιση των δυο μεθόδων αξιολόγησης εμφανίζει μέγιστο σκορ 0,056 και 0,799 αντιστοίχως, επομένως θεωρείται ασήμαντη για την ακρίβεια των αποτελεσμάτων μας.

11.1.3 Σύγκριση Nearest Centroid και K-Nearest Neighbors

Setup:

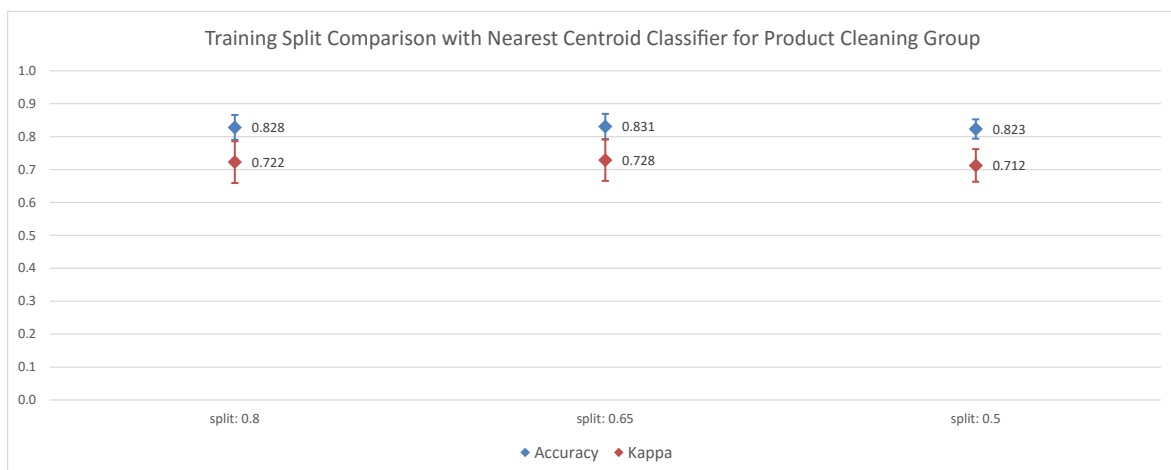
- *Algorithms:*
 - Nearest Centroid Classifier
 - k-Nearest Neighbors Classifier
- *Attributes:*
 - Product Cleaning Group
 - Product Group
- *Splits:*
 - 80% - 20%
 - 65% - 35%
 - 50% - 50%
- *Distance: Average of*
 - Euclidean Total
 - Cosine Vector
 - KL - Divergence Diagonal

Στο πείραμα αυτό συγκρίνουμε το αποτέλεσμα των δύο αλγορίθμων, Nearest Centroid και k-Nearest Neighbors, στα δεδομένα μας, εξετάζοντας τις ακόλουθες παραμέτρους:

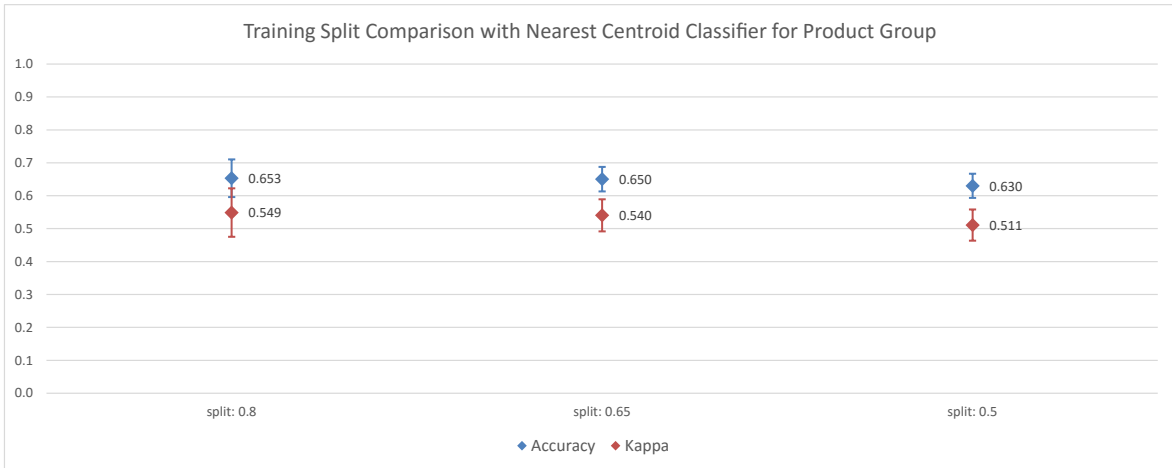
- Χρόνος Εκτέλεσης
- Ποσοστό διαχωρισμού train και test δεδομένων
- Accuracy και Kappa Coefficient του αποτελέσματος

Η απόσταση μεταξύ των πινάκων μεταβάσεων κάθε αντικειμένου chunk διατηρήθηκε σταθερή σε αυτό το κομμάτι και υπολογίστηκε σαν ο μέσος όρος των τριων καλύτερων σκορ από τις αποστάσεις που χρησιμοποιήθηκαν στο προηγούμενο πείραμα: Euclidean Total, Cosine Distance- Vector και KL - Divergence Diagonal.

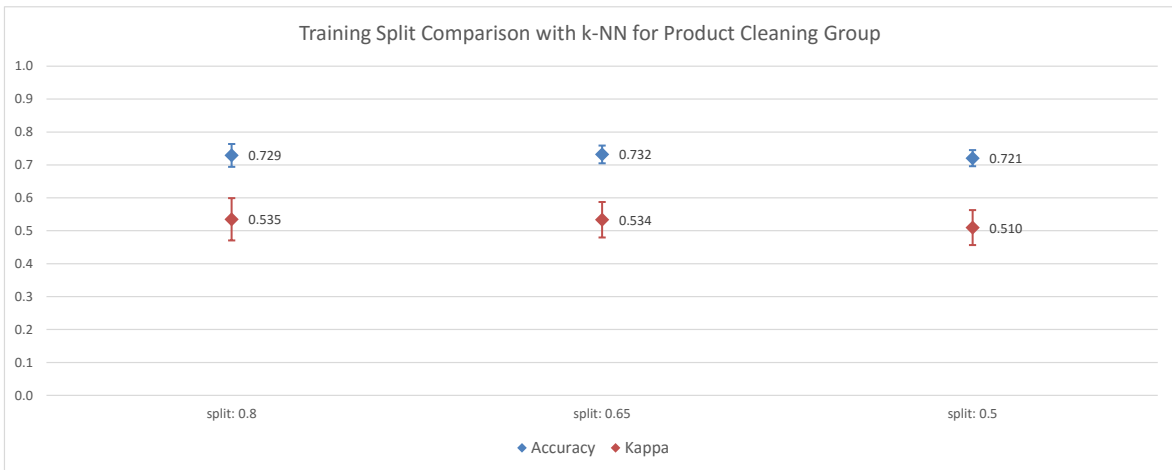
Τα αποτελέσματα της εφαρμογής του αλγορίθμου παρουσιάζονται στα 11.3 and 11.4, ενώ τα αποτελέσματα του k-NN στα 11.5 and 11.6.



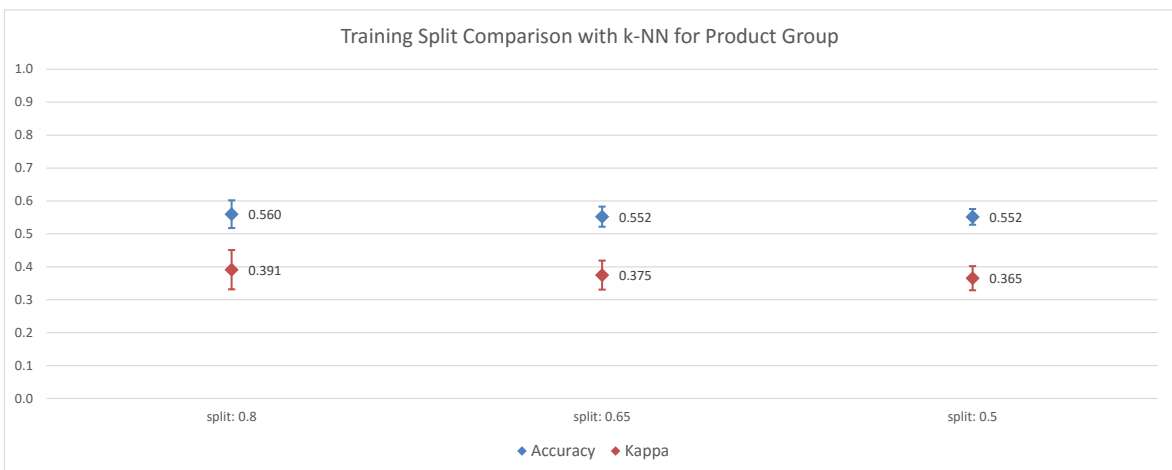
Σχήμα 11.3: Nearest Centroid Classifier για Product Cleaning Group



Σχήμα 11.4: Nearest Centroid Classifier για Product Group



Σχήμα 11.5: k-Nearest Neighbors Classifier για Product Cleaning Group



Σχήμα 11.6: k-Nearest Neighbors Classifier για Product Cleaning Group

Είναι προφανές ότι, παρά τα τρία διαφορετικά ποσοστά διαχωρισμού, τα αποτελέσματα του classification είναι αρκετά κοντινά, από τα διαγράμματα των δύο αλγορίθμων. Όσον αφορά στο Product Cleaning Group, ο πρώτος αλγόριθμος δίνει μέγιστα σκορ στο διαχωρισμό 65%, με Accuracy 83,1% και Kappa

Coefficient 72,8%, ενώ τα αντίστοιχα σκορ για τον k-NN είναι 73,2% και 53,4%.

Τα συμπεράσματα είναι παρόμοια όταν εξετάζουμε το Product Group. Και εδώ οι διαφορές είναι πολύ μικρές, ανάμεσα στα διαφορετικά ποσοστά διαχωρισμού, με το πρώτο να μας δίνει ένα ελαφρώς υψηλότερο ποσοστό από τα υπόλοιπα. Για τον Nearest Centroid, αυτό το μέγιστο είναι 65,3% Accuracy και 54,9% Kappa Coefficient, ενώ για τον k -NN the αντίστοιχα σκορ είναι 56% και 39,1%.

Από τα αποτελέσματα αυτά, είναι εμφανές ότι ο αλγόριθμος Nearest Centroid είναι ανώτερος του k-NN, για το συγκεκριμένο πρόβλημα και τις παραμέτρους.

Ένα ισχυρό επιχείρημα που υποστηρίζει το παραπάνω είναι επίσης και το γεγονός ότι η εκτέλεση του πρώτου χρειάστηκε 20,15 δευτερόλεπτα, ενώ του δεύτερου 161,46, δηλαδή τουλάχιστον 8 φορές περισσότερο χρόνο.

11.1.4 Αλγόριθμος K-Nearest Neighbors για διαφορετικά k

Setup:

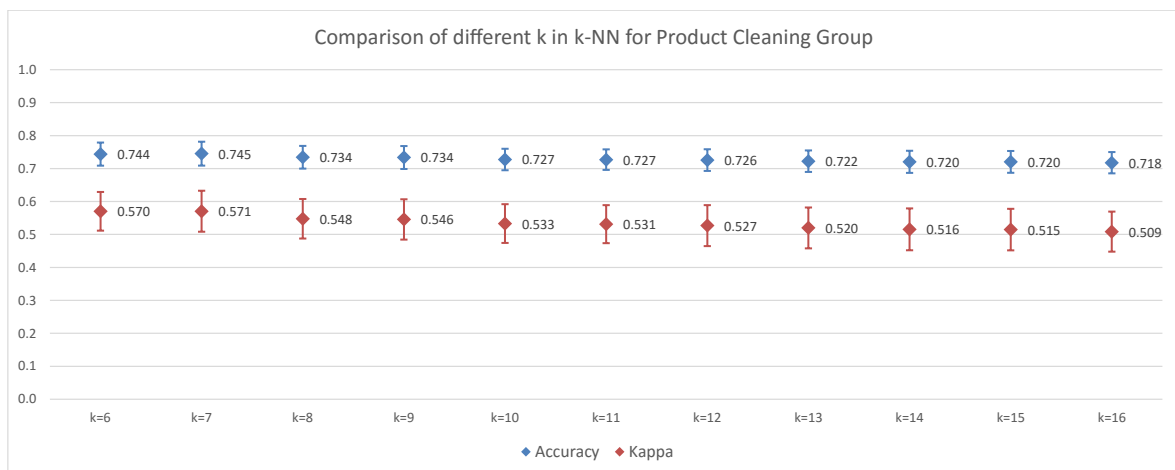
- *Algorithms:* k-Nearest Neighbors Classifier
- *Attributes:*
 - Product Cleaning Group
 - Product Group
- *Split:* 80% - 20%
- *Distance:* Average of
 - Euclidean Total
 - Cosine vector
 - KL-Divergence diagonal

Εφαρμόστηκε ο αλγόριθμος k-NN για τις παραπάνω παραμέτρους, ώστε να συγκρίνουμε την επίδραση της επιλογής του k στη γενικότερη απόδοση του αλγορίθμου.

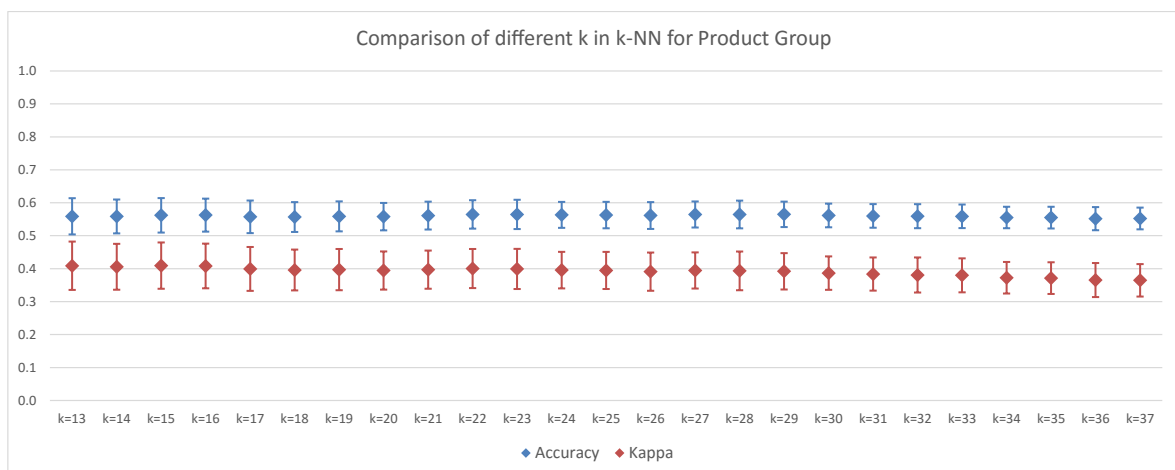
Ο k-NN υλοποιήθηκε για τα δύο γνωρίσματα που εξετάστηκαν και προηγουμένως, τα Product Cleaning Group και Product Group, με ια ποσοστα train και test στα 80% και 20%.

Η απόσταση διατηρήθηκε σταθερή επίσης, και, όπως και πριν, ίση με το μέσο όρο των τριων καλύτερων σκορ από τις Euclidean Total, Cosine Distance- Vector και KL - Divergence Diagonal.

Τα αποτελέσματα παρουσιάζονται στα 11.7 και 11.8



Σχήμα 11.7: k-NN για διαφορετικά k στο Product Cleaning Group



Σχήμα 11.8: k-NN για διαφορετικά k στο Product Group

Παρατηρώντας τα διαγράμματα, είναι εμφανές ότι η επιλογή του k έχει πολύ μικρή επίδραση στην ακρίβεια του αλγορίθμου και τα σκορ, ακόμα και για διαφορετικές ή μακρινές τιμές του k .

11.2 Clustering

11.2.1 Baseline

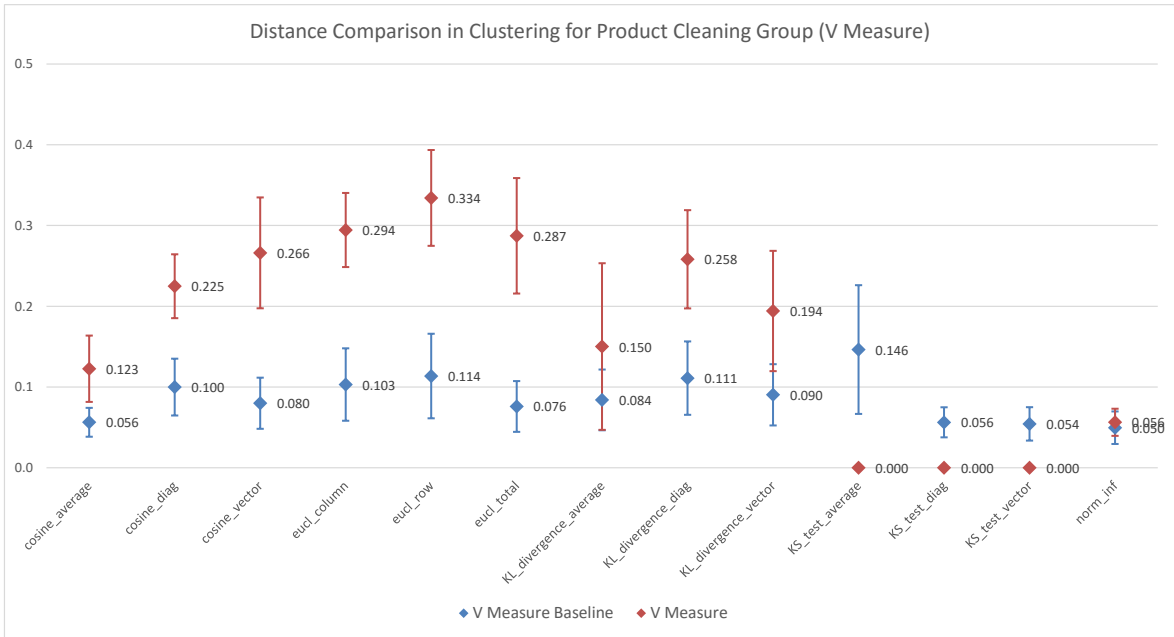
Η τιμή baseline για το Clustering υπολογίζεται από την αρχικοποίηση του αλγορίθμου για κάθε πείραμα, το οποίο μας δίνει δύο ξεχωριστά σκορ. Το αποτέλεσμα κάθε εκτέλεσης συγκρίνεται με αυτά, όπως φαίνεται και σε όλα τα διαγράμματα παρακάτω, έτσι ώστε να εξετάσουμε την ακρίβεια των αποτελεσμάτων.

11.2.2 Μέθοδοι Υπολογισμού Αποστάσεων

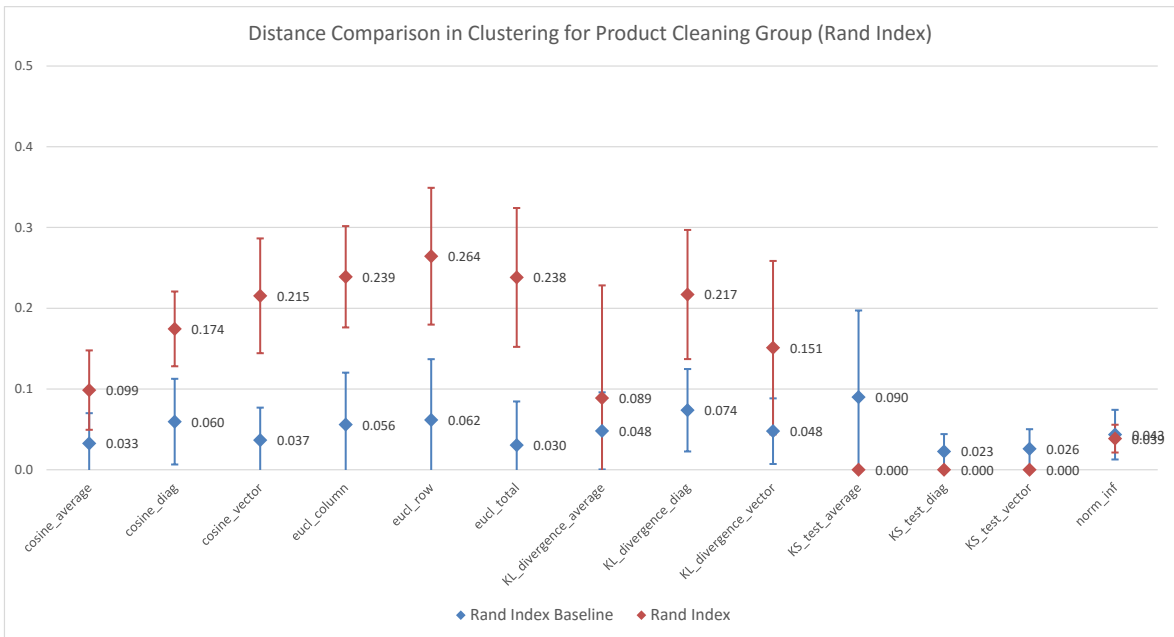
Setup:

- *Algorithms:*
 - k-Means
 - Baseline
- *Attributes:*
 - Product Cleaning Group
 - Product Group
- *Initial Centroid Sets Type:*
 - All centroids of each set belonged to different clusters (Alldiff)
 - * Average of 20 sets
 - All centroids of each set belonged to the same cluster (Allsame)
 - * Average of 20 sets
- *Distances:*
 - Euclidean Total
 - Euclidean Row
 - Euclidean Column
 - Cosine Average
 - Cosine Vector
 - Cosine Diagonal
 - KL - Divergence Average
 - KL - Divergence Vector
 - KL - Divergence Diagonal
 - KS - Test Average
 - KS- Test Vector
 - KS- Test Diagonal
 - Infinity Norm
- *Evaluation Methods:*
 - V-Measure
 - Rand Index

Στα 11.9 και 11.10, ο αλγόριθμος k-Means εκτελέστηκε για 20 διαφορετικές ομάδες αρχικών κέντρων, τα οποία επιλέχθηκαν έτσι ώστε όλα της ίδιας ομάδας να ανήκουν είτε στο ίδιο είτε σε διαφορετικά αρχικά clusters, τα οποία είναι τα 5 Product Cleaning Groups. Παρατηρούμε ότι οι Euclidean μέθοδοι αποδίδουν το μέγιστο, με σκορ 33,4% V-Measure και 26,4% Rand Index.

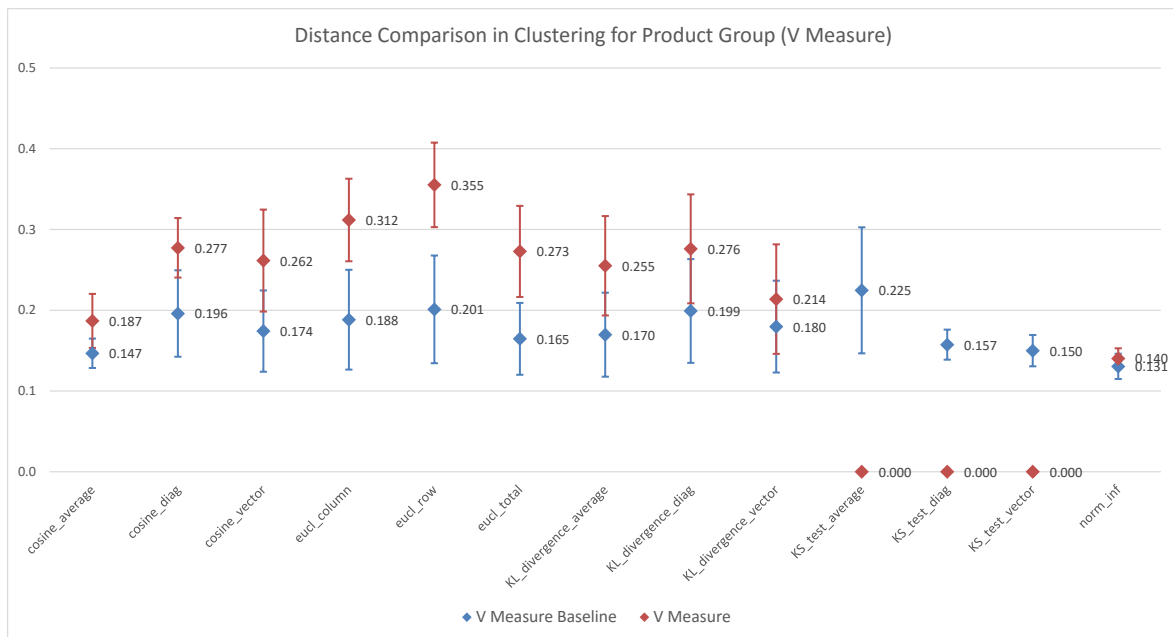


Σχήμα 11.9: k-Means Clustering στο Product Cleaning Group (V-Measure)

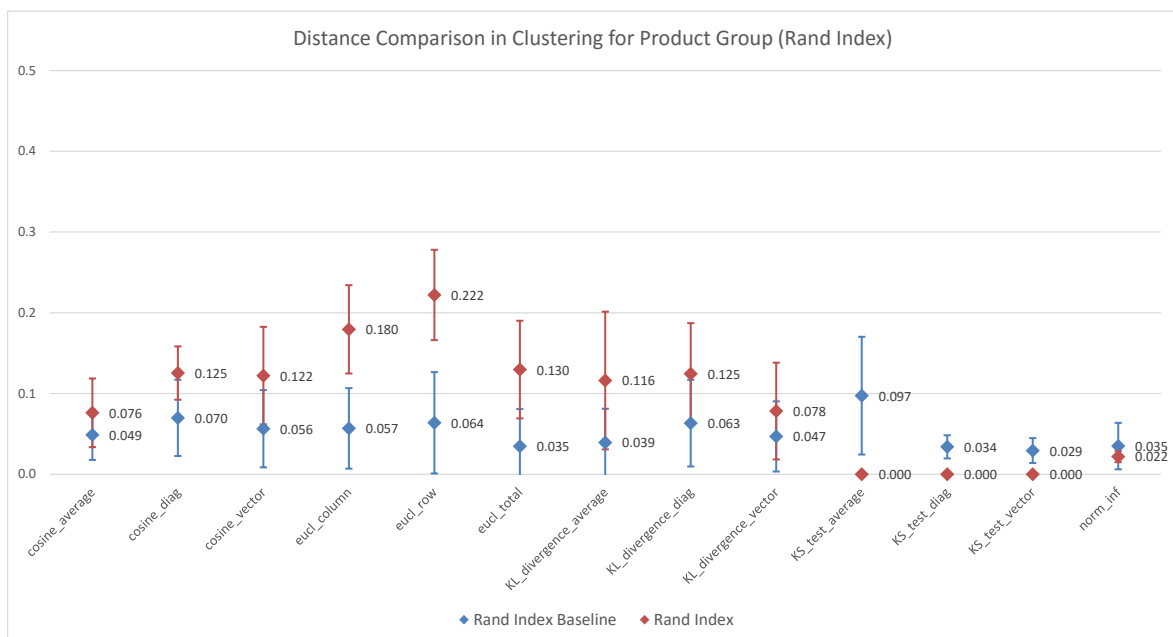


Σχήμα 11.10: k-Means Clustering στο Product Cleaning Group (Rand Index)

Τα 11.11 και 11.12 ακολουθούν την ίδια λογική με τα προηγούμενα, αλλά εδώ το clustering πραγματοποιείται για το Product Group. Έχουμε δηλαδή 12 πιθανά clusters αντί για τα 5 που είχαμε στα προηγούμενα δύο πειράματα.



Σχήμα 11.11: k-Means Clustering στο Product Group (V Measure)



Σχήμα 11.12: k-Means Clustering στο Product Group (Rand Index)

Οι Euclidean αποστάσεις έχουν και εδώ τα κυρίαρχα σκορ, με τη μέγιστη ακρίβεια να επιτυγχάνεται για τη Euclidean Row, που μας δίνει V-Measure 35,5% και Rand Index 22,2%. Τα σκορ KL-Divergence είναι επίσης αρκετά υψηλά.

Σε όλες τις περιπτώσεις, τα χαμηλότερα σκορ δίνονται από τα KS-Test και Infinity Norm, και πολλές φορές δεν ξεπερνούν το baseline.

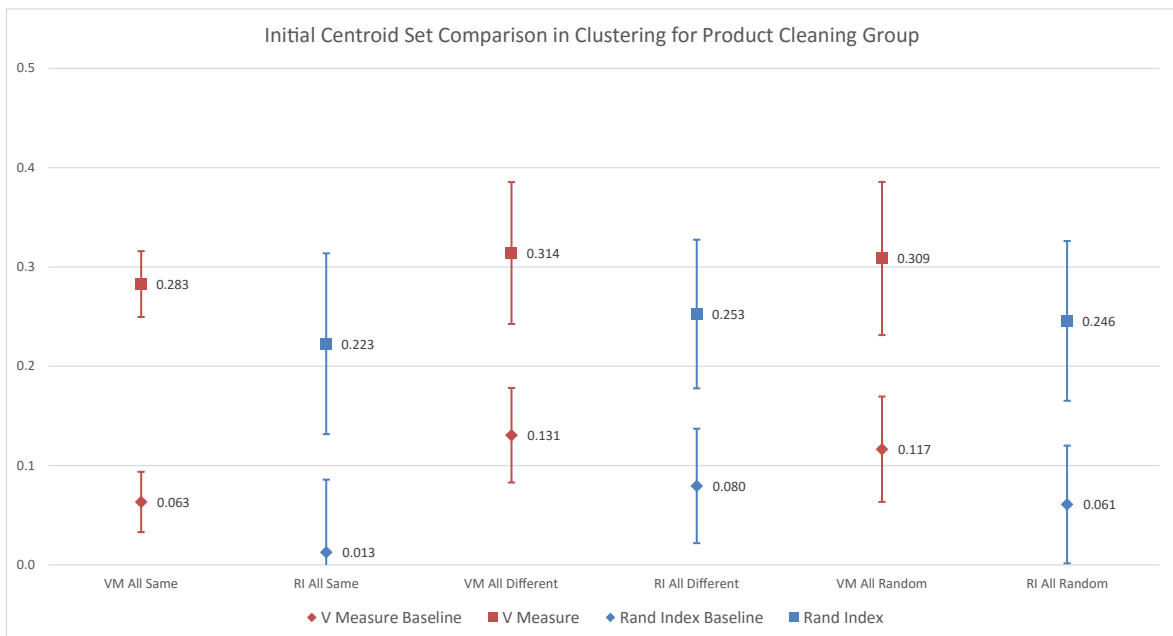
11.2.3 Αρχικά κέντρα

Setup:

- Algorithms:

- k-Means
- Baseline
- *Attributes:*
 - Product Cleaning Group
 - Product Group
- *Initial Centroid Sets Type:*
 - All centroids of each set belonged to different clusters (Alldiff)
 - * Average of 100 sets
 - All centroids of each set belonged to the same cluster (Allsame)
 - * Average of 100 sets
 - All centroids of each set belonged to a random cluster (Allrand)
 - * Average of 100 sets
- *Distance: Average Of*
 - Euclidean Total
 - Euclidean Row
 - Euclidean Column
- *Evaluation Methods:*
 - V-Measure
 - Rand Index

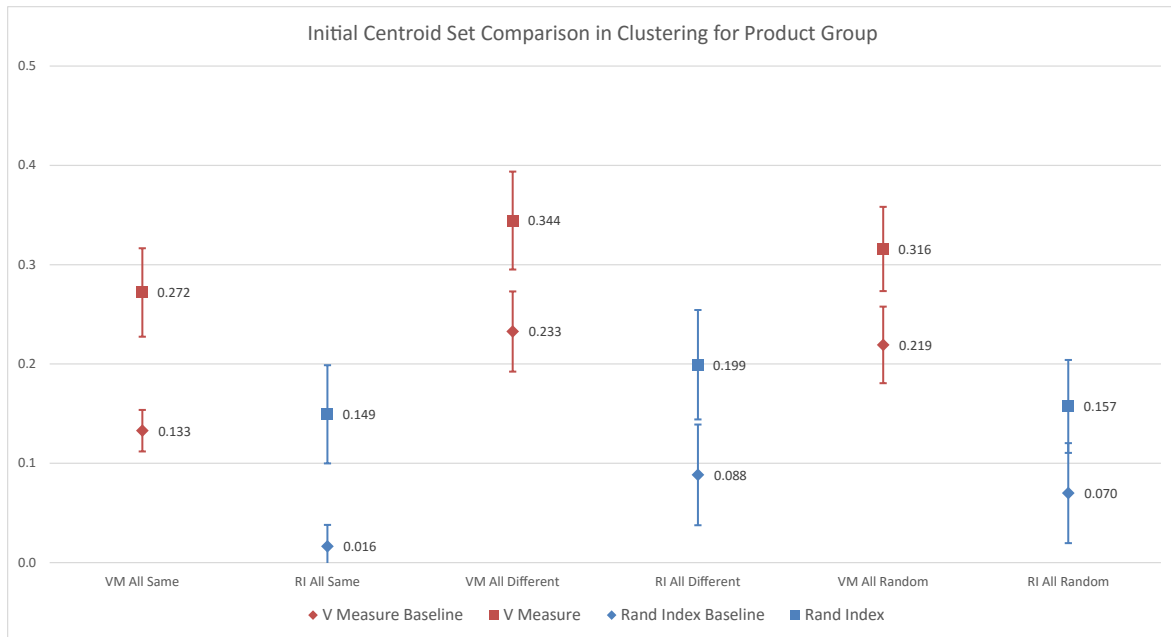
Το 11.13 παρουσιάζει τα αποτελέσματα για το παραπάνω σενάριο ένταξης σε 5 clusters, τα Product Cleaning Groups.



Σχήμα 11.13: k-Means Clustering(Product Cleaning Group)

Τα μέγιστα σκορ παρατηρούνται όταν τα αρχικά κέντρα ανήκουν σε διαφορετικά clusters, οπότε και έχουμε V-Measure 31,4% και Rand Index 25,3%. Αντίθετα, το χαμηλότερο σκορ εμφανίζεται όταν τα αρχικά κέντρα ανήκουν όλα στο ίδιο cluster.

Το 11.14 ακολουθεί την παραπάνω μεθοδολογία, με τη διαφορά ότι η διαδικασία εστιάζει σε 12 clusters, μιας και εξετάζουμε τα Product Groups.



Σχήμα 11.14: k-Means Clustering(Product Group)

Διαπιστώνουμε και εδώ ότι το υψηλότερο σκορ παράγεται στην περίπτωση των αρχικών κέντρων σε διαφορετικά clusters, και το χειρότερο όταν ανήκουν στο ίδιο.

Κεφάλαιο 12

Επίλογος

Στο κεφάλαιο αυτό παρουσιάζονται τα συμπεράσματα από τα πειράματα που εκτελέσαμε, και οι προτάσεις μας για μελλοντική εργασία στο παρόν πρόβλημα.

Στο 12.1 παρατίθενται τα συμπεράσματά μας από τα πειράματα του Κεφαλαίου 11, και συνοψίζεται η απόδοση των Classification και Clustering, αξιολογώντας την επίδραση κάθε παραμέτρου.

Στο 12.2 περιγράφονται πιθανά θέματα που δεν είχαμε την ευκαιρία να εξετάσουμε κατά την παρούσα εργασία, και θα μπορούσαν να ερευνηθούν περισσότερο στο μέλλον.

12.1 Συμπεράσματα

Εδώ συνοψίζουμε τις παρατηρήσεις από τα πειράματά μας και συζητάμε τα τελικά συμπεράσματα.

12.1.1 Classification

Αναφορικά με την αποδοτικότητα των μεθόδων υπολογισμού απόστασης, οι Euclidean μας δίνουν πάντα τα υψηλότερα σκορ, μαζί με τις KL-Divergence και τις Cosine. Ήταν αρκετά δύσκολο να αναγνωρίσουμε κατα πόσο οι τεχνικές average, diagonal ή vector απέδωσαν καλύτερα, καθώς, για κάθε μέθοδο, το βέλτιστο σκορ εμφανίζεται για κάποια διαφορετική τεχνική, χωρίς να υπάρχει κάποιο γενικό μοτίβο. Τα σκορ των KS-Test και Infinity Norm αποδείχθηκαν ακατάλληλα για τα δεδομένα μας, και τα αποτελέσματα σπάνια ξεπερνούσαν τις τιμές baseline. Υποθέτουμε ότι αυτό συνέβαινε λόγω της φύσης των κατανομών μας, και των δεδομένων μας στη μορφή των πινάκων μετάβασης.

Κατά τη σύγκριση των δύο αλγορίθμων, των Nearest Centroid και K-Nearest Neighbours, ο πρώτος εμφανίζει σημαντικά καλύτερη απόδοση, λαμβάνοντας υπόψη τόσο το Accuracy όσο και το Kappa. Όπως αναφέρθηκε μάλιστα και στο Κεφάλαιο 11, ο χρόνος εκτέλεσης του k-NN ήταν 8 φορές υψηλότερος από του άλλου αλγορίθμου, γεγονός που ενισχύει την ανωτερότητα του δεύτερου στο πρόβλημά μας. Αυτό αποδίδεται και εν μέρει στο γεγονός ότι οι υπολογισμοί που απαιτείται να πραγματοποιήσει ο k-NN είναι αρκετά πιο πολύπλοκοι, καθώς περιλαμβάνουν όλους τους υπολογισμούς αποστάσεων από τους k-κοντινότερους γείτονους. Ο διαχωρισμός των δεδομένων σε train και test είχε πολύ μικρή επίδραση στο αποτέλεσμα, κάτι ιδιαίτερα αξιοσημείωτο, αν σκεφθούμε ότι ο όγκος των δεδομένων που χρησιμοποιήθηκαν ήταν πολύ μικρός, και θα μπορούσε να συντελέσει σε σημαντικές διαφορές για τα ξεχωριστά ποσοστά διαχωρισμού.

Για την εκτέλεση του k-NN για διαφορετικά k δεν παρατηρήθηκαν σημαντικές διαφορές, κάτι που είναι πολύ θετικό, καθώς υποδεικνύει ότι τα δεδομένα μας δεν είναι ευαίσθητα στην επιλογή του k, πιθανώς λόγω της μορφής της εισόδου του αλγορίθμου.

12.1.2 Clustering

Σε όλα τα πειράματα, ανεξάρτητα από την παραμετροποίηση του k-Means, τα σκορ της ακρίβειας ήταν αρκετά χαμηλά. Επομένως, γενικά θεωρούμε ότι η απόδοση του αλγορίθμου στο clustering των προϊόντων-chunks δεν ήταν ικανοποιητική.

Όσον αφορά στις αποστάσεις, οι καλύτερες μέθοδοι υπολογισμού ήταν και εδώ οι Euclidean, Cosine και KL-Divergence, με τις Euclidean να ξεχωρίζουν σχετικά από τις άλλες δύο. Και εδώ, τα

σκορ των KS-Test και Infinity Norm είναι ελάχιστα υψηλότερα από τα baseline σκορ, αποδεικνύοντας ότι αποτυγχάνουν να συσχετίσουν τα δεδομένα με έναν αποδοτικό τρόπο.

Παρατηρήσαμε επίσης, στο κομμάτι της αρχικοποίησης του αλγορίθμου με διαφορετικές κατηγορίες κέντρων, ότι το υψηλότερο σκορ μας έδινε η αρχική τοποθέτησή τους σε διαφορετικά clusters, και έπειτα η αρχική τοποθέτησή τους σε τυχαία. Αυτό ήταν αναμενόμενο, καθώς ο αλγόριθμος εξ ορισμού χρειάζεται να κατηγοριοποιήσει όλα τα δεδομένα σε διαφορετικά clusters, επομένως, ξεκινώντας με αυτά τα clusters σαν κέντρα, επιτυγχάνουμε πιο ακριβές αποτέλεσμα. Αντιθέτως, όταν τα αρχικά κέντρα βρίσκονταν όλα στο ίδιο cluster, χρειάστηκαν πολλές επαναλήψεις για να απομακρυνθούμε από αυτό, να υπολογίσουμε σταδιακά τα υπόλοιπα και να αναθέσουμε τα δεδομένα σωστά σε αυτά.

12.2 Μελλοντική Εργασία

Εδώ παρουσιάζουμε εν συντομία τα ζητήματα που δεν είχαμε το χρόνο να ερευνήσουμε κατά τη διάρκεια της διπλωματικής εργασίας, και επίσης θέματα που ανέκυψαν κατά την έρευνά μας, και θα μπορούσαν να εξετασθούν περαιτέρω.

- **Υπολογισμός Αποστάσεων**

Επιπλέον μέθοδοι υπολογισμού θα μπορούσαν να χρησιμοποιηθούν για την απόσταση των πινάκων μετάβασης. Θα μπορούσαμε να ερευνήσουμε μεθόδους εξειδικευμένες στην εύρεση απόστασης ανάμεσα σε διδιάστατους πίνακες, ή και να μετατρέψουμε με διαφορετικό τρόπο τον πίνακα σε διάνυσμα (παραδείγματος χάρη, ανά στήλη, αντί για ανά σειρά όπως πραγματοποιείται τώρα).

- **Αλγόριθμοι Clustering** Η χαμηλή απόδοση του k-Means θα μπορούσε να είναι ένα ισχυρό κίνητρο να πειραματιστούμε και με άλλους αλγορίθμους clustering, που μπορεί να αποδειχθούν πιο κατάλληλοι για τα δεδομένα μας.

- **Classification με τη χρήση διανύσματος συχότητας and N-διάστατου πίνακα**

Οι αλγόριθμοι Classification μπορούν να τροποποιηθούν ώστε να λαμβάνουν διαφορετικές παραμέτρους σαν είσοδο. Αντί να κατηγοριοποιήσουμε τα δεδομένα με βάση τους κωδικούς των μηνυμάτων, και τον πίνακα μετάβασης αυτών, μπορούμε να υπολογίσουμε ένα απλό διάνυσμα συχότητας για κάθε αντικείμενο, το οποίο να απεικονίζει το ποσοστό εμφάνισης καθενός από τα 45 μηνύματα στο συγκεκριμένο προϊόν. Αυτό το διάνυσμα είναι που θα δοθεί σαν είσοδος στον αλγόριθμο, και οι αποστάσεις μεταξύ των διανυσμάτων θα υπολογισθούν όπως και πριν. Με τον τρόπο αυτό αξιολογούμε την απόδοση κάθε αλγορίθμου, για όλες τις μεθόδους υπολογισμού της απόστασης, με ένα μονοδιάστατο πίνακα, που προφανώς μας παρέχει λιγότερες πληροφορίες από τον αντίστοιχο διδιάστατο.

Με τον ίδιο τρόπο, θα μπορούσαμε να γενικεύσουμε τη μέθοδο και σε 3 ή N διαστάσεις, παρέχοντας έτσι αυξημένη πληροφορία στον αλγόριθμο, και αναμένοντας διαφορετικά αποτελέσματα και αποδοτικότητα της διαδικασίας.

- **Επιλογή διαφορετικών μεταβλητών**

Στην εργασία αυτή, όλες οι υλοποιήσεις μας βασίστηκαν στον κωδικό μηνύματος κάθε αντικειμένου chunk, όπως εξηγήθηκε και στο Κεφάλαιο που αφορούσε στην προ-επεξεργασία των δεδομένων. Σε μια άλλη προσέγγιση, όλες οι μεταβλητές που χαρακτηρίζουν τη διαδικασία μας (Θερμοκρασία, Πίεση, Ομογενοποίηση και ούτω καθεξής) θα μπορούσαν να ληφθούν υπόψη κατά την εκτέλεση των αλγορίθμων. Κάτι τέτοιο θα απαιτούσε διαφορετικό διαχωρισμό των αρχικών δεδομένων, και την εύρεση ενός τρόπου να ενσωματώσουμε τις επιπλέον μεταβλητές σε γνωρίσματα των προϊόντων. Η είσοδος του κάθε αλγορίθμου θα χρειαζόταν επίσης τροποποίηση, μια και οι νέες μεταβλητές είναι αριθμητικές. Το ίδιο ισχύει και για τις μεθόδους υπολογισμού της απόστασης. Όλα τα πειράματα θα ήταν απαραίτητο να επανασχεδιαστούν,

ώστε να εξετάζουν πλέον την επίδραση κάθε μεταβλητής ξεχωριστά, κρατώντας τις υπόλοιπες σταθερές, ή πραγματοποιώντας κατάλληλους συνδυασμούς για την εύρεση συσχετίσεων.

- **Βαθμολόγηση Παραγωγικής Διαδικασίας** Τα αντικείμενα chunks αναφέρονται σε ξεχωριστά προϊόντα που παράγονται. Κατά τη διαδικασία τοποθέτησης ετικέτων, παρατηρήσαμε ότι η παραγωγή ίδιων προϊόντων δεν ήταν μια καθόλα ίδια διαδικασία, όσον αφορά στη διάρκεια, στη σειρά και στο είδος των ενεργειών που πραγματοποιούνταν. Αυτό σημαίνει ότι είχαμε προϊόντα με τα ίδια χαρακτηριστικά (Product Cleaning Group, Product Code, Product Group κλπ), αλλά με διαφορετικούς πίνακες μεταβάσεων, καθώς τα μηνύματα και η σειρά τους στο καθένα ήταν ανόμοια.

Για να αντιμετωπίσουμε το παραπάνω πρόβλημα, χρειάζεται μια καταγραφή της ιδανικής-βέλτιστης διαδικασίας, για όλα τα προϊόντα, έτσι ώστε να συγκρίνουμε το πραγματικό με αυτό και να εξαλείψουμε τις διαφορές. Για κάθε προϊόν, θα γίνεται μια καταγραφή της ιδανικής σειράς ενεργειών, και θα δημιουργούμε το αντικείμενο chunk για αυτή. Έπειτα, η σύγκριση των πινάκων μεταβάσεων αυτών, με εκείνους από τα πραγματικά προϊόντα θα δίνουν την απόδοση κάθε διαδικασίας, και θα μας επιτρέπουν να βαθμολογούμε τα προϊόντα και να αποφασίζουμε κατά πόσο προσεγγίζουν τα ιδανικά. Αυτό αποτελεί ένα πολύ σημαντικό πρώτο βήμα στην κατεύθυνση βελτιστοποίησης της παραγωγής, ιδιαίτερα όσον αφορά στην ποιοτική αναβάθμιση των προϊόντων και την ελαχιστοποίηση του χρόνου παραγωγής.

- **Μέθοδοι Αξιολόγησης** Στην αξιολόγηση των αποτελεσμάτων, χρησιμοποιήσαμε κάποιες ενδεικτικές μεθόδους, που εφαρμόζονται συχνά σε τέτοιες διαδικασίες. Υπάρχουν όμως πολλές περισσότερες, οι οποίες μάλιστα παρέχουν και διαφορετικές συσχετίσεις και συμπεράσματα. Δύο παραδείγματα είναι η Silhouette Coefficient, η οποία μετρά τη συνοχή μέσα σε ένα cluster ή μια class, και το Student's T- Test, το οποίο υπολογίζει την πιθανότητα δύο διαφορετικά σύνολα δεδομένων να είναι δείγματα της ίδιας κατανομής. Με τον τρόπο αυτό, η αξιολόγηση των αποτελεσμάτων θα μπορούσε να μας οδηγήσει σε ακόμα πιο ενδιαφέροντα συμπεράσματα.

Παράρτημα Α

Scripts

A.1 Classification

A.1.1 Nearest Centroid Classifier

```
1 # Nearest Centroid Performance
2
3
4 attribute_list = ["pr_cl_group", "pr_group"]
5 split_percentages = [0.8, 0.65, 0.5]
6 num_iters = 30
7 dist_methods = ["eucl_total", \
8                 "cosine_vector", \
9                 "KL_divergence_diag"]
10
11 for attribute in attribute_list:
12     for split in split_percentages:
13         for iteration in range(num_iters):
14
15             (train_set, test_set) = split_random(chunk_list, attribute, split)
16             centers = calculate_centers(train_set, attribute)
17
18             for dist_method in dist_methods:
19                 for chunk in test_set:
20                     assign_to_nearest_center(chunk, centers, dist_method)
21
22             (accuracy, kappa) = evaluation(test_set, attribute)
```

Listing A.1: Nearest Centroid Classifier

A.1.2 K - Nearest Neighbors Classifier

```
1 # Classification Algorithm K- Nearest Neighbors
2
3
4 attribute_list = ["pr_cl_group", "pr_group"]
5 split_percentages = [0.8, 0.65, 0.5]
6 num_iters = 30
7 dist_methods = ["eucl_total", \
8                 "cosine_vector", \
9                 "KL_divergence_diag"]
10 unique_pr_cl_gr = 5
```

```

11 unique_pr_gr      = 12
12
13
14 for attribute in attribute_list:
15
16     if attribute=="pr_cl_group":
17         min_k = unique_pr_cl_gr+1
18         max_k = unique_pr_cl_gr*3+1
19     elif attribute=="pr_group":
20         min_k = unique_pr_gr+1
21         max_k = unique_pr_gr*3+1
22
23     for split in split_percentages:
24         for iteration in range(num_iters):
25
26             (train_set,test_set) = split_random(chunk_list,attribute,split)
27
28             for dist_method in dist_methods:
29
30                 for chunk in test_set:
31                     distances = calculate_distances(chunk,train_set,dist_method)
32                     sorted_distances = sort(distances)
33
34                     for k in range(min_k,max_k+1):
35                         assign_to_most_frequent(chunk,k,sorted_distances[:k])
36
37                     for k in range(min_k,max_k+1):
38                         (accuracy,kappa) = evaluation(test_set,k,attribute)

```

Listing A.2: K - Nearest Neighbors Classifier

A.2 Clustering

A.2.1 K - Means Implementation

```

1 # Clustering Algorithm, Implementation of k-means
2
3
4 attribute_list      = ["pr_cl_group", "pr_group"]
5 centers_sets_types = [centers_sets_all_random,\
6                       centers_sets_all_diff,\
7                       centers_sets_all_same]
8 num_sets           = 100
9 dist_methods       = ["eucl_total",\
10                      "eucl_row",\
11                      "eucl_column"]
12
13 for attribute in attribute_list:
14     import_centers_sets(num_sets,attribute)
15     for centers_set_type in centers_sets_types:
16
17         for centers in centers_set_type:
18
19             for dist_method in dist_methods:

```



```

20     assign_to_nearest_center(chunk_list,centers,dist_method)
21     evaluate_clustering_baseline(chunk_list,attribute)
22     centers_have_changed = True
23
24
25     while centers_have_changed:
26         old_centers = centers
27         centers = calculate_clusters_centers(chunk_list)
28         assign_to_nearest_center(chunk_list,centers,dist_method)
29
30         if (centers == old_centers) :
31             centers_have_changed = False
32
33     evaluate_clustering(chunk_list,attribute)

```

Listing A.3: Clustering Implementation K-means

A.3 Distances

A.3.1 2D to Vector Conversion

```

1 #Convert 2D matrix to Vector using the Vector algorithm
2
3 def m2v_vector(matrix):
4     ret_vect = []
5     [ret_vect.extend(row) for row in matrix]
6     return ret_vect

```

Listing A.4: Vector algorithm

```

1 #Convert 2D matrix to Vector using the diagonal algorithm
2
3 def m2v_diag(matrix):
4     ret_vect = []
5     for i in range(len(matrix)):
6         for j in range(len(matrix[0])):
7             if (i<j):
8                 ret_vect.append((matrix[i][j]+matrix[j][i])/2.0)
9             elif (i==j):
10                ret_vect.append(matrix[i][j])
11     return ret_vect

```

Listing A.5: Diagonal algorithm

A.3.2 Distance Algorithms Implementation

```

1 #Euclidean Total distance algorithm
2 from scipy.spatial.distance import euclidean

```

```

3
4 def dist_eucl_total(chunk1, chunk2):
5     return euclidean(m2v_vector(chunk1.TM), m2v_vector(chunk2.TM))

```

Listing A.6: Distance Euclidean Total

```

1 #Euclidean Total distance algorithm
2 from scipy.spatial.distance import euclidean
3
4 def dist_eucl_row(chunk1, chunk2):
5     tt_sum=0
6     for i in range(len(chunk1.TM)):
7         tt_sum+=euclidean(chunk1.TM[i], chunk2.TM[i])
8     return tt_sum

```

Listing A.7: Distance Euclidean Row

```

1 #Euclidean Total distance algorithm
2 from scipy.spatial.distance import euclidean
3 import numpy as np
4
5 def dist_eucl_column(chunk1, chunk2):
6     tt_sum=0
7     for i in range(len(chunk1.TM)):
8         t_sum = 0
9         for j in range(len(chunk1.TM[0])):
10            t_sum+= np.square(chunk1.TM[j][i]-chunk2.TM[j][i])
11            tt_sum+=np.sqrt(t_sum)
12    return tt_sum

```

Listing A.8: Distance Euclidean Column

```

1 from scipy.spatial.distance import cosine
2 import numpy as np
3
4 def dist_cosine_average(chunk1, chunk2):
5     results = []
6     for i in range(len(chunk1.TM)):
7         cosrow = cosine(chunk1.TM[i], chunk2.TM[i])
8         if not np.isnan(cosrow):
9             results.append(cosrow)
10    return np.average(results)

```

Listing A.9: Distance Cosine Average

```

1 from scipy.spatial.distance import cosine
2
3 def dist_cosine_vector(chunk1, chunk2):

```

```
4 return cosine(m2v_vector(chunk1.TM), m2v_vector(chunk2.TM))
```

Listing A.10: Distance Cosine Vector

```
1 from scipy.spatial.distance import cosine
2
3 def dist_cosine_diag(chunk1, chunk2):
4     return cosine(m2v_diag(chunk1.TM), m2v_diag(chunk2.TM))
```

Listing A.11: Distance Cosine Diagonal

```
1 from scipy.stats import entropy
2 import numpy as np
3
4 def dist_KL_divergence_average(chunk1, chunk2):
5     results = []
6     for i in range(len(chunk1.TM)):
7         tmp_ch_1 = chunk1.TM[i]
8         tmp_ch_2 = chunk2.TM[i]
9         for j in range(len(tmp_ch_1)):
10            if tmp_ch_1[j] == 0 :
11                tmp_ch_1[j] = 0.001
12            if tmp_ch_2[j] == 0 :
13                tmp_ch_2[j] = 0.001
14            results.append(entropy(tmp_ch_1, tmp_ch_2))
15 return np.average(results)
```

Listing A.12: Distance Kullback–Leibler Divergence Average

```
1 from scipy.stats import entropy
2
3 def dist_KL_divergence_vector(chunk1, chunk2):
4     tmp_ch_1 = m2v_vector(chunk1.TM)
5     tmp_ch_2 = m2v_vector(chunk2.TM)
6
7     for i in range(len(tmp_ch_1)):
8         if tmp_ch_1[i] == 0 :
9             tmp_ch_1[i] = 0.001
10            if tmp_ch_2[i] == 0 :
11                tmp_ch_2[i] = 0.001
12
13 return entropy(tmp_ch_1, tmp_ch_2)
```

Listing A.13: Distance Kullback–Leibler Divergence Vector

```
1 from scipy.stats import entropy
2
3 def dist_KL_divergence_diag(chunk1, chunk2):
```

```

4 tmp_ch_1 = m2v_diag(chunk1.TM)
5 tmp_ch_2 = m2v_diag(chunk2.TM)
6
7 for i in range(len(tmp_ch_1)):
8     if tmp_ch_1[i] == 0 :
9         tmp_ch_1[i] = 0.001
10    if tmp_ch_2[i] == 0 :
11        tmp_ch_2[i] = 0.001
12
13 return entropy(tmp_ch_1, tmp_ch_2)

```

Listing A.14: Distance Kullback–Leibler Divergence Diagonal

```

1 from scipy.stats import ks_2samp
2 import numpy as np
3
4 def dist_KS_test_average(chunk1, chunk2):
5     results = []
6     for i in range(len(chunk1.TM)):
7         Y = ks_2samp(chunk1.TM[i], chunk2.TM[i])
8         results.append(Y[0]/Y[1])
9     return np.average(results)

```

Listing A.15: Distance Kolmogorov–Smirnov test Average

```

1 from scipy.stats import ks_2samp
2
3 def dist_KS_test_vector(chunk1, chunk2):
4     Y = ks_2samp(m2v_vector(chunk1.TM), m2v_vector(chunk2.TM))
5     return Y[0]/Y[1]

```

Listing A.16: Distance Kolmogorov–Smirnov test Vector

```

1 from scipy.stats import ks_2samp
2
3 def dist_KS_test_diag(chunk1, chunk2):
4     Y = ks_2samp(m2v_diag(chunk1.TM), m2v_diag(chunk2.TM))
5     return Y[0]/Y[1]

```

Listing A.17: Distance Kolmogorov–Smirnov test Diagonal

```

1 import numpy as np
2
3 def dist_norm_inf(chunk1, chunk2):
4     return np.linalg.norm(np.array(chunk1.TM) - np.array(chunk2.TM), np.inf)

```

Listing A.18: Infinity Norm

A.3.3 Evaluation of Distance Methods in Classification

```
1 # Classification Algorithm to compare differnt distance methods
2
3
4 attribute_list = ["pr_cl_group", "pr_group"]
5 split_percentages = [0.8]
6 num_iters = 30
7 dist_methods = ["eucl_total", \
8                 "eucl_row", \
9                 "eucl_column", \
10                "cosine_vector", \
11                "cosine_diag", \
12                "cosine_average", \
13                "KL_divergence_vector", \
14                "KL_divergence_diag", \
15                "KL_divergence_average", \
16                "KS_test_vector", \
17                "KS_test_diag", \
18                "KS_test_average", \
19                "norm_inf"]
20
21 for attribute in attribute_list:
22     for split in split_percentages:
23         for iteration in range(num_iters):
24
25             (train_set, test_set) = split_random(chunk_list, attribute, split)
26             centers = calculate_centers(train_set, attribute)
27
28             for dist_method in dist_methods:
29                 for chunk in test_set:
30                     assign_to_nearest_center(chunk, centers, dist_method)
31
32             (accuracy, kappa) = evaluation(test_set, attribute)
```

Listing A.19: Classification Distance Methods Evaluation

A.3.4 Evaluation of Distance Methods in Clustering

```
1 # Clustering Algorithm to compare differnt distance methods
2
3
4 attribute_list = ["pr_cl_group", "pr_group"]
5 centers_sets_types = [centers_sets_all_diff, \
6                      centers_sets_all_same]
7 num_sets = 20
8 dist_methods = ["eucl_total", \
9                 "eucl_row", \
10                "eucl_column", \
11                "cosine_vector", \
12                "cosine_diag", \
13                "cosine_average", \
14                "KL_divergence_vector", \
15                "KL_divergence_diag", \
```

```

16         "KL_divergence_average", \
17         "KS_test_vector", \
18         "KS_test_diag", \
19         "KS_test_average", \
20         "norm_inf"]
21
22 for attribute in attribute_list:
23     import_centers_sets(num_sets, attribute)
24     for centers_set_type in centers_sets_types:
25
26         for centers in centers_set_type:
27
28             for dist_method in dist_methods:
29
30                 assign_to_nearest_center(chunk_list, centers, dist_method)
31                 evaluate_clustering_baseline(chunk_list, attribute)
32                 centers_have_changed = True
33
34                 while centers_have_changed:
35                     old_centers = centers
36                     centers = calculate_clusters_centers(chunk_list)
37                     assign_to_nearest_center(chunk_list, centers, dist_method)
38
39                 if (centers == old_centers) :
40                     centers_have_changed = False
41
42     evaluate_clustering(chunk_list, attribute)

```

Listing A.20: Clustering Distance Methods Evaluation

Bibliography

- [1] C.C Taylor D. Michie D.J. Spiegelhalter. Machine Learning, Neural and Statistical Classification. Αγγλικά. Ellis Horwood, Φεβ. 1994, σσ. 9–11, 14–16, 20.
URL: <http://www1.maths.leeds.ac.uk/~charles/statlog/whole.pdf>.
- [2] Stuart J. Russell και Peter Norvig. Artificial Intelligence: A Modern Approach. Αγγλικά. Prentice-Hall, Inc, 1995, σσ. 26–28.
- [3] Rokach Lior Maimon Oded. Data Mining and Knowledge Discovery Handbook. Αγγλικά. Second. Springer US, 2005.
- [4] Julia Hirschberg Andrew Rosenberg. «Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning».
Στο: V-Measure: A conditional entropy-based external cluster evaluation measure.
Lecture Notes in Artificial Intelligence.
Prague: Association for Computational Linguistics, Ιούν. 2007.
- [5] S. B. Kotsiantis.
Supervised Machine Learning: A Review of Classification Techniques. Informatica 31:249–268. 2007.
- [6] Elena Deza Michel Marie Deza. Encyclopedia of Distances. Αγγλικά. Springer, 2009.
- [7] Petra Perner. «Industrial Conference on Data Mining (ICDM 2009)».
Στο: Advances in Data Mining Applications and Theoretical Aspects.
Lecture Notes in Artificial Intelligence. Leipzig, Germany: Springer-Verlag, Ιούλ. 2009.
- [8] Ethem Alpaydm. Introduction to Machine Learning. Αγγλικά. Second.
The Adaptive Computation and Machine Learning. The MIT Press, 2010.
- [9] Morven Leese Brian S. Everitt Sabine Landau. Cluster Analysis. Αγγλικά. Fifth.
WILEY SERIES IN PROBABILITY AND STATISTICS. Wiley, 2011.
- [10] Prof. Kimito Funatsu. Knowledge-Oriented Applications in Data Mining. Αγγλικά.
InTech, Ιαν. 2011. URL: <http://www.intechopen.com/books/knowledge-oriented-applications-in-data-mining/data-mining-industrialapplications>.
- [11] Mark A. Hall Ian H. Witten Eibe Frank.
Data Mining: Pactical Machine Learning Toold and Techniques. Αγγλικά. Third.
Morgan Kauffman, Elsevier, 2011, σσ. 3–9.
- [12] URL: <http://www.pmean.com/definitions/kappa.htm>.

Κατάλογος σχημάτων

2.1	Initial Data-Set format in CSV format	22
2.2	Example MsgNumbers	23
2.3	Product Cleaning Group Distribution	24
2.4	Product Group Distribution	24
2.5	Visualization Tool - Landing Page	28
2.6	Visualization Tool - Temperature, Water - 4 days	29
2.7	Visualization Tool - Temperature, Water - 1 day	29
2.8	Visualization Tool - Four Variables Correlation	30
2.9	Visualization Tool - Raw Materials Addition	30
2.10	Visualization Tool - Action Message Presentation	31
2.11	Visualization Tool - Variable Values Presentation	31
2.12	Visualization Tool - Values and Set Points	32
5.1	Distance Comparison for Product Cleaning Group Classification	49
5.2	Distance Comparison for Product Group Classification	50
5.3	Nearest Centroid Classifier for Product Cleaning Group	51
5.4	Nearest Centroid Classifier for Product Group	52
5.5	k-Nearest Neighbors Classifier for Product Cleaning Group	52
5.6	k-Nearest Neighbors Classifier for Product Cleaning Group	52
5.7	K - NN for Different k in Product Cleaning Group	54
5.8	K - NN for Different k in Product Group	55
5.9	k-means Clustering in Product Cleaning Group (V-Measure)	56
5.10	k-means Clustering in Product Cleaning Group (Rand Index)	57
5.11	k-means Clustering in Product Group (V-Measure)	57
5.12	k-means Clustering in Product Group (Rand Index)	58
5.13	K means Clustering(Product Cleaning Group)	59
5.14	K means Clustering(Product Group)	59
8.1	Αρχική μορφή του αρχείου δεδομένων CSV	71
8.2	Παραδείγματα MsgNumbers	71
8.3	Κατανομή του Product Cleaning Group	72
8.4	Κατανομή του Product Group	72
8.5	Εργαλείο Οπτικοποίησης - Landing Page	77
8.6	Εργαλείο Οπτικοποίησης - Temperature, Water - 4 days	77
8.7	Εργαλείο Οπτικοποίησης - Temperature, Water - 1 day	78
8.8	Εργαλείο Οπτικοποίησης - Four Variables Correlation	78
8.9	Εργαλείο Οπτικοποίησης - Raw Materials Addition	79
8.10	Εργαλείο Οπτικοποίησης - Action Message Presentation	79
8.11	Εργαλείο Οπτικοποίησης - Variable Values Presentation	80
8.12	Εργαλείο Οπτικοποίησης - Values and Set Points	80
11.1	Σύγκριση Μεθόδων Απόστασης για Classification στο Product Cleaning Group	96
11.2	Σύγκριση Μεθόδων Απόστασης για Classification στο Product Group	97

11.3 Nearest Centroid Classifier για Product Cleaning Group	98
11.4 Nearest Centroid Classifier για Product Group	99
11.5 k-Nearest Neighbors Classifier για Product Cleaning Group	99
11.6 k-Nearest Neighbors Classifier για Product Cleaning Group	99
11.7 k-NN για διαφορετικά k στο Product Cleaning Group	101
11.8 k-NN για διαφορετικά k στο Product Group	101
11.9 k-Means Clustering στο Product Cleaning Group (V-Measure)	103
11.10k-Means Clustering στο Product Cleaning Group (Rand Index	103
11.11 k-Means Clustering στο Product Group (V Measure	104
11.12k-Means Clustering στο Product Group (Rand Index	104
11.13k-Means Clustering(Product Cleaning Group)	105
11.14k-Means Clustering(Product Group)	106