

## **Abstract**

Banks deal with huge amounts of customer's data and thus needs tremendous efforts to improve the understanding of the accumulated data in order to detect customer's behavior and accordingly enable the executive managers to make the right decision and avoid any possible losses, wasting time and effort .

The main aim of this thesis is to distinguish between borrowers who pay back loan from those who don't . therefore the executive managers can easily reduce the costs of non-payment borrowers and decrease the high number of bad loans in order to serve the bank and its customers by using data mining techniques.

The dataset of this research was obtained from the UCI machine learning repository website. In order to improve the accuracy of our classification and gain useful results some preprocessing techniques were applied such as : removed any irrelevant and correlated data , implemented data discretization ,data cleaning ,and target class balancing as well to achieve a suitable dataset for our Algorithms. Then five data mining classification techniques were conducted which are: Naive bayes , J48, IBK, Multilayer Perceptron (MLP) and Sequential minimal optimization (SMO).The Weka software from Waikato university with (10-cross validation) was used to model and validate the proposed models.

Experiments in this research were conducted in two stages. Firstly, J48 classifier was applied on full dataset, the results carried out in this stage show that: applying of the preprocessing techniques on the data set improved the performance of the classifier. Secondly, five classification techniques were applied to the preprocessed datasets. The results carried out in this stage showed that the performance of the five classification algorithms are nearly same . Out of these five classification algorithms, J48classifier had the highest accuracy (84.35%) .

## **مستخلص البحث**

تعامل البنوك مع كميات هائلة من بيانات العملاء وهذا يحتاج إلى جهود هائلة لتحسين فهم البيانات المتراكمة ، لاكتشاف سلوك العملاء وبالتالي تمكن المدراء من اتخاذ القرار الصحيح ، وتجنب أي خسائر محتملة ، من إضاعة الوقت والجهد.

الهدف الرئيسي من هذه الرسالة هو التمييز بين أصحاب القروض المسترددة وأصحاب القروض الهاكلة مما يمكن المدراء من تقليل تكاليف أصحاب القروض غير المسترددة وبالتالي العدد الكبير من القروض المعدومة من أجل خدمة البنك وعملائه باستخدام تقنيات تنقيب البيانات.

تم الحصول على مجموعة البيانات لهذا البحث من موقع (UCI). من أجل تحسين دقة تصنيفنا والحصول على نتائج مفيدة ، تم تطبيق بعض تقنيات المعالجة المسبقة مثل: إزالة أي بيانات غير ذات صلة ومتراقبة ، وتغيير البيانات المنفصلة ، وتنظيف البيانات ، وموازنة الطبقات المستهدفة أيضًا لتحقيق مجموعة بيانات مناسبة لخوارزمياتنا. ثم تم إجراء خمس تقنيات لتصنيف بيانات التعدين وهي Bayes Naive و J48 و IBK و (MLP) Multilayer Perceptron و (SMO) Sequential minimal optimization. تم استخدام برنامج Weka من جامعة Waikato مع استخدام الدالة (10-fold validation) لعمل التماذج المقرحة والتحقق من صحتها.

أجريت التجارب في هذا البحث على مرحلتين. في المرحلة الأولى ، تم تطبيق التصنيف J48 على مجموعة البيانات الكاملة من البيانات ، وأظهرت النتائج التي أجريت في هذه المرحلة ما يلي: أن تطبيق تقنيات المعالجة المسبقة على مجموعة البيانات قد حسن أداء المصنف. في المرحلة الثانية ، تم تطبيق خمسة أساليب تصنيف على مجموعات البيانات التي تم تجهيزها مسبقاً. أظهرت النتائج التي أجريت في هذه المرحلة أن أداء خوارزميات التصنيف الخمسة مماثل تقريباً. من خوارزميات التصنيف الخمسة هذه ، كان التصنيف J48 يتمتع بأعلى دقة (84.35٪).

## **Dedications**

*This dissertation is lovingly dedicated to :*

*My mother*

*A strong and gentle soul who taught me to trust in ALLAH , believe in  
hard work and that so much could be done with little*

*My father*

*for earning an honest living for us and for being my first teacher who  
support and encourage me to believe in myself*

## **Acknowledgement**

*Thanks first and foremost to God Almighty who honored me in accomplishing this work, and I would like to express my sincere gratitude to Sudan University of Science and Technology Collage of Graduate Studies for letting me fulfill my dream of being a student , my supervisor, Dr. Shaza Mirghani , I am extremely grateful for her assistances and suggestions throughout this research .I would also like to thank my parents and friends who helped and supported me a lot in finalizing this research.*

## Table of Contents

CHAPTER ONE: INTRODUCTION .....	1
1.1 Overview .....	1
1.2 Background.....	1
1.3 Problem statement.....	1
1.4 Research significance.....	2
1.5 Research Questions .....	2
1.6 Research objectives.....	2
1.7 Research scope.....	2
1.8 Thesis Structure .....	2
CHAPTER TWO: LITERATURE REVIEW.....	3
2.1 Introduction .....	3
2.2 The Knowledge Discovery .....	3
2.3 Data Mining .....	4
2.3.1 Classification.....	5
2.3.2 Prediction .....	5
2.4 Data Mining Applications in Banking .....	6
2.4 .1 Risk Management in Banks .....	6
2.5 Previous studies and related works.....	7
2.5 .1 Developing Prediction Model for banking Loans Risk in Banks Using Data Mining techniques.....	7
2.6 Summary of the related works.....	8
CHAPTER THREE: METHODOLOGY .....	11
3.1 Introduction .....	11
3.2 Data Set Description .....	11
3.3 Data Preprocessing.....	11
3.4 Classification Using Weka .....	13
3.4.1 Naïve Bayes .....	14
3.4.2. IBK .....	15
3.4.3 J48.....	16
3.4.4 Multilayer Perceptron (MLP).....	17
3.4.5 Support Vector Machine (SVM) .....	18

3.5 Evaluation of classifiers .....	19
3.6 The Experiments .....	20
FOUR: RESULT CHAPTER ANALYSIS AND DISCUSSION.....	21
4.1 Introduction .....	21
4.2 The first Experiments .....	21
4.3 The second Experiments .....	24
CHAPTER FIVE: CONCLUSION AND RECOMMENDATION.....	31
5.1 Introduction .....	31
5.2 Conclusion.....	31
5.3 RECOMMENDATION .....	31
References:.....	32
Appendix A.1 .....	34
Appendix A.2 .....	35

## List of Tables

Table 2. 1 Summary of the related works.....	8
Table 3. 1 An overview of the dataset.....	11
Table 4. 1 Confusion Matrix J48 Classifier with Raw data set .....	21
Table 4. 2 detailed Accuracy by Class J48 Classifier with Raw data set .....	21
Table 4. 3 Confusion Matrix J48 Classifier with preprocessed data set .....	22
Table 4. 4 detailed accuracy by class J48 classifier with preprocessed data set.....	22
Table 4. 5: the evaluation of J48 classifier on different datasets with Cross-validation mode. ....	23
Table 4. 6 Confusion Matrix Naive bayes.....	24
Table 4. 7 Detailed Accuracy by Class Naive bayes.....	25
Table 4. 8 Confusion Matrix J48 .....	25
Table 4. 9 Detailed Accuracy by Class using J48.....	25
Table 4. 10 Confusion Matrix IBK .....	26
Table 4. 11 Detailed Accuracy by Class using IBK.....	26
Table 4. 12 Confusion Matrix SMO.....	27
Table 4. 13 Detailed Accuracy by Class using SMO .....	27
Table 4. 14 Confusion Matrix MLP .....	28
Table 4. 15 Detailed Accuracy by Class using MLP .....	28
Table 4. 16 Evaluation of classifiers on credit card dataset with Cross validation.....	29

## List of Figures

Figure 2. 1 Steps of the KDD Process.....	4
Figure 2. 2: Data mining model and tasks .....	5
Figure 3. 1: the preprocessed dataset .....	13
Figure 3. 2: (Naïve Bayes) classifier with preprocessed dataset.....	15
Figure 3. 3: (IBK) classifier with preprocessed dataset .....	16
Figure 3. 4: (J48) classifier with preprocessed dataset.....	17
Figure 3. 5: (MLP) classifier with preprocessed dataset.....	18
Figure 3. 6: ( SMO) classifier with preprocessed dataset .....	19
Figure 4. 1detailed accuracy by class using J48 classifier with full data set.....	22
Figure 4. 2 detailed accuracy by class J48 classifier with preprocessed data set .....	23
Figure 4. 3comparing between two data sets.....	24
Figure 4. 4 Detailed Accuracy by Class Naive bayes .....	25
Figure 4. 5 Detailed Accuracy by Class using J48.....	26
Figure 4. 6 Detailed Accuracy by Class using IBK .....	27
Figure 4. 7 Detailed Accuracy by Class using SMO.....	28
Figure 4. 8 Detailed Accuracy by Class using MLP .....	29
Figure 4. 9 comparing between the classifiers .....	30

