

Acknowledgments

I thank my advisor Dr. Saeed Ul Hassan for being the supportive advisor. His precious insights, guidance and support throughout the research journey has not only helped me become a better student but also examine data for research and contribute to previously done research. Though he is very busy, he always made me feel like a priority. He took time to review my findings, answer questions, and give thoughtful feedback for which I am extremely grateful. This thesis would not have been possible without them.

I thank Dr. Mohsin Ali for his help, advice and support throughout my work and for agreeing to serve on my thesis committee.

Moreover, I thank my family especially my parents and my friends for me letting me follow my dreams and for their encouragement which has been the vital part of the research.

Table of Contents

Chapter	Title	Page
	Acknowledgments	i
	Table of Contents	ii
	List of Figures	iv
	List of Tables	v
	Abstract	vi
1	Introduction	1
	1.1 Overview	1
	1.2 ICD, HIPPA and Complications	2
	1.3 Objectives	2
	1.4 Research Methodology	3
	1.5 Contribution	3
	1.6 Thesis Outline	4
2	Literature Review	5
	2.1 Detailed Overview of Literature Review	5
3	MIMIC-III Critical Care Database	8
	3.1 MIMIC-III	8
	3.2 MIMIC-III Construction	9
	3.3 MIMIC-III derived Concepts	10
4	Methodology	14
	4.1 System Overview	14

	4.2 System Design	20
5	Experimental Results	27
	5.1 Statistics	27
	5.2 Results	27
6	Conclusion and Future Investigation	32
	6.1 Conclusion	32
	6.2 Future Investigation	33
6	References	34
	References	34
7	Appendices	37

List of Figures

Figure	Title	Page
1.1	Methodology.	4
3.1	MIMIC-III Construction Model.	9
4.1	Technical Chain of Steps.	15
4.2	Batch Processing for Huge Files Using Python.	16
4.3	One Hot Encoding for Categorical Features	18
4.4	SMOTE (Sampling) .	19
4.5	Sigmoid Function.	22
4.6	Logistic Vs Linear SVM	23
4.7	Artificial Neural Network (Perceptron Basic Model)	24
4.8	Artificial Neural Network Model Configuration	25
4.9	Training and Validation Error	25
4.10	Training and Validation Loss	26
5.1	Patient Age Distribution.	28
5.2	Patient Length of Stay Distribution.	28
5.3	Insurance Types Distribution.	29
5.4	AUC for RF.	29
5.5	AUC for LR.	30
5.6	AUC for SVC.	30

List of Tables

Table	Title	Page
3.1	MIMIC-III Tables Summary	12
3.2	Class Distribution of data for MIMIC-III Dataset	13
3.3	Derived Concepts	13
4.1	Diagnoses and Procedures Count.	17
4.2	Class Count.	18
5.1	Results with ADASYN Data Sampling.	29
5.2	Results with SMOTE Data Sampling.	31
5.3	Complications vs Expiry.	31

Abstract

Analyzing the health of patients using Electronic Health Records (EHR) and take precise decisions is an important problem in healthcare research. To predict outcomes and take such decisions patients admitted in ICU requires specific features of medical records: worth, capacity, access and dimensionality. A substantial amount of data needs to be stored in proper infrastructure and comprehensive analysis which can aid doctors, clinicians, medical experts and families is required. In this Thesis, an investigation of medical data is carried out. We propose an ETL approach to extract features and train machine learning models to predict the post-procedural complications seeking if those complications can lead to further complications or mortality. To derive insights for that, we used known dataset named as Medical Information Mart for Intensive Care III (MIMIC-III) with two types of data sampling techniques: SMOTE and ADASYN. For both techniques our results showed that Random Forrest outperforms other linear models with an accuracy of $> 80\%$ and AUROC of 0.83.

Chapter 1

Introduction

1.1 Overview

With advancements in digital technologies in medical sciences, different techniques making it possible to explore big data precisely and predict tasks of clinicians, labs and doctors using patient records of diagnoses or procedures. Intensive care unit (ICU) is where severely ill patients are admitted and requires accurate predictors that can help doctors and take respective measures. Diagnostic and medical technologies have evolved rapidly and doctors have to take complex decisions. In the book [1] it is being said that unfortunately, majority of clinical decisions are not based on evidence. According to 2012 Institute of Medical Committee Report, a small percentage of decision were evidence based. This problem ranges to not have proper clinical guidelines. In this thesis we investigate the different methodologies for extracting, transforming and ETL techniques [2] [3], which obtain data from original source to perform informative analysis and features extraction to aid model to predict post-procedural (diagnoses and procedure) complications of critically ill patients and investigating those complications if those can lead to mortality of patient or not. The methods use demographics, data from different hospital system, lab events, diagnoses, notes and other engineered information regarding each patient.

The database used for the study is Medical Information Mart for Intensive Care MIMIC-III [4] which is de-identified data and abide by protecting health information (PHI). Other researches about MIMIC-III data is also presented to motivate our problem, establish understanding of dataset, key findings and recommendations for future investigations. The question of predicting post-procedural complications from data science perspective and critical health perspective is not only important for doctors, administrators but also for the patient as well. For administrators this would help managing patients and required resources. Avoiding predicted complications can further be avoided if such information is known during the stay of patient at ICU.

On MIMIC-III researchers widely contributed to support the cause. For any severely sick patient obesity is not so important feature to consider. [5] Presented an investigation of BMI against obesity. They conjecture that illness scores would give results differently if BMI was considered. Their setting included documented weight and height. The assessment was based on 184402 from 184 different ICU's across United States and assessment showed that 4% were classified as under-weight, 30% as normal and same for over-weight. Apart from these 28% were mentioned as obese and 8% as morbidly obese. To further explore the obesity paradox [10] presented a study recently where they evaluated the relation between BMI with survival and ensured the gender interactions. The design of their study was retrospective with total cohort samples N=139 from June 2014 to September 2016.

The main results were Survival (OS) and (PFS). Analysis was performed using Random Forests. Their findings showed that the paradox of obesity exists under overweight in the real world. The result showed that sarcopenia body mass composition can be used as predictors.

A study from Journal of intensive care medicine [6] having the design of retrospective study and setting of single tertiary academic medical center examine the impact of overstay of patients and discharge delays on in hospital mortality and morbidity. For the interventions, for all patients, from the bed request in ward to discharge time was calculated. Created bins for greater and less than 24 hours discharge delays. To find out the relationship between delays and ICU outcome they used multivariate linear regression and logistic regression. There was no such association between the long delays and subsequent mortality. Many other researcher have contributed in prediction of other critical factors related to medicine and the health of patient. A study from 2018 [7] with a design of retrospective study examined that severity scores may also lead to misclassification of critically ill obese patients using the weight of patient and the obesity score having documented results of their diagnoses prior to hospital admissions. They further compared the laboratory results between the normal patients having less obesity scores and high obesity score. They showed that an unconventionality is found within the scores of creatinine, white blood cells and urea nitrogen in blood and compared between them. These results could also be improved and used for the mortality and also in predicting other important clinical tasks.

1.2 ICD, HIPAA and Complications

The International Classification of Diseases (ICD) [8] is the foundation which identifies and evaluates the statistics of diseases globally. Creates standards for all the diagnoses, diseases which further can be used for research purposes. Under revision of ICD9 codes, the code 996 defines complications particular to certain specified procedures and diagnoses. Most complications are caused due to cardiac, vascular or other used devices and some of them relates to reaction caused due to a procedure performed. In our work, we are focusing on such complications and investigating if those can lead to the mortality of patient. HIPAA (Health Insurance Portability and Accountability Act of 1996) [9] provides privacy and protection of health records for safeguarding and protecting (PHI) medical information. To protect health information MIMIC-III provided anonymized data, still we need to make sure we are following HIPAA compliance rules so that our research does not conflict with any of the standards defined.

1.3 Objectives

A predictor that can evaluate complications using the EHR can help hospital administrators, experts, families and doctors to take decisions that are required for the safety of the patient. For any hospital

if the management already knows the precise expected stay of patient they can plan better to use their resources to provide the best possible care. For doctors predictors can evidence, statistics and labels which can be used for taking valuable decisions, for families who will pay for the expenses against the patient can better plan their billing using such predictors.

Extended diagnoses and patients stay at the hospital due to complications caused during the stay at hospital or after is associated with not only the health of patient, cost, increased number of deaths but also increased number of readmissions. Each of these parameters defines the hospital performance. So, our objective is to produce insights that can complement these parameters ranging from cost to patient health. To be specific, our objective is to extract features against not so rare diagnoses and procedures with defined ETL process, balancing of our dataset using distributions and sampling techniques and predict 996 (complication or no complication) (In general). Furthermore, to derive insights that can determine whether those complications can lead to death or not.

1.4 Research Methodology

As open sourced MIMIC-III data is available for research purposes for free. Our main focus was to use tools and technologies in a way that can complement already done research. Hence, we used Python, R and PostgreSQL for ETL process, features selection and modeling our features to extract results which I will explain in respective chapters in details. We have used the basic Data Science (mining) model as full cycle of development in our research.

The model divides several phases which are executed in iteration to get to goal which we can see in Figure 1.1 In the first phase data extraction and understanding is built, then comes the phase of data processing, handling the information as per required business transformations. The next step is the integration of all the transformations into one source to further use by next step of feature space extraction and model building. The later part of process includes models evaluations and optimization.

1.5 Contribution

The prediction of complications is subjective to type of complication, data availability, resources used by management to record EHR. Against each type of diagnoses, subclasses of complications are already defined and regularly updated by ICD. Domain experts have contributed to individual complications defined by ICD, but we generalized all the complications and implemented models to predict a general class of whether a patient will have complication or not. Moreover, one must consider the existence of unstructured data available in MIMIC-III which lead to extensive search for variables without access to actual system deployed at hospitals. Because of that a considerable part was dedicated to the ETL and a novel approach was implemented to extract meaningful features which can

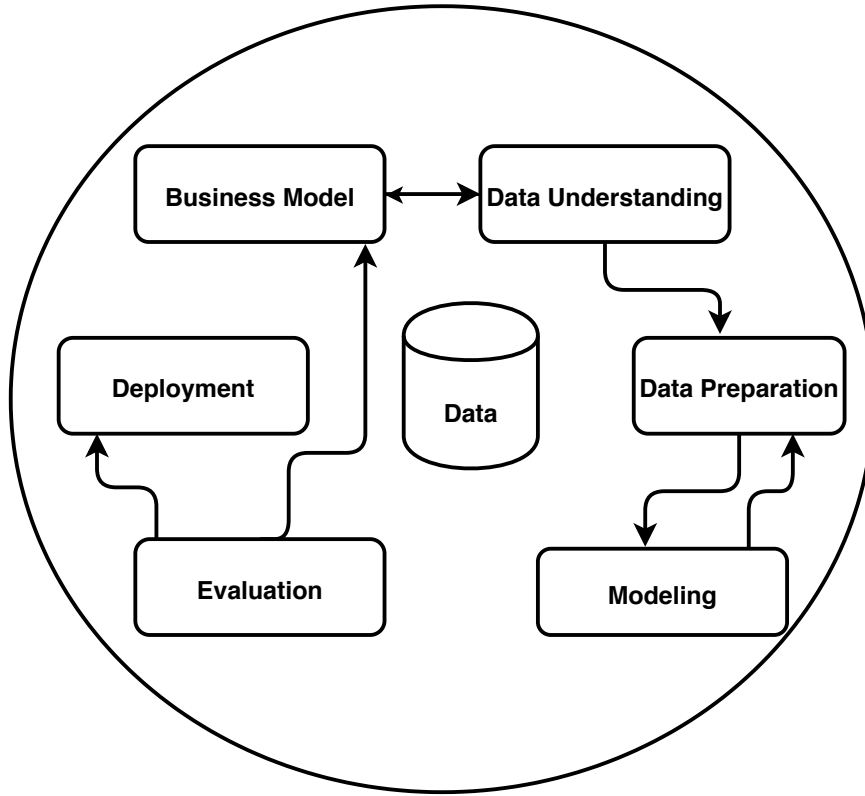


Figure 1.1: Methodology

lead to reliable prediction models. We used PostgreSQL for the SQL implementation and Python as it provide large number of libraries and models to implement. To be specific, we implemented an ETL approach, generalized all the complications to one class, implemented models to predict that class. (If 1, patient will have complication else no complication)

1.6 Thesis Outline

I organize the rest of this dissertation as follows.

In Chapter 2, we reviewed the research and studies already done by researchers and contributors.

In Chapter 4, I propose my methodology.

In Chapter 5, I present the experimental results.

Finally, in Chapter 6, I conclude my thesis.

Chapter 2

Literature Review

2.1 Detailed Overview of Literature Review

Health and medicine are one of the key sectors that requires use of new technologies to produce new possibilities and cause a greater impact in the society. Some of the recent research studies mentioned above have been conducted on MIMIC-III data set. Numerous researches have also been conducted on the subject matter and overall on the usage of MIMIC-III for creating new possibilities of research and scientific areas. Following are some of the researches that helped us motivate our problem and contribute by applying information retrieval and data science techniques.

Before we take the recent contributions related to predicting complications and other clinical tasks into account, let us first explain the fact that structured data may not always contain accurate and required information. As we are working with a huge medical data mart, the volume, diversity, dimensions and ETL matters. To identify the patient cohort from whole mart by searching the structured tables for diagnoses, procedures, chart events, demographics and other features is a burdensome task keeping the scoring systems [10] and authenticate the results of predictors with cost and time complexity [11]. The studies with EHR databases that have both structural and unstructured data have not only help researchers in identifying the new possibilities but also helped them to recognize patients having higher risks [12]. Several studies [13] [14] [15] [16] [17] [18] upon structured data related to ICD (International Classification of Diseases) have helped us carry our research as they extracted structured information related to patient, diagnoses and procedure and shown to have good precision and AUROC. Working with large data sources, retrieval of information is time taking and difficult when multiple data sources are involved in the ETL process [19]. From IEEE International conference of Healthcare Informatics a study [20] explored the use of SVM on assigning the ICD-9 code using the clinical notes.

Specific to MIMIC-III, from learning about the relationship between healthcare processes [21], the study [22] which made us realize that laboratory research is based on information retrieval irrespective of the design of study whether retrospective or prospective. All the management, doctors and clinicians are not involved in utilities of labs and how to proceed with MIMIC data and its structures. Another study [23] combined description of notes from the EHR and clinical trials using the state of the art NLP models to evaluate find relationships within trials and notes. Which, is our future plan to incorporate the textual features combining with complications to increase the performance of predicting complications accurately. LOS prediction for the patient can provide valuable information to the management of the hospital but also for patient's health. [24] Explored the NN usage for predicting the LOS prediction within time range of (<5) days or (>5) days after the patient left the Intensive care unit. They used a subset of MIMIC-III and written all their models in R and PostgreSQL using

a supercomputer provided by the Florida Polytechnic University. Their model predictions achieved 80% accuracy and outperformed any other linear models previously used of predicting the length of stay. They used a neural net of 2 layers with 5 and 3 nodes at each level. In their final dataset, they had total 31018 data points, from which, they randomly chose subsets of 15000 points and trained the NN on the Florida Polytechnic University supercomputer.

Another study in 2019 [25] mapped ICD-9 codes using the clinical notes from physicians, doctors and other staff against each patient automatically using deep learning. Their main focus was to evaluate the performance of deep learning for mapping the ICD-9 codes. Their pipeline had extraction of data from Notes and Diagnoses, removing special characters and stop words from each note. They extracted non-sequential features and applied tfidf and word2vec. They baseline models LR, RF, FRNN on non-sequential features and on sequential features from word sequences + embedding matrix they applied Conv1D, LSTM and GRU. Overall they applied multiple experiments and showed that deep learning outperforms linear models to map ICD-9 codes with 0.6957 F1 score and accuracy of 0.8967.

Research published in BMC Medical Informatics and Decision Making [26] proposed a new approach by combining rule-based features and DL model with prior knowledge for effective disease classification. They evaluated their method on i2b2 obesity challenge and demonstrated that their model outperform other method used for disease classification. They based their method on Solt's Systems [27] (The Perl's ¹ Implementation of Solt provided by author) for first identifying the major triggers terms and phrases from notes, then predicting the classes with rare terms and phrases and then trained convolutional neural network. Another mortality prediction case study published in Machine Learning for HealthCare Conference [28] demonstrated large variability in studies of mortality prediction. They reviewed the performance of the related studies based on MIMIC-III and compared it to Gradient Boosting and LR ran on extracted data.

Apart from above mentioned researches a lot of other researchers contributed in predicting the complications. A study [29] explored the model performance in predicting after operation complications following anterior cervical discectomy and fusion (ACDF). They applied Logistic Regression, Random Forrest and ANN to achieve their results. 20,879 patients and ANN outperformed LR ($p < 0.05$). Similarly, [30] used machine learning to derive and validate the hospital readmission. Their model reached (area under the receiver operating curve, 0.76). [31] A study with design of retrospective study, on predicting the real time complications used RNN against mortality, renal failure, and post-operative bleeding and operation revision. 47,559 samples of ICU were included. Their deep learning models produced PPV 0.90 and sensitivity 0.85 for mortality, 0.87 and 0.94 for renal failure, and 0.84 and 0.74 for bleeding.

Another study using MIMIC-III [32] implemented a deep rule based fuzzy systems for predicting the

¹https://github.com/yao8839836/obesity/tree/master/perl_classifier

accurate in-hospital mortality. Their main contribution to the system was proposing a system which can handle multiple categorical data.

Chapter 3

MIMIC-III Critical Care Database

This chapter explains the structure, context and development researchers have done on data source MIMIC-III.

3.1 MIMIC-III

Over the past decade, much have been written about the field of data science regarding the explosion of big data. In health care, every decision made for a critical patient requires precision by clinicians and doctors. To carry out research to aid clinicians and doctors to make better, reliable and quick decision using research and applications. This demands privacy existence of data and wide-ranging analysis.

Researcher and developer community is contributing in the MIMIC-III (Medical Information Mart for Intensive Care III) database. MIT developed MIMIC-III for research purposes and improve the already done research in the field of Bioinformatics using different techniques of data analytics and sciences. MIMIC-III is large and freely available relational DB comprising de-identified [38] health related data which is associated with over 40,000 patients stayed in ICU at Beth Israel Deaconess Medical Center (Boston, Massachusetts). The data ranges from June 2001 to October 2012.

MIMIC-III is a comprehensive collection of de-identified data from 53423 distinct critical care hospital admissions from 38597 distinct adult patients. The data has been compiled into 26 tables which contain, for example, an average of 4579 noted values and 380 lab charted measurements for each hospital admission as well as a total of 3.8 gigabytes of unstructured textual data from various health-care provider notes and analyses. In addition to de-identifying patient data, MIT requires training in the protection of patient data for anyone requesting access to the MIMIC dataset. After completing the prescribed training, data can be downloaded as 26 comma separated files (csv) files explained in Table 3.1 representing the 26 tables in the MIMIC-III database. Additionally, there is a published data dictionary which can be found at given link ¹. There is inconsistency in the usage of unique attributes and definition of primary keys between the sample SQL code and the published data dictionary. For example, every table has an attribute called ROW_ID and the sample SQL code consistently declares this attribute as unique and/or as a primary key for every table despite the fact that tables like the Patient table have a unique identifier (SUBJECT_ID) that is intended to be the primary key and serve as foreign key in child relations that refer to the PATIENTS table.

After downloading and analyzing the MIMIC source tables, implementation occurs in 5 additional steps:

¹<https://mimic.physionet.org/mimictables/admissions/>

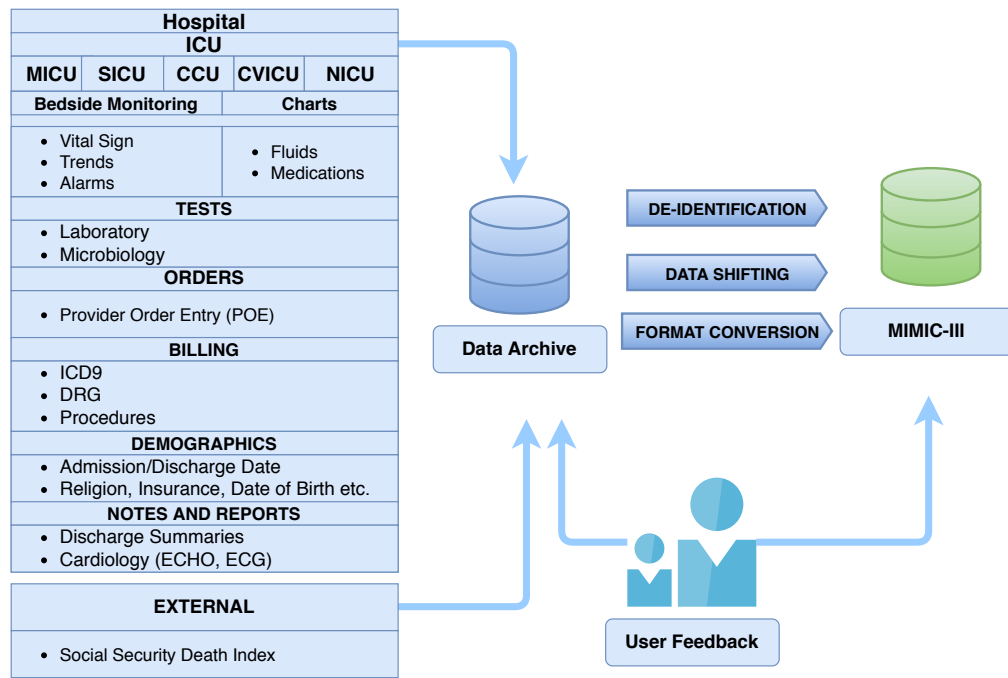


Figure 3.1: MIMIC-III Construction Model.

- Create tables with attribute rules (data types) and identify the primary key for each table. 1. Load records from csv files into each table.
- Declare the indexes for each table.
- Define foreign keys in each table and establish table relationships.
- Implement user interface (with appropriately granted permissions) for the database.

This database includes patients hospital level demographic data, ICU level laboratory test results, Labs and resources measurements, procedures, diagnoses, clinical notes, staff, notes and services. The class distribution of data shown in table 3.2.

3.2 MIMIC-III Construction

MIMIC-III was constructed based upon hospital level, patient level, ICU level & used systems level. It includes billing, notes and reports as shown in below figure 3.1.

The data collected from different ICU centers MICU, SICU, CCU, CVICU, NICU. The division of data includes bedside monitoring, charts, tests, orders, billing, demographics and notes.

The archived data then passes through the steps of transformation which include de-identification, data shifting and format conversion along with user feedback to form the MIMIC-III database.

3.3 MIMIC-III derived Concepts

The active researchers have contributed to already given data with additional scripts to generate new concepts and insights at MIMIC code repository which includes views and tables as well. They also encourage other researchers to contribute to derived insights which helps to distinct between the original data and derived data and one can use as per the problem they are solving and contribute as well. Table 3.3 shows the major concepts that are being used frequently by researchers.

Table	Summary
Admission	The admission file provides the information of each patient admission to the hospital with entry-level information
Care Givers	The information related to caregivers which examined the patients is provided in this file
Callout	It contains the information related to discharges
Chart Events	All the lab events or fluids calculation and different charted events against each patient is provided in Chart Event file
CPT Events	The CPT file includes all the procedural codes that were applied for each patient
D_CPT	This file provide a dictionary for CPT codes containing brief information
D_ICD_Procedures	This file contains ICD-9 codes for classification of procedures and their respective information
D_ICD_Diagnoses	This file contains ICD-9 codes for classification of diagnoses and their respective information
D_ITEMS	D_ITEMS provides the information of each items performed on patient
D_LabItems	D_LabItem provides the information of items related to labs
DateTimeEvents	All the date-time events are provided in DateTimeEvents file including the time of each diagnosis, procedures, items or lab event is performed
Diagnoses_ICD	This file defines the all the diagnoses performed with sequence on individual patient
Procedures_ICD	This file defines the all the procedure performed with sequence on individual patient
DRG_Code	All the codes related to drugs are provided in DRG.Code file
ICU_Stays	This file includes information of each ICU stay
InputEvents_CV	This file contains information of input events at each level of fluid or lab event for CareVue database
InputEvents_MC	This file contains information of input events at each level of fluid or lab event for MetaCare database
LabEvents	This file includes all the lab events that are performed on each patient
Microbiology Events	It contains the microbiology tests and procedures
Note Events	Contains all the notes and textual information entered by each care giver
Output Events	It contains all the output events for each lab item, diagnoses and procedure
Patients	This file includes all the demographics of the patient
Prescriptions	All the prescriptions and medication entries are included in this file
ProcedureEvents_MV	It contains all the procedures against each patient

Services	This file includes all the services previously performed and current services for each patient
----------	--

Table 3.1:: MIMIC-III Tables Summary

These tables are linked by the identifiers which mostly have the suffix ID. For example this id(HADM_ID) referred as uniquely identified hospital admission and SUBJECT_ID is refers to a unique patient. ROW_ID an exception, which is a row-identifier unique to that table. These tables are pre-fixed with D_ are dictionaries and provide definitions for identifiers. For example, in every row of OUTPUTEVENTS is related with a single ITEMID which represents the measured concept, but it does not have the actual name of the drug. By joining these two columns(OUTPUTEVENTS and D_ITEMS) on ITEMID, it is possible for the identification of what this ITEMID represents.

Class of Data	Description
Biling	It contains the coded data for administrative and billing purposes. It includes CPT codes, ICD codes, Diagnosis-related group codes.
Interventions	It includes different procedure such placement of lines, imaging studies and dialysis.
Descriptive	Information related to patient like demographic detail, discharge and admission detail and death date.
Laboratory	Test results of hematology, urine analysis, blood chemistry and microbiology.
Medications	Hospital management records of medications and medication orders.
Notes	Notes made by doctor such as patient progress and hospital discharge summaries.
Physio logic	Report generated by nurse after approximately hourly visit to patient. Reports like BP, respiratory and etc.
Reports	Patient text reports of imaging and electrocardiogram studies.
Dictionary	Look up table for different code referencing like codes with their associated laels.

Table 3.2: Class Distribution of data for MIMIC-III Dataset

Class of Data	Summary
Comorbidity	Derived scripts for specify the existance of various comorbidities using the ICD-9 codes(billing) for patient at the time of hospital discharge.
First day	Sub-folder holds different scripts which was used for calculation various clinical concepts(Highest BP, temperature, etc) of a patient on its very first day. It also contains scripts that can adapt data outside of 1st day.
Durations	Time of start and start of various treatment or duration of phenomena that includes medical agents that have a vasoactive effect on the circulatory system of a patient and others.
Sepsis	Details of sepsis, the most common mortality cause for patients in the ICU.
Severity Scores	The seriousness of illness scores which summed up the condition of a patients illness at the time of admission to ICU (usually in the rst 24 hours).
Organ Failure	This script derives binary flags for major organ failures.

Table 3.3: Derived Concepts

Chapter 4

Methodology

In this section, we introduce the ETL followed by us to derive certain insights which lead us to conclusion of stated problem. The section is divided into extraction, transformation and loading sections to reach to our features.

4.1 System Overview

4.1.1 Technical Process

Given the size of data mart and the volume of raw data, we devoted most of time to extraction and transformation of data. In the first step prior for requesting the access to MIMIC, you will have to complete the CITI “Data or Specimens Only Research” course by registering yourself on CITI program. After getting data access we are provided links to the 26 comma separated file containing patient, hospital and ICU related data¹.

Following are key steps covered in technical process to engineer features:

- Tables Creation
- Relationship Mapping (Indexes and Keys)
- Materialized views from already given tables
- Trim down values for ICD-9 Codes
- Filter rows with subject id lookup and pass it to items lookup for certain diagnoses and procedure
- ICD-9 Codes for class complications which is 996
- Making it to binary classes with 1 and all other classes to 0
- Extracting derived features from chart events and lab events with batch processing
- Consolidate all other features with derived concept
- Format all features, fill out invalid fields and normalize features for model training

¹<https://mimic.physionet.org/gettingstarted/access/>

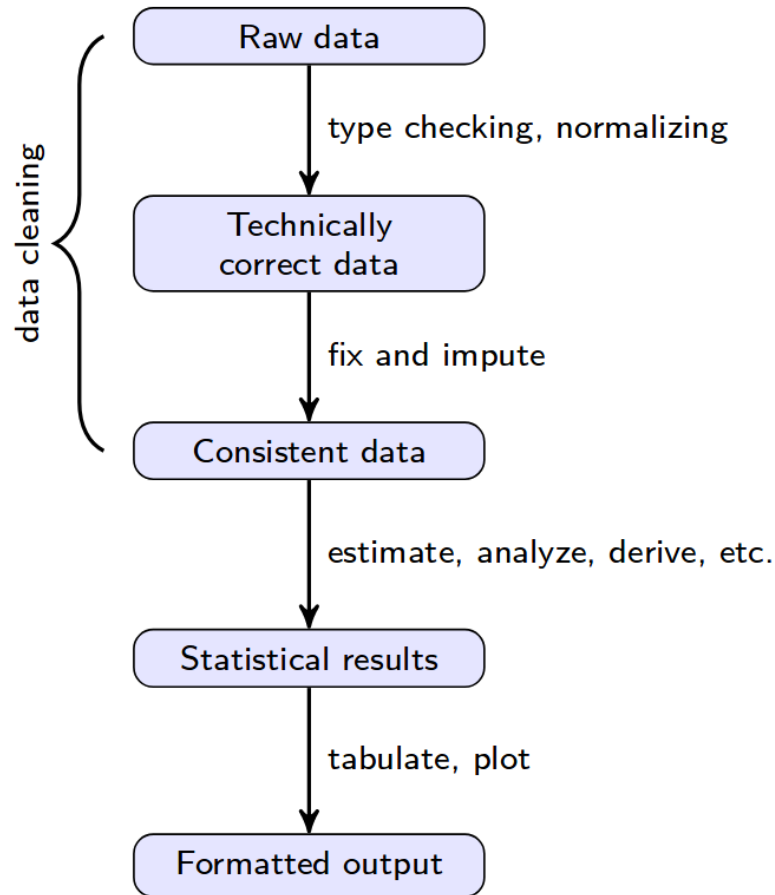


Figure 4.1: Technical Chain of Steps.

4.1.2 Relational Mapping, Batch Processing and Preprocessing

All 26 files are relationally mapped with each other². After getting these files we created a database of all those file and created respective tables. To improve the performance indexes and constraints were added. Some of them are very huge and requires pre and post processing. The smaller files were dealt with PSQL but on other hand, the large files caused problem for not only creating table but also of processing those files in RAM. To handle such problems with huge files, we implemented Python script for asyn batch processing using Pandas³ which is an open sourced library to manipulate structured data and very highly efficient because of its reliable data frame objects along with transformation tools available with it.

As the data sources and research work is now publicly available. Researchers have contributed in the

²<https://mitlcp.github.io/mimicschemaspy/>

³<https://pandas.pydata.org/>

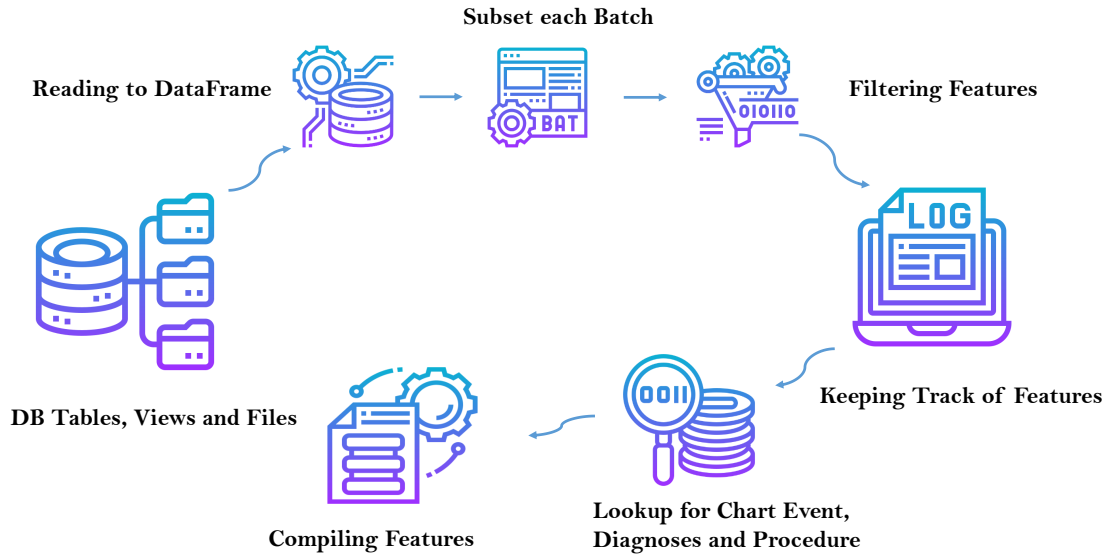


Figure 4.2: Batch Processing for Huge Files Using Python.

form of code, new concepts, and optimization of previously written script and in many other ways. Similarly we have used and created features set containing top diagnoses and procedures performed on ICU patients.

To complete ETL process, PostgreSQL and Python played important role. Multiple SQL scripts were written for creation of tables, indexes, materialized views and derived tables. All of which are presented on a public repository [19] ⁴. Extraction of major chart events and lab events against each patient involved filtering of specific patients, lookup against particular diagnoses and procedures.

4.1.3 Diagnoses

From 14328 unique diagnoses that were available in MIMIC-III, we selected those diagnoses which were more common with a threshold of happening more than 30. The number is randomly selected and experimented with. We can change number with other experiments to see if the models can outperform earlier built models. The reason behind choosing this particular number was that using these number of groups, if we use any aggregated function to fill out missing information, there will be a low probability of creating features containing near to zero variance.

⁴<https://github.com/faisalmaqbool94/Thesis-Bioinformatics-MIMICIII->

Table 4.1: Diagnoses and Procedures Count.

Item	Counts
Diagnoses	14328
Procedure	3882

4.1.4 Procedures

From 3882 unique procedures registered in MIMIC-III, we selected more common ones. Same as we did for diagnoses. For procedures there were only few fields that required any preprocessing. In the Figure 4 the subset creations and filtering involved lookups where we created separate files for segregating the subjects which are only appearing in above criteria of procedures and diagnoses.

4.1.5 Feature Space

Now that we have explored and discussed about the dataset. Now we discuss the features that are valuable to models. Once these features are identified we have to define the processes of our models and reach our goals. From potential features which can impact the predictor are included, physicians from PIMS hospital Islamabad and Islamabad Diagnostic Centre, selected and helped us engineer our features they know the important from their medical knowledge, field experience or intuition. Following are our selected variables:

- General: Insurance, Martial status, Hospital Expire Flag, Length of Stay, Calculated Bicarbonate, TotalCo2, Chloride, Free Calcium, Glucose, Hematocrit, Hemoglobin, Lactate, Oxygen, Oxygen Saturation, PCO2, PH, Potassium, Sodium, Temperature, Calcium Total, Centromere, Creatinine, Globulin, Blood Glucose, Blood Lipase, Blood Magnesium, Blood Potassium, Blood Sodium, Platelets Counts, Red Blood Cells, White Blood Cells, Lymphocytes
- Engineered Concepts (Derived from Table 3): Hypothyroidism, Renal Failure, Liver Disease, Peptic Ulcer, Aids, Lymphoma, Metastatic Cancer, Solid Tumor, Rheumatoid Arthritis, Coagulopathy, Obesity, Weight Loss, Fluid Electrolyte, Blood Loss Anemia, Alcohol Abuse, Drug Abuse, Psychoses, Depression, Congestive Heart Failure, Cardiac Arrhythmias, Valvular Disease, Pulmonary Circulation, Peripheral Vascular, Hypertension, Paralysis, Other Neurological, Chronic Pulmonary, Diabetes Uncomplicated, Diabetes Complicated

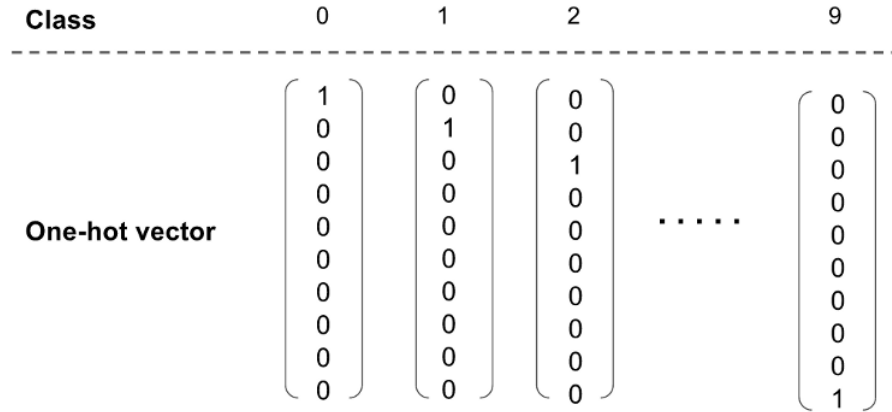


Figure 4.3: One Hot Encoding for Categorical Features.

Table 4.2: Class Count.

Class	Count
Complication	2754
No Complication	30491

4.1.6 Features Engineering

Once all the featured got extracted for certain subjects and against hospital admissions. We had to distinguish between numerical and categorical variables. For the categorical features performed One Hot Encoding technique.

Missing values got treated by the average of all the diagnoses and same goes for procedures. Average value is taken because we have extracted diagnoses and procedures which are commonly occurring and average of each group was taken. Our target variable ICD9 Code got converted into binary variable and mapped to (0, 1) where 0 indicates the non-complication and 1 indicates the occurrence of complication. All the complications are further subdivided into thousands of categories but we were just interested in the main class of complication which is indicated by code 996.

4.1.7 Sampling

As we have considered a very sensitive topic which requires a lot of domain knowledge and predicting a complication requires precision. Although we collected and engineered our dataset for targeting our goals but as the table 5 shows that we clearly have class imbalance problem. To tackle this problem we applied over sampling and down sampling techniques which are explained below.

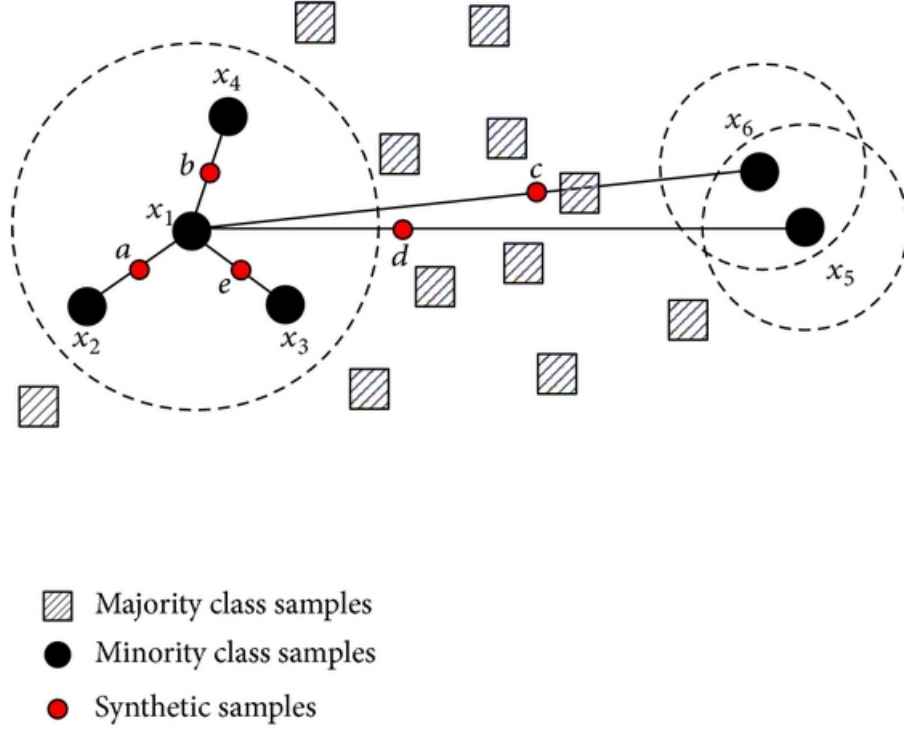


Figure 4.4: SMOTE (Sampling)

SMOTE: Synthetic Minority Over-Sampling Technique

SMOTE [33] proposed by Chawla et al., is an oversampling method. This method interpolate the minority class neighbors to construct new minority class samples randomly. The method can be described as follows. First, for each minority class sample x one gets its k nearest neighbors from other minority class samples. Second, one chooses one minority class x' sample among k neighbors. Finally, one generates the synthetic sample $x(\text{new})$ by interpolating between x and x' as follows:

$$X_{\text{new}} = x + \text{rand}(0, 1) * (x' - x) \quad (\text{Equation 4.1})$$

Where $(0, 1)$ refers to random number between $(0, 1]$. In view of geometry, SMOTE can be regarded as interpolating between two minority class samples. The decision space for the minority class is expanded that allows the classifier to have a higher prediction on unknown minority class samples. The SMOTE algorithm is simple and effective while generating synthetic samples, and the overfitting problem is avoided. It expands the decision space for the minority class.

ADASYN: Adaptive Synthetic Sampling

ADASYN [34] is another oversampling method which interpolates new minority class samples by first calculating the number of synthetic samples of minority class and for each minority sample find the nearest neighbors by calculating the Euclidean distance.

- **Input:** D with m samples with, $i = 1$ to m , where x is an n-dimensional vector in feature space and y is the corresponding class. Let $m(r)$ and $m(x)$ be the number of minority and majority class samples respectively, such that $m(r)$ less than or equal to $m(x)$ and $m(r) + m(x) = m$
- **Procedure:**
 - Calculate the Degree of Imbalance, $d = m_r/m_x$
 - If $d < d_x$ (where d_x is the preset threshold for maximum tolerated imbalance) then:
 - * Calculate the number of synthetic samples to be generated from the minority class:, $G = (m_x - m_r) * \beta$, β is the balance level of the synthetic samples generated. $\beta = 1$ means there is a total balance between two classes.
 - * For each $x_i \in$ minority samples, find the k-nearest neighbors based on Euclidean distance and calculate the ratio r_i , $r_i = \Delta_i / K$
 - * Normalize $r_x \leftarrow r_i / \sum r$, such that r_x is now a density distribution.
 - * Calculation of synthetic sample generated for each minority data point $g_i = r_x * G$ where G is the total number of synthetic data examples that need to be generated for the minority class as defined in aforementioned Equation.
 - * For each minority class data example x_i , generate g_i synthetic data examples according to the following steps:
 - Do the Loop from 1 to g_i
 - – Randomly choose one minority data example, x_u , from the K nearest neighbors for data x_i .
 - – Generate the synthetic data example: $s_i = x_i + (x_u - x_i) * \lambda$ where $(x_u - x_i)$ is the difference vector in ndimensional spaces, and λ is a random number: $\lambda \in [0, 1]$. **END Loop.**

The major difference between SMOTE and ADASYN is the difference in the generation of synthetic sample points for minority data points. In ADASYN, we consider a density distribution r which thereby decides the number of synthetic samples to be generated for a particular point, whereas in SMOTE, there is a uniform weight for all minority points.

4.2 System Design

4.2.1 Binary Classification

As we have defined our problem as binary class, either the patient would have complication or not, and defined a target feature which will be used to identify that. Depending upon the number of samples in each class, we dealt with balanced and unbalanced labeled dataset. For unbalanced dataset we applied

above mentioned two over-sampling techniques. To notion of metric performance is called accuracy which we used to validate our models defined in below section. Accuracy is defined as:

$$Accuracy = \frac{1}{n} \sum_{i=1}^n I(i_i=i'_i) \quad (\text{Equation 4.2})$$

Here i_i is the predicted class label for the i th iteration using (a defined function), the total number of samples are defined by n and the index represents the each sample. I is the indicator variable that equals one if classified correctly and zero if the result is negative. One can also define the accuracy metric as follow:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (\text{Equation 4.3})$$

Where TP and TN are positive and negative instances correctly classified by the model and FP are negative instances classified positive by model, similarly FN are positive instances classified negative by the model.

4.2.2 Selected Models

For our prediction task, we applied following models and are explained below:

- Logistic Regression (LR)
- Linear SVN (Linear Support Vector Classification)
- Random Forrest (Decision Tree)
- ANN (Artificial Neural Network)

To apply all these models, Python sklearn, imblearn, matplotlib, Pandas and Numpy libraries were used.

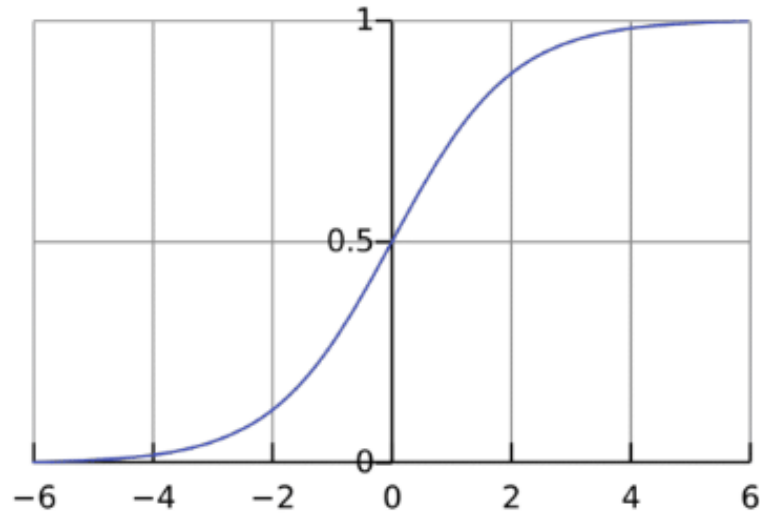


Figure 4.5: Sigmoid Function: Taken from [35]

Logistic Regression

Logistic regression is a predictor which provide the output by mapping the dependent variables and uses a logistic function. The range of the model is as follow and being continuous and differentiable. Property the sigmoid function, which we denote as follows:

$$P(x) = \frac{1}{1 + e^{-x}} \quad (\text{Equation 4.4})$$

Figure 4.5 illustrate the sigmoid function.

Random Forrest

Random forests [36] are predictors which on given labeled data create random forests for reaching towards the target of each instances and map all the features to those trees. The actual output is then generated by averaging the output of all the forests. Tin Kam Ho initially presented the concept of Random Forrests in the year of 1995. In 2001 [37] where Breimen extend the concept by extending already built algorithm and work of Amit [38]. The random forest creates many decision trees from which it calculates the actual and optimal path to find the expected label. It average the performance from all the forests, the major concepts due to which model was named Random Forrest are as follow:

- While building new forests, sample the data in random fashion
- While splitting the parent and leaf nodes, creation of subsets of multiple features

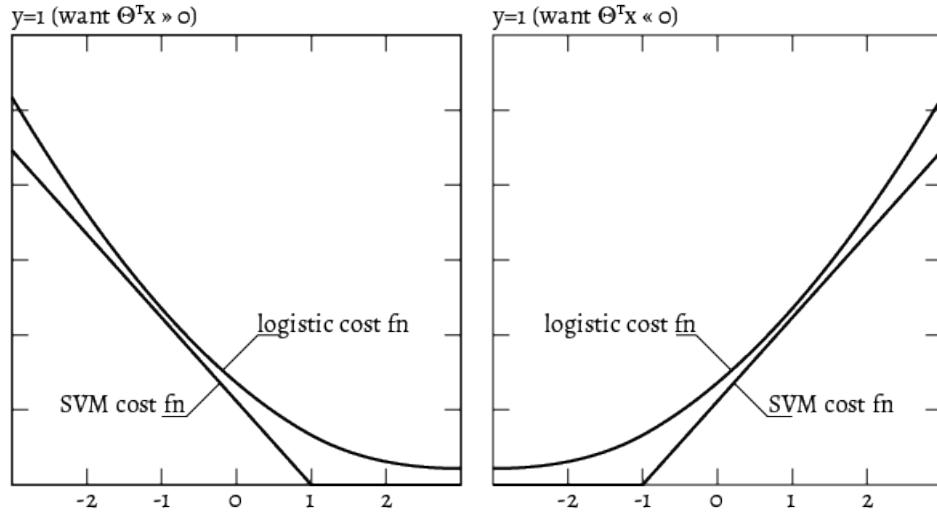


Figure 4.6: Logistic Vs Linear SVM.

Linear SVC

A linear SVC support vector classifier (SVM) a supervised predictor which given a labeled dataset finds the "best fit" hyperplane which divides and segregate the data into different classes. It has the cost function like logistic regression defined below:

$$P_{\Theta}(x) = \begin{cases} 1 & \text{if } \Theta^T x \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (\text{Equation 4.5})$$

A comparison is given with the logistic cost function in figure 4.6.

Artificial Neural Network (ANN)

A comparison is given with the logistic cost function below:

- Inputs: All the features available in the training dataset become the input for a perceptron. Also, an extra value called a bias value is fed as one of the inputs.
- Weights: The value of weights are initiated randomly (most of the times zero for all) and these values are updated accordingly by reviewing the training error
- Weighted sum: This is the summation of all the values obtained after multiplying each weight with its associated input value and adding the bias at the end.
- Activation function: These functions convert an input signal of a node to an output signal. Some of the commonly used activation functions are **tanh**, **sigmoid**, **relu** [39], **softmax**, exponential and linear. The flexibility of these activation functions is one of the reasons neural networks perform better than traditional multilinear models.

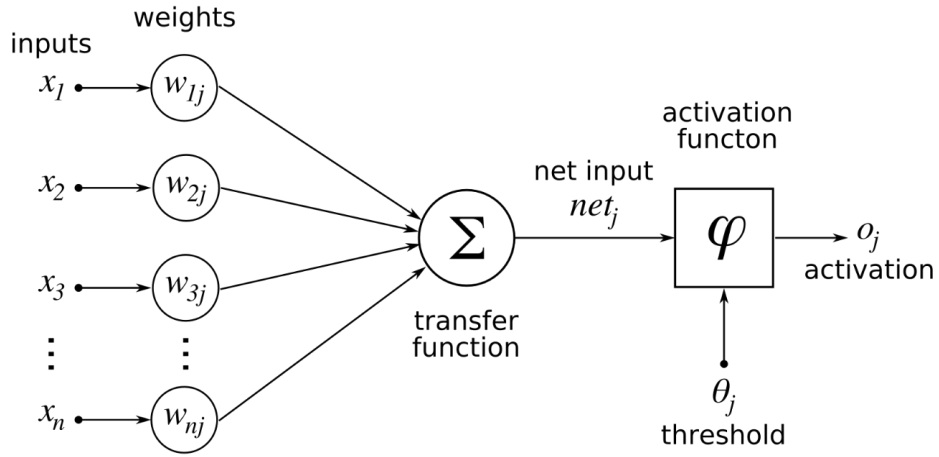


Figure 4.7: Artificial Neural Network (Perceptron Basic Model).

- **Output:** The weighted sum is passed into the activation function and becomes the input 13 value of the next layer. As a first step, the weight vector is initialized. All the features available in the training dataset are fed as input to the perceptron. These input features are then multiplied with the corresponding weights and the values are summed up including the bias value. The new computed value is fed to the activation function in order to get the predicted output. If the predicted value doesn't match with the actual value, the error is calculated and the weights are updated in order to reduce the error for the next iteration. This process is repeated until the error is reduced to a prescribed level, or if a certain number of steps is achieved.

Our ANN model in compiled form given below shows number of nodes and hidden layers. Training of our ANN model are shown in figure 4.8 and the respective plot for training and validation error is represented in Figure 4.9. Similarly the Figure 4.10 shows the training and validation loss respectively.

Model: "sequential_9"

Layer (type)	Output Shape	Param #
dense_11 (Dense)	(None, 2048)	135168
dropout_9 (Dropout)	(None, 2048)	0
dense_12 (Dense)	(None, 2048)	4196352
dropout_10 (Dropout)	(None, 2048)	0
dense_13 (Dense)	(None, 2048)	4196352
dropout_11 (Dropout)	(None, 2048)	0
dense_14 (Dense)	(None, 2048)	4196352
dropout_12 (Dropout)	(None, 2048)	0
dense_15 (Dense)	(None, 1)	2049
Total params: 12,726,273		
Trainable params: 12,726,273		
Non-trainable params: 0		
model compiled		

Figure 4.8: Artificial Neural Network Model Configuration.

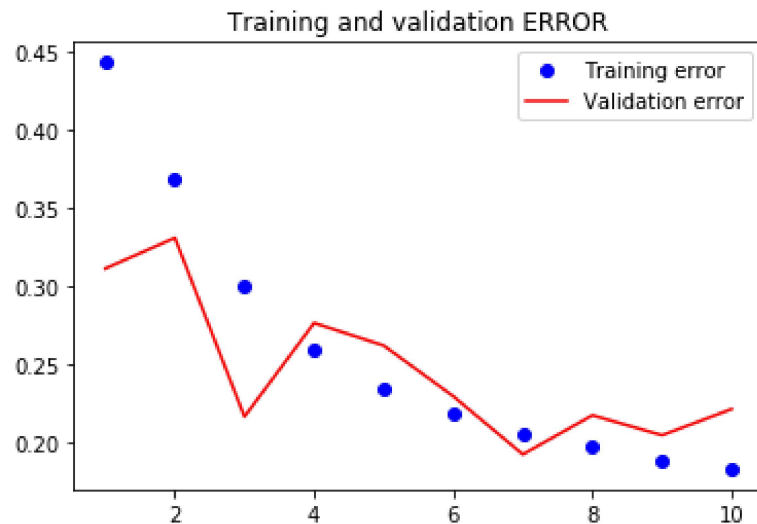


Figure 4.9: Training and Validation Error.

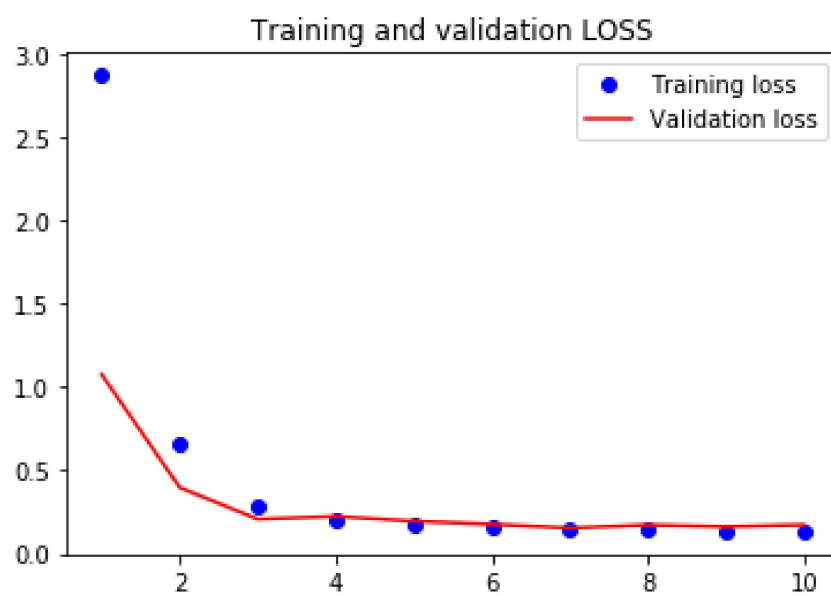


Figure 4.10: Training and Validation Loss.

Chapter 5

Experimental Results

5.1 Statistics

In this section, we present some of the basic numerical descriptors of our dataset and the results of our initial analysis. Following are some statistics from our data.

The figure 5.1 shows the patient age distribution which clearly impacted every feature of patient. Similarly, we plotted the length of stay distribution for patients. The figure 5.2 shows the distribution

As we have considered the insurance types as well. Insurance type do impact in cases of complications because the complications are not only based on certain medication only, they also are based upon certain tools, systems or an items used on patients due to which a particular complication occurred. For example, complications related to stents. The insurance covers such type of things that a patients will get what kind of services hence we did not neglect this feature. Following is the distribution on types of insurances.

5.2 Results

After the data preparation, problem statement and the decision of performance metric there were several model candidates to run. The candidate models that we selected are Logistic Regression, Linear SVC (SVM), Random Forrest and ANN.

First we standardized our data and normalized our data frames to be passed to models to predict the complications. We used standard min-max Scaler for normalization of our data. For NN we used a multilayer perceptron for Complications predictions which used ReLU as activation function. Before each model, two data sampling processes were executed. SMOTE and ADASYN as we had class imbalance problem. After getting the interpolated data we passed that dataset to each of the model to predict the complications. The different result metric resulted from each of the model against each data sampling technique is shown in below table:

The Table 5.1 shows result after applying the ADASYN data sampling technique. After that we implement SMOTE as well to see if the models vary with the type of static sampling instead of interpolating data based on the distribution as in ADAYSN

The Table 5.2 shows result with SMOTE data sampling technique. It clearly shows that RF still outperforms other models.

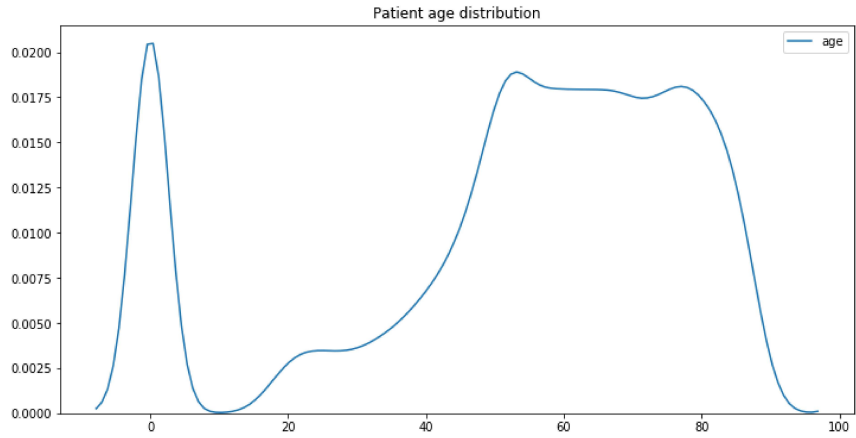


Figure 5.1: Patient Age Distribution.

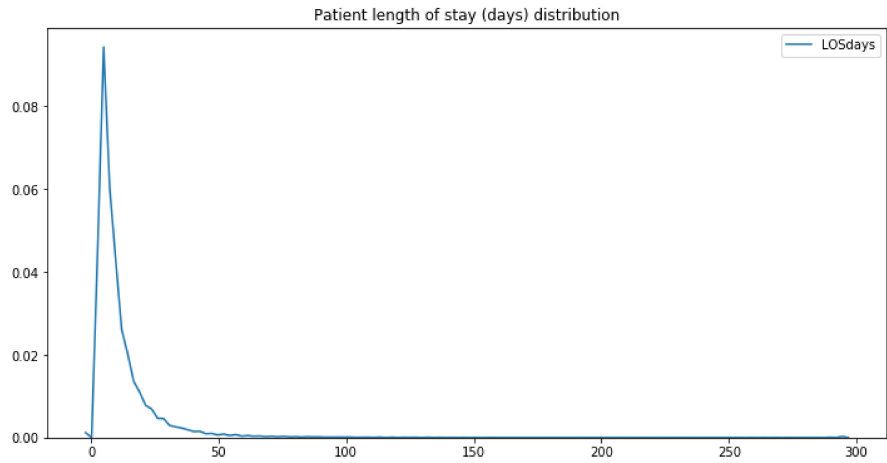


Figure 5.2: Patient Length of Stay Distribution.

We cannot evidently say that, complications lead to death but there can certainly be some parameters which leads to mortality based on complications. In our future development we will try to map those features with deep domain knowledge and expand our study to contribute more. Table 5.3 shows the distribution of expired patients against complications.

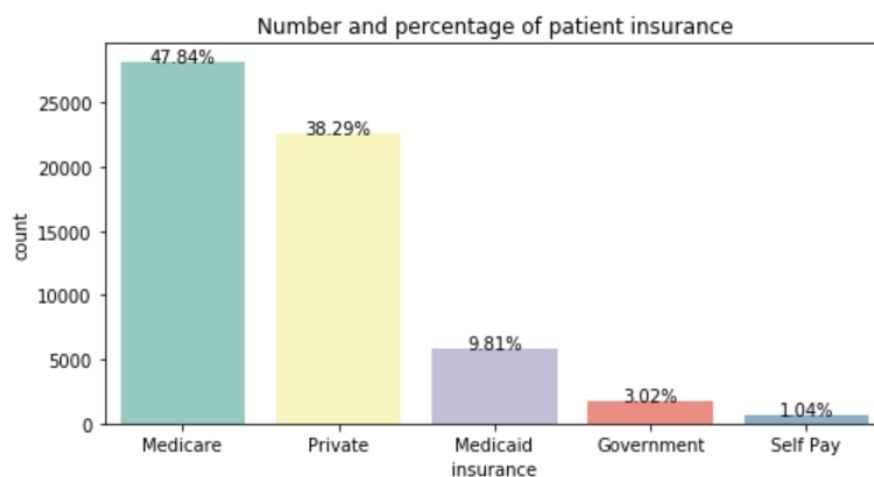


Figure 5.3: Insurance Types Distribution.

Table 5.1: Results with ADASYN Data Sampling.

Model	Accuracy	AUROC
Logistic Regression	65%	0.72
Linear SVC	66%	0.72
Random Forrest	86%	0.83
ANN	81%	0.82

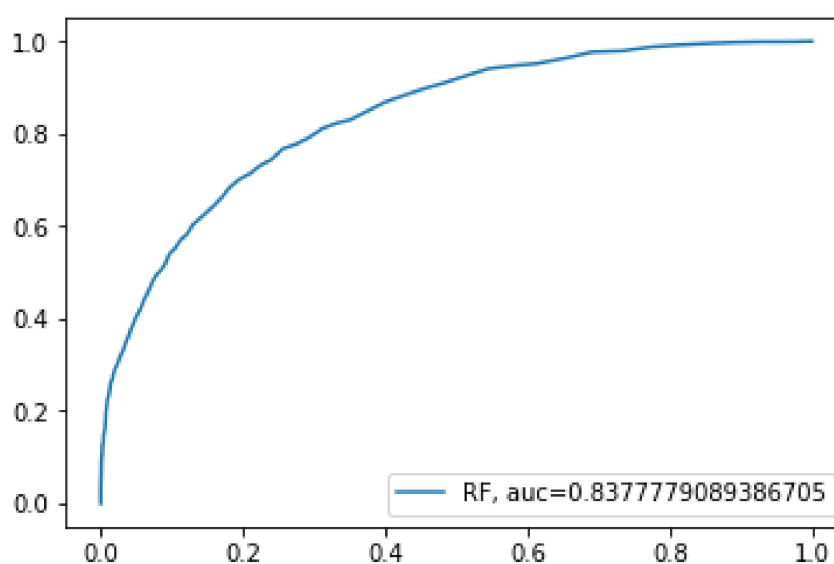


Figure 5.4: AUC for RF.

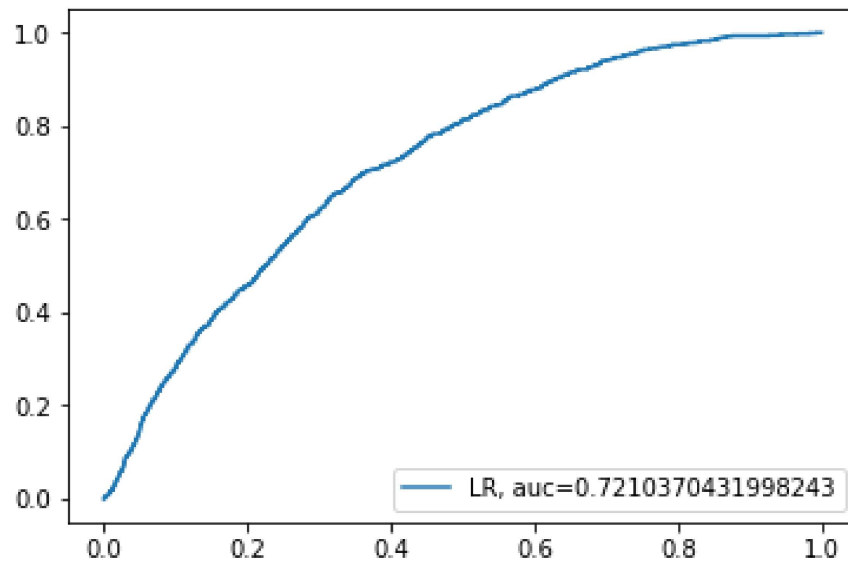


Figure 5.5: AUC for LR.

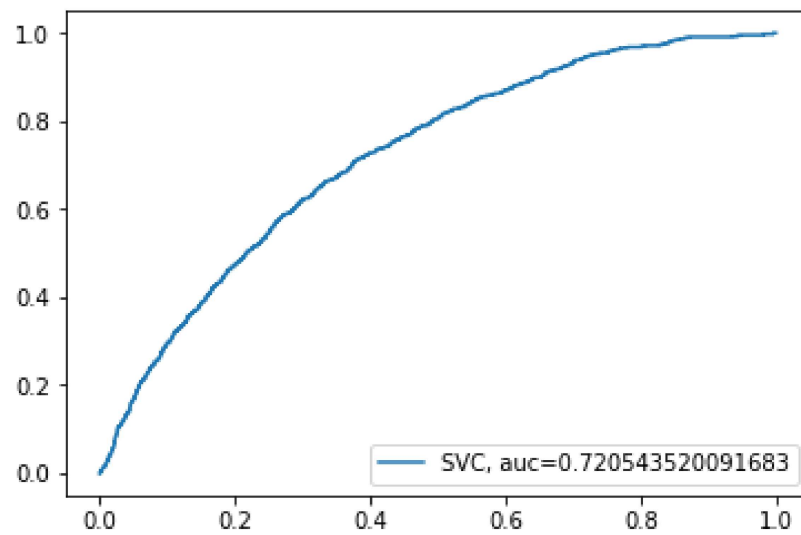


Figure 5.6: AUC for SVC.

Table 5.2: Results with SMOTE Data Sampling.

Model	Accuracy	AUROC
Logistic Regression	67%	0.72
Linear SVC	67%	0.73
Random Forrest	85%	0.86
ANN	84%	0.83

Table 5.3: Complications vs Expiry.

Patients with Complications	Patients Died at Hospital
2754	576

Chapter 6

Conclusion and Future Investigation

6.1 Conclusion

Extended diagnoses and patients stay at the hospital is associated with not only the health of patient, cost, increased number of deaths but also increased number of readmissions. Each of these parameters defines the hospital performance. So, our focus was to produce insights that can complement these parameters ranging from cost to patient health. In order to propose a new method or optimize models that are already built. We can improve the data quality enhance methods of feature engineering and tune our models to outperform earlier methods. To do that we applied following steps:

- For high quality: feature engineering with ETL, batch processing and data sampling
- Apply different sampling techniques and implement different models

These are presented in the order in above sections. The data is important for each predicting task, if the data is not available, there is no point of using complex model and optimize with resources and conversely if the data is available but the models are not complex enough to utilize that data can create performance drawbacks [40]. First, the use of MIMICIII database to study electronic health records is probably the best to develop studies and researches that can contribute to the society. As it is being shared with the researchers, educational community and scientists, people are aggressively contributing to the cause and creating a huge impact. But, one must consider the existence of unstructured data available in MIMIC-III which lead to extensive search for variables without access to actual system deployed at hospitals. Because of that a considerable part was dedicated to the ETL to extract meaningful features which can lead to reliable prediction models.

We used Python as it provide large number of libraries and models. Finally, in our study, Random Forest Classifier which gives the best prediction after both the sampling of data ADASYN and SMOTE. The accuracy also indicated that by including the derived concept combined with chart and lab events against each patient. By extraction of derived features of organ failure expanded our problem and features space which lead to the analysis of the importance of the variables that have relationship between each other and to complications, We have achieved models with great prognostic capacity using demographic, concepts, lab events, chart events features and interpolating the minority class with different techniques which are not only intuitive for management's view, for patient's health and for doctors as well.

6.2 Future Investigation

The future studies regarding the complications and mortality of patients due to those complications are evident. As we can do multiclass classification which drills down the complications related to specific types instead of binary class of having complications or not. It was more adequate to start off with linear models and then further moved towards complex models, so, to improve our models we can adapt complex models in our future work. Another line of research would be to engineer more features, create new concepts, and combine NLP techniques for textual features and applying complex models to contribute more the already done research.

References

- [1] M. C. Data, *Secondary Analysis of Electronic Health Records*. Springer International Publishing, 2016.
- [2] S. Bergamaschi, “A semantic approach to etl technologies,” in *Data Knowledge Engineering*, 2011, pp. 717–731.
- [3] P. Vassiliadis, “Towards a benchmark for etl workflows,” 2007.
- [4] A. E. Johnson, “Mimic-iii, a freely accessible critical care database,” in *Scientific Data 3 (CD)*, 2016. [Online]. Available: <https://mimic.physionet.org/>
- [5] R. O. Deliberato, “An evaluation of the influence of body mass index on severity scoring,” in *Critical care medicine*, 2019, pp. 247–253.
- [6] S. Bose, “Impact of intensive care unit discharge delays on patient outcomes: a retrospective cohort study,” *Journal of intensive care medicine*, vol. 34, no. 11–12, pp. 924–929, 2019.
- [7] R. O. Deliberato, “Severity of illness scores may misclassify critically ill obese patients,” in *Critical care medicine*, vol. 46, no. 3, 2018, pp. 394–400.
- [8] International classification of diseases. [Online]. Available: <https://www.cdc.gov/nchs/icd/icd9.htm>
- [9] Health insurance portability and accountability act (hipaa). [Online]. Available: <https://www.hhs.gov/hipaa/for-professionals/index.html>
- [10] A. G. Rapsang and D. C. Shyam., “Scoring systems in the intensive care unit: a compendium.” *Indian journal of critical care medicine: peer-reviewed, official publication of Indian Society of Critical Care Medicine*, vol. 18, no. 4, p. 220, 2014.
- [11] V. H. Kury, Fabrício SP and J. J. Cimino., “Reproducing a prospective clinical study as a computational retrospective study in mimic-ii.” in *AMIA Annual Symposium Proceedings*, 2015, p. 804.
- [12] D. W. Bates, “Big data in health care: using analytics to identify and manage high-risk and high-cost patients,” in *Health Affairs*, 2014, pp. 1123–1131.
- [13] J. B. Segal and N. R. Powe., “Accuracy of identification of patients with immune thrombocytopenic purpura through administrative records: a data validation study,” *American journal of hematology*, pp. 12–17, 2004.
- [14] A. F. Eichler and E. B. Lamont., “Utility of administrative claims data for the study of brain metastases: a validation study,” *Journal of neuro-oncology*, pp. 427–431, 2009.

- [15] A. Perotte, “Diagnosis code assignment: models and evaluation metrics,” pp. 231–237, 2013.
- [16] M. T. Mullen, “Icd9 codes cannot reliably identify hemorrhagic transformation of ischemic stroke,” p. 505, 2013.
- [17] L. V. Lita, “Large scale diagnostic code classification for medical patient records,” in *Proceedings of the Third International Joint Conference on Natural Language Processing*, 2008.
- [18] T. Baumel, “Multi-label classification of patient notes: case study on icd code assignment,” in *Workshops at the Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [19] S. M. Bache, Richard and A. Taweel., “An adaptable architecture for patient cohort identification from diverse data sources.” *Journal of the American Medical Informatics Association*, pp. e327–e333, 2013.
- [20] J. C. Ferrao, “Using structured ehr data and svm to support icd-9-cm coding,” in *IEEE International Conference on Healthcare Informatics*, 2013.
- [21] V. Mandalapu, “Understanding the relationship between healthcare processes and in-hospital weekend mortality using mimic iii,” in *Smart Health*, 2019.
- [22] T. B. Huang, Yuan-Lan and Z.-D. Hu., “Using freely accessible databases for laboratory medicine research: experience with mimic database,” *Journal of Laboratory and Precision Medicine*, 2017.
- [23] R. G. Shao, Jianyin and J. C. Facelli., “Semantic characterization of clinical trial descriptions from clinicaltrials. gov and patient notes from mimic-iii,” *Journal of Clinical and Translational Science*, pp. 12–12, 2017.
- [24] T. Gentimis, “Recognizing human action in time-sequential images using hidden Markov model,” in *IEEE 15th Intl Conf on Dependable, Autonomic and Secure Computing, 15th Intl Conf on Pervasive Intelligence and Computing, 3rd Intl Conf on Big Data Intelligence and Computing and Cyber Science and Technology Congress (DASC/PiCom/DataCom/CyberSciTech)*, 2017, pp. 1194–1201.
- [25] C. O. Huang, Jinmiao and L. W. Sy., “An empirical evaluation of deep learning for icd-9 code assignment using mimic-iii clinical notes,” in *Computer Methods and Programs in Biomedicine*, 2019, pp. 141–153.
- [26] C. M. Yao, Liang and Y. Luo., “Clinical text classification with rule-based features and knowledge-guided convolutional neural networks.” in *BMC medical informatics and decision making*, 2019, p. 71.
- [27] I. Solt, “Semantic classification of diseases in discharge summaries using a context-aware rule-based classifier,” *Journal of the American Medical Informatics Association*, pp. 580–584, 2009.

- [28] T. J. P. Johnson, Alistair EW and R. G. Mark., “Reproducibility in critical care: a mortality prediction case study,” in *Machine Learning for Healthcare Conference*, 2017, pp. 361–376.
- [29] V. Arvind, “Predicting surgical complications in adult patients undergoing anterior cervical discectomy and fusion using machine learning,” in *Neurospine*, 2018, p. 329.
- [30] J. C. Rojas, “Predicting intensive care unit readmission with machine learning using electronic health record data,” in *Annals of the American Thoracic Society*, 2018, pp. 846–853.
- [31] A. Meyer, “Machine learning for real-time prediction of complications in critical care: a retrospective study,” in *The Lancet Respiratory Medicine*, 2018, pp. 905–914.
- [32] R. Davoodi and M. H. Moradi., “Mortality prediction in intensive care units (icus) using a deep rule-based fuzzy classifier,” *Journal of biomedical informatics*, pp. 48–59, 2018.
- [33] N. V. Chawla, “Smote: Synthetic minority over-sampling technique,” *Journal of artificial intelligence research*, pp. 321–357, 2002.
- [34] H. He, “Adasyn: Adaptive synthetic sampling approach for imbalanced learning,” in *IEEE International Joint Conference on Neural Networks*, 2008.
- [35] C. Nwankpa, “Activation functions: Comparison of trends in practice and research for deep learning,” in *arXiv preprint arXiv*, 2018.
- [36] T. K. Ho, “Random decision forests.” in *Proceedings of 3rd international conference on document analysis and recognition.*, 1995.
- [37] L. Breiman, “Random forests,” in *Machine learning*, 2001.
- [38] Y. Amit and D. Geman., “Shape quantization and recognition with randomized trees,” in *Neural computation*, 1997.
- [39] K. Eckle and J. Schmidt-Hieber., “A comparison of deep networks with relu activation function and linear spline-type methods.” in *Neural Networks*, 2019.
- [40] P. N. Halevy, Alon and F. Pereira., “The unreasonable effectiveness of data,” 2009.

Appendix A

Resources

Available code repository

- <https://github.com/faisalmaqbool94/Thesis-Bioinformatics-MIMICIII->

Repository of code shared development community provided by MIT Laboratory for Computational Physiology

- <https://github.com/MIT-LCP/mimic-code>