

DECLARATION

I **PERCY PAULIN J (38110403)** hereby declare that the Project Report entitled **“PREDICTION OF DIABETES USING DATA SCIENCE TECHNIQUE”** done by me under the guidance of **Dr. S. L. JANY SHABU M.Phil., Ph.D.**, and Sathyabama Institute of Science and Technology is submitted in partial fulfillment of the requirements for the award of Bachelor of Engineering degree in Computer Science and Engineering.

DATE:

PLACE:

SIGNATURE OF THE CANDIDATE

ACKNOWLEDGEMENT

I am pleased to acknowledge my sincere thanks to **Board of Management of SATHYABAMA INSTITUTE OF SCIENCE OF SCIENCE AND TECHNOLOGY** for their kind encouragement in doing this project and for completing it successfully. I am grateful to them.

I convey my thanks to **Dr. T.SASIKALA M.E, Ph.D., Dean, School of Computing** and **Dr. S.VIGNESHWARI M.E., Ph.D., and Dr. L.LAKSHMANAN, M.E., Ph.D., Head of the Department, Dept. of Computer Science and Engineering** for providing me necessary support and details at the right time during the progressive reviews.

I would like to express my sincere and deep sense of gratitude to my Project Guide **Dr. S. L. JANY SHABU M.Phil., Ph.D.**, for his valuable guidance, suggestions and constant encouragement paved way for the successful completion of my project work.

I wish to express my thanks to all Teaching and Non-teaching staff members of the **Department of Computer Science and Engineering** who were helpful in many ways for the completion of the project.

ABSTRACT

Diabetes is a chronic (long-lasting) health condition that affects how your body turns food into energy. Most of the food you eat is broken down into sugar (also called glucose) and released into your bloodstream. When your blood sugar goes up, it signals your pancreas to release insulin. Nowadays machine learning is applied to healthcare system where there is a chance of predicting the disease early. The main necessity of Artificial intelligence is data. The past dataset is collected and that dataset is used to build a machine learning model. The necessary pre-processing techniques are applied like univariate analysis and bivariate analysis are implemented. The data is visualised for better understanding of the features and based on that a classification model is built by using machine learning algorithm and comparison of algorithms are done based on their performance metrics like accuracy, F1 score recall etc.

TABLE OF CONTENTS

| CHAPTER NO | TITLE | PAGE NO |
|------------|--------------------------------|---------|
| | ABSTRACT | V |
| | LIST OF FIGURES | VIII |
| | LIST OF ABBREVIATION | IX |
| 1 | INTRODUCTION | 1 |
| 1 | 1.1 Data Science | |
| 2 | 1.1.1 Data Scientist | |
| 3 | 1.2 Artificial Intelligence | 3 |
| 3 | 1.3 Machine learning | |
| 4 | 1.4 Overview | |
| 5 | 2 LITERATURE SURVEY | |
| 8 | 3 AIM AND SCOPE OF THE PROJECT | |
| 8 | 3.1 Aim | |
| 8 | 3.2 Scope | |
| 8 | 3.3 Objective | |

| | | |
|-----------|-------------------------------------|-----------|
| | 3.4 Proposed system | |
| 8 | | |
| | 3.4.1 Advantage | |
| 9 | | |
| 4 | ALGORITHM AND METHODS | |
| 10 | | |
| | 4.1 Model description | |
| 10 | | |
| | 4.1.1 Data Pre-processing | |
| 10 | | |
| | 4.1.2 Data Visualization | |
| 11 | | |
| | 4.1.3 Comparing Algorithm | |
| 12 | | |
| | 4.1.4 Deployment using Flask | |
| 33 | | |
| | vi | |
| | 4.2 Requirement and technology used | 34 |
| | 4.2.1 Functional requirement | |
| 34 | | |
| | 4.2.2 Non Functional requirement | |
| 34 | | |
| | 4.2.3 Environmental requirement | |
| 35 | | |
| 5 | CONCLUSION AND FUTURE WORK | 36 |
| | 5.1 Conclusion | |
| 36 | | |
| | 5.2 Application | |
| 36 | | |

| | |
|----|------------------|
| 37 | 5.2 Future work |
| 38 | REFERENCE |
| 39 | APPENDIX |
| 39 | A. Sample Code |
| 56 | B. Screen short |

LIST OF FIGURES

FIGURE NO

TITLE

PAGE NO

| | | | |
|----|------|--|----|
| | 3.1 | System Architecture | 9 |
| 10 | 4.1 | After inputting all the library and data | |
| | 4.2 | Bar graph of Age wise positive cases | 11 |
| 12 | 4.3 | Positive and Negative cases in heat map | |
| | 4.4 | Positive and Negative predictive value of Logistic Regression. | 16 |
| | 4.5 | Accuracy result of Logistic Regression | 17 |
| 19 | 4.6 | Positive and Negative predictive value of Decision Tree. | |
| 21 | 4.7 | Accuracy result of Decision Tree. | |
| | 4.8 | Positive and Negative predictive value of Random Forest. | 23 |
| | 4.9 | Accuracy result of Random Forest. | 25 |
| | 4.10 | Positive and Negative predictive value of Support Vector Machine. | 29 |
| | 4.11 | Accuracy result of Support vector machine | 30 |
| | 5.1 | Pie chart of positive and negative cases | 36 |

LIST OF ABBREVIATION

- 1 AI - Artificial intelligence
- 2 ML - Machine Learning
- 3 SML – Supervised Machine Learning

CHAPTER 1

INTRODUCTION

1.1 DATA SCIENCE

Data science is an interdisciplinary field that uses scientific methods, processes, algorithms and systems to extract knowledge and insights from structured and unstructured data, and apply knowledge and actionable insights from data across a broad range of application domains. The term "data science" has been traced back to 1974, when Peter Naur proposed it as an alternative name for computer science. In 1996, the International Federation of Classification Societies became the first conference to specifically feature data science as a topic. However, the definition was still in flux.

The term "data science" was first coined in 2008 by D.J. Patil, and Jeff Hammerbacher, the pioneer leads of data and analytics efforts at LinkedIn and Facebook. In less than a decade, it has become one of the hottest and most trending professions in the market. Data science is the field of study that combines domain expertise, programming skills, and knowledge of mathematics and statistics to extract meaningful insights from data. Data science can be defined as a blend of mathematics, business acumen, tools, algorithms and machine learning techniques, all of which help us in finding out the hidden insights or patterns from raw data which can be of major use in the formation of big business decisions.

1.1.1 Data Scientist:

Data scientists examine which questions need answering and where to find the related data. They have business acumen and analytical skills as well as the ability to mine, clean, and present data. Businesses use data scientists to source, manage, and analyse large amounts of unstructured data.

1.2 ARTIFICIAL INTELLIGENCE

Artificial intelligence (AI) refers to the simulation of human intelligence in machines that are programmed to think like humans and mimic their actions. The term may also be applied to any machine that exhibits traits associated with a human mind such as learning and problem-solving. Artificial intelligence (AI) is intelligence demonstrated by machines, as opposed to the natural intelligence displayed by humans or animals. Leading AI textbooks define the field as the study of "intelligent agents" any system that perceives its environment and takes actions that maximize its chance of achieving its goals.

Learning processes. This aspect of AI programming focuses on acquiring data and creating rules for how to turn the data into actionable information. The rules, which are called algorithms, provide computing devices with step-by-step instructions for how to complete a specific task.

Reasoning processes. This aspect of AI programming focuses on choosing the right algorithm to reach a desired outcome.

Self-correction processes. This aspect of AI programming is designed to continually fine-tune algorithms and ensure they provide the most accurate results possible.

1.3 MACHINE LEARNING

Machine learning is to predict the future from past data. Machine learning (ML) is a type of artificial intelligence (AI) that provides computers with the ability to learn without being explicitly programmed. Machine learning focuses on the development of Computer Programs that can change when exposed to new data and the basics of Machine Learning, implementation of a

simple machine learning algorithm using python. Process of training and prediction involves use of specialized algorithms. It feed the training data to an algorithm, and the algorithm uses this training data to give predictions on a new test data. Machine learning can be roughly separated in to three categories. There are supervised learning, unsupervised learning and reinforcement learning. Supervised learning program is both given the input data and the corresponding labeling to learn data has to be labeled by a human being beforehand. Unsupervised learning is no labels. It provided to the learning algorithm.

This algorithm has to figure out the clustering of the input data. Classification is the process of predicting the class of given data points. Classes are sometimes called as targets/ labels or categories. Classification predictive modeling is the task of approximating a mapping function from input variables(X) to discrete output variables(y). In machine learning and statistics, classification is a supervised learning approach in which the computer program learns from the data input given to it and then uses this learning to classify new observation. This data set may simply be bi-class (like identifying whether the person is male or female or that the mail is spam or non-spam) or it may be multi-class too. Some examples of classification problems are: speech recognition, handwriting recognition, bio metric identification, document classification etc.

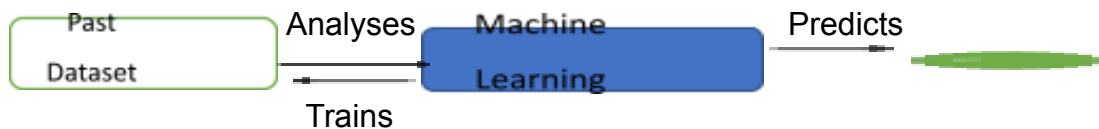


Fig: Process of Machine learning

Supervised Machine Learning is the majority of practical machine learning uses supervised learning. Supervised learning is where have input variables (X) and an output variable (y) and use an algorithm to learn the mapping function from the input to the output is $y = f(X)$. The goal is to approximate the mapping function so well