



# Disease Outcome Prediction

<sup>1</sup>Sakethram Marpu, <sup>2</sup>Rafeeq Shaik, <sup>3</sup>Aileen Peekka

<sup>1</sup>Student, <sup>2</sup>Student, <sup>3</sup>Student

<sup>1</sup>Computer Science Engineering,

<sup>1</sup>Vellore Institute of Technology, Vellore, India

**Abstract :** This study is the next frontier of data analytics and healthcare that will revolutionize predictive disease management. We describe user-friendly type interfaces using Decision Trees Random Forests, Naive Bayes and k-nearest Neighbour algorithms for prognostic disease prediction together with the clinical application abstract. The following application is implemented using Tkinter library in Python for patients to input their symptoms and used four different machine learning models that were trained on the desired datasets. The training and testing dataset consists of a series of symptoms matched to distinct diseases. This facilitates a clearer and more accurate understanding in the results by allowing users to explore disease-symptom associations. The application uses a database to record patient data and forecast predictions. It allows data tracking which helps ease in further study. In conclusion, this abstract demonstrates that models can be incorporated in an intuitive user interface and shows how these applications could evolve to improve predictive disease analytics towards improved healthcare outcomes.

**IndexTerms - Machine Learning, Disease outcome prediction, Healthcare, Clinical application, Forecasting.**

## INTRODUCTION

And year after year, the volume of patients and diseases lifts up the capacity in the medical system leading to a price increase over time on health cost levels throughout many countries. Most of the diseases require one to obtain medical advice so they can be in a position to treat. Predicting disease using an algorithm can be remarkably simple and cheap provided it has sufficient data. This is where prevention as treatment comes with a major symptom-based illness prediction component. Here, we have tried to predict the disease as accurately as possible in this work, based on symptoms of a patient. This type of system could represent a major improvement for future medical treatments. Nowadays, one of the significant effects that a disease analysis has is modern health care as it influences essential parts like prevention and early detection for better treatment. Data Analytics — Healthcare providers can explore their patients' information for trends and patterns to diagnose diseases early. Prognostic Disease Analytics that identifies early signals of disease worsening, predicts how patients will respond to treatment and anticipates adverse events all enable clinicians to intervene pro-actively care optimizing patient outcomes as well as allocate resources more effectively. Prognostic disease analytics are essential in the era of precision medicine and value-based care. The promise that healthcare technologies can be used to improve patient outcomes and disrupt the way care is delivered, as well save money in health system use has driven tremendous research innovation and investment into it. With increasing involvement in this nuanced area, it is important to examine the applications of AI/ML in clinical research; describe the existing challenges and considerations associated with leveraging these technologies for healthcare purposes; discuss overviews on ethical use cases specific examples relevant to various stakeholders. The K-Nearest Neighbours is a simple but powerful algorithm for regression and classification applications. It is based on the closeness principle where whatever be the class of a valid data point, that can be determined by considering its k nearest neighbors level i.e. majority vote or even average value. This is because KNN has a user-friendly background, but more importantly its behavior is intuitive and this approach works best for local pattern or cluster based data. The RF learning technique uses a great number of decision trees in order to establish an acceptable predictive model that is also accurate. Each tree in the forest is trained to predict on its own 2 using a bootstrap sample of original data. Finally, a final forecast is made by aggregating the forecasts from all trees, often in an ensemble process like majority vote. RF has been doing well with high-dimensional data, finding nonlinear associations and over-fitting reduction. Decision Trees: DT is an interpretable framework which efficiently partitions the feature space based on features which best segregate cases into homogenous groups w.r.t to target variable. Only the inner nodes represent a decision based on some feature and each leaf combines to class label or regression value depending upon input features. A decision tree is intuitive, interpretable where we make decisions at every node and are used for both classification and regression applications. Naive Bayes is a probabilistic classifier with the basis on Bayesian theorem and it assumes that features are independent given their class label. Naive Bayes often works well in practice, even though it's a simplistic algorithm that can offer surprisingly good results especially for text classification and other high-dimensional datasets.

## NEED OF THE STUDY.

Moreover, the growing global population and upward trend of chronic diseases has placed incredible pressure on health systems around the globe. Disease diagnosis and management based on classic methods are long-term, laborious processes that can also be subject to bias due to human error. Considering this, there is an imperative need for novel solutions that may help improving the efficiency and precision of disease prediction and management. The study attempts to meet this need by combining sophisticated machine learning models with healthcare data into an intuitive application that can predict diseases based on patient symptoms. It should be noted that relevance of the study is based on

1. **Early Detection of Disease:** Detecting disease accurately and promptly is important for better patient results with the use of less healthcare resources. The project is intended to predict diseases in a systematic manner based on algorithms such as Decision Trees, Random Forests, Naive Bayes and k-Nearest Neighbors which will help conduct medical interventions early.
2. **Big Data Healthcare:** As the world becomes more data driven, interpreting large amounts of patient related info will be a make-or-break for any healthcare entity that aims to benefit from this technology. This project illustrates the application of data analytics in healthcare that is aiding doctors towards personalized medicine which would empower more evidence based.
3. **Resource Optimisation:** The application predicted in this study will aid in lessening the healthcare and medical infrastructure by easing down on diagnostic procedures. Faster and more accurate diagnoses of disease pave the way for improved resource allocation, lessened costs and upgraded healthcare delivery.
4. **Advancements in technology:** The union between machine learning models and the user-friendly interface founds our capacity of imagining another form of healthcare. The study demonstrates that advancements in technology have the potential to generate multi-dimensional insights pertaining current or future healthcare delivery.

This expands the importance of this study's results in today's data-driven healthcare and prolonged catch up-care issues many hospital systems are facing over natural disasters or pandemics. Sharing this project in IJNRD publication to enhance the discussion around innovative health solutions, which could be beneficial for researchers and practitioners as well policy makers.

## RESEARCH METHODOLOGY

The methodology section outlines the plan and method that how the study is conducted. This includes the Universe of the study, sample of the study, Data and Sources of Data, study's variables and analytical framework. The details are as follows;

### 3.1 Population and Sample

The population could be either all patients in a healthcare system or individuals at risk of specific diseases, regardless of the place where they live and their age.

The sample in this project is maybe 10,000 patient records (patient history) along with symptoms and diagnosed counts. The following sample for training the algorithms of machine learning (Decision Trees, Random Forests, Naive Bayes, k-Nearest Neighbors) in order to do disease prediction by symptoms.

### 3.2 Data and Sources of Data

**Data:** The data in the "Prognostic Disease Analytics" project describes details on how to prepare and train a model (machine learning) for predicting disease. These data generally include patient records, symptoms and diagnosis etc. common to medical practice

#### Sources of Data

Datasource for this project can be extracted from:

1. **EHRs (Electronic Health Records):** Patient digests stored in EHR systems at healthcare organizations, including fine-grained and highly detailed information on symptoms bias illnesses treatments outcomes.
2. **Medical Databases:** Archives of medical information (for example, patient health records) hosted by a healthcare institution or a research collaboration.
3. **Public Datasets:** Publicly available datasets that are designed or created for machine learning and data science tasks The popular source is the Kaggle, one of known platforms for datasets used in variety machine learning projects.

### 3.3 Theoretical framework

In this paper, we introduce the theoretical foundation of using machine learning to model diseases based on symptoms presented by patients. The framework unites methods of healthcare, data science and machine learning to illustrate the functionality through which these areas interplay yet integrate accurate predictions of disease.

Components of the Theoretical Framework

#### 1. Health and Disease Forecasting

Idea: Healthcare should strive to first and foremost effectively diagnose, treat diseases. It would be greatly beneficial to health care if we could predict diseases earlier and more accurately.

Application in Project: Predicting Disease from Symptoms Reported by Patients — This project is targeting early diagnosis and intervention within healthcare, It maps well with this hook.

#### 2. Machine Learning Models

Concept: Machine learning algorithms learn from data and then use that knowledge to make predictions or decisions. Patterns in data are analyzed using modeling methods like Decision Trees, Random Forests, Naive Bayes and k-Nearest Neighbors.

Application in Project :These models are used to predict diseases by studying patterns in the symptoms reported. Both models have their advantages, and using them together helps to make predictions more robust yet accurate.

#### 3. Data Analytics

- Noun: Data analytics is the process of analyzing data to get useful insights. Used to retrieve valuable knowledge from existing data, it serves as an aid in terms of making sound decisions.

- Application in Project: In this work developers utilize the patient data to do their predictions by using Data Analytics. As such, machine learning is performed on a very large scale using big data.

#### 4. Intro to Human-Computer Interaction (HCI)

Concept: UI is the place HCI to build user friendly systems which levels human computer interaction. It makes technology available and easy for the people to use.

Project Application : Here, we have a user-interface of the project using Tkinter library that is developed in python which helps Health care professionals or Patients to View input Symptoms and get Disease Prediction easily.

#### 5. Database Management

Application in Project: Implementing a database to store patient data and predictions, we are able to conduct further analysis while making patient information traceable for the future.

### 3.4 Statistical tools and econometric models

Statistical tools and econometric models matter greatly in the realm of (Big) Data analysis, predicting outcomes, etc... For example with this project "Prognostic Disease Analytics" The key tools and models are summarized:

#### 1. Descriptive Statistics

- Idea: Descriptive statistics are there to summarize and describe the features of a dataset and tell about central tendency, dispersion & distribution of data.

Application in the Project: Descriptive statistics are applied to study how symptoms and diseases are spread through data, that is what interests such kind of analysis before training (case) machine learning models war.

Example: Mean, median, mode, standard deviation & frequency distributions

#### 2. Inferential Statistics

Idea: Inferential Statistics to do predictions about a population based on only one sample of data.

Implementation in Project: This tells us how well the predictions were considering whatever data being trained on, and predicts based on the same set of human population.



It includes stuff like hypothesis testing, confidence intervals and regression analysis.

### 3. Correlation Analysis

Concept: Correlation analysis quantifies the strength and direction of a relationship between two variables.

Application in Projects: It is used to determine how one disease can be diagnosed on the basis of some symptom and we build a predictive model around it.

Examples : Pearson correlation coefficient, Spearman's rank-correlation respectful

### 4. Regression Analysis

Concept: Simply put, regression analysis assesses the relationship between a dependent variable and one or more predictors (independent variables), using data samples.

Project Application: using algorithms to predict the probability of getting a disease given certain symptoms as input.

Types: KNN, Random Forest, Decision Tree and Naive Bayes

#### 1. Logistic Regression

Logistic regression: Used to predict the probability of a binary outcome based on one or more predictors.

Application in Project It can be applied to find out whether a patient has disease or not (Binary outcome) based on the symptoms inflicted.

Logistic regression is widely used in econometrics as a classification problem because of the nature of problems more appropriate related to medical diagnosis and disease prediction.

#### 2. Probit and Logit Models

Concept: Regressions models to estimate binary or categorical outcomes can be a logit regression model and probit regression model

Application in Project:- These models can be applied to disease prediction where the resulting outcome is categorical, for example predicting a specific category of diseases on symptoms..

#### 3. Time Series Analysis

Category: Time series analysis revolves around the set of statistical techniques are used to model and forecast data points which were collected in a specific time interval.— Concept:

Application to Project: Time series analysis is a more elusive type of analysis, probably not common in all fields and it might usually be left out by some people completely. For instance using time series you can model trends within disease outbreaks or how symptoms evolve over time (The basic example we did this project on).

#### 4. Panel Data Analysis

Concept : Panel data analysis is when we have observations on a variable for many entities over time which allows us to study dynamics of the variables across-time and cross-ordinal (ie. etc).

Usage in Project: The model is appropriate for longitudinal data where patient symptoms or disease outcomes change over time.

#### 3.4.1 Descriptive Statistics

Descriptive statistics get used in statistical analysis to describe the main features of a dataset quantitatively. Its purpose is to give terse information about the sample and measures in order that these will be understood at a glance: indications, indicators of central tendency values or distribution, variability within the data set.

##### 1. Measures of Central Tendency

Mean: Mean is the average of all data points in a dataset. The sum of all the values and then divided by no.

For our example in the project "Prognostic Disease Analytics," this might represent an average age of patients within a dataset.

**Median:** The middle number in a data set when it has been ranked from least to greatest. It is more robust to outliers compared the mean

**Likely example:** Measures such as median, the central tendency of how many symptoms patients have outputs for this use case you would probably `contra_identifierGetMethod(median)` with a `descriptionpseudo_execution(prefix)("how manySymptoms")`.

**Mode:** The most common value in a dataset.

The mode could be the one which identifies the most frequently reported symptom by patients, e.g.

## 2. Dispersion( Variability Measures)

**Range-** the maximum and minimum in a data set. It is a measure of dispersion or spreadoutness in the data.

**Variance:-** Amount of variation between the data points and mean. It is the sum of squares deviation from the mean.

**Standard Deviation:** The square root of variance. It shows how much a data point varies from the average.

For example, low values of standard deviation in frequency represented most patients had nearly similar symptom counts.

## 3. Measures of Shape

**Skewness** — Describes the lack of symmetry in data. A distribution is positively skewed if it has a long right tail, and negatively skewed if it has a long left-tail.

For example Skewness might tell us whether a particular symptom may be reported by the majority of patients rarely if at all, or frequently.

**Kurtosis:** the tailedness of the distribution. High kurtosis represents heavy tails and Low kurtosis represents light-tailed.

For instance: If I want to understand how symptom severity is distributed across the dataset, kurtosis might be a good indicator of this.

## 4. Frequency Distribution

**Point:** A frequency distribution provides the number of times each value appears in a dataset. Tables or charts for that can look at histogram or bar It definitely helps when you. This is how we are able to see data through symbols.

## 5. Percentiles and Quartiles

**Percentiles:** Percentile shows the standing of a value within his dataset. The 50th percentile is the same as a median.

For example: One might use percentiles for the distribution of patient ages where, 25th percentile would indicate that age below which lies in lower quartile (i.e. only 25 % patients lesser than this), ....

**Quartiles:** Quartiles are positions in the ordered datasets that divide it into four equal parts with Q1 (first quartile, 25th percentile), beangQ2(median) and Q3(Third-quartile,75th percintle).

### 3.4.2 Literature Survey:

#### **DISEASE PREDICTION SYSTEM USING SUPPORT VECTOR MACHINE AND MULTI-LINEAR REGRESSION :-**

It is a technique of machine learning intended to predict the disease with symptoms[1]. However, the machine learning one of dozens available could serve to save a billion lives. Now mit ML-verfahren können wir sehr viele krankheits-symptome bereits im voraus vorhersagen. However, the model does not employ many techniques to increase the accuracy of illness prediction.

**CROP YIELD PREDICTION BY DATA MINING:-** Here, the investigation was that single disease prediction by more monitored methods of ML with a substantial searching for papers [2]. The search was based on different key terms (Table 1) and performed in PubMed as well the Scopus databases. This study aims to identify the main trends among various categories of guided machine learning algorithms, across their performance and utilization in predicting disease risk.

#### **A HYBRID ARTIFICIAL INTELLIGENCE BASED SYSTEM FRAMEWORK FOR HEART DISEASE PREDICTION:**

Often called the killer disease [4], is one of the most complex, intricate and deadliest diseases that affect human beings today. In this disease, the heart is usually not able to effectively pump blood throughout other parts of body in order for normal function of

body. Consequently, heart failure Ok arouses (the attenuate carts the report of physical examination and medical history indicates that 4 suspects) on a basis for clinicians to analyze patients with penetrative methods for experimental cardiac disorder. These all are the procedures which fail to provide accurate diagnosis and delay in the output due ti human mistakes. Further, they take more time for exams and are costly as well computationally heavy on the system.

**MACHINE LEARNING BASED DISEASE PREDICTION:** — An excel sheet was generated form an open-source dataset, and we listed each symptom with its related disease. After that, the information of age and gender set up according to diseases. We included nearly 230 unique disorders comprising more than1000 individual symptoms in our list for the machine learning training process. input would be the age, sex and symptoms of an individual, output will predict whether he is sick or not [4]. A modified version of this KNN. We used an integer parameter, K to find the nearest predicted values by, for example in our case if k=17 so we want mot important 17 predection value. However, the method gets much more affected if K's value is set way too small and it significantly increased support to anomaly points in comparison with nucleus.

**MAJOR COMPONENTS AND ARCHITECTURE OF A DATA MINING SYSTEM:** Source of data, The Server that is responsible for accessing a database (or some other knowledge servers) via the Repository server and it gets this information from the repository filtered by user request KnowledgeBase acts to search guide considering already set limitations Mining Engines consist series basic modules such as characterisation (post-processing), classification, clustering association rules or simply Apriori algorithm, regression evolution analysis [5] Pattern evaluation modulecommunicate with the DM-modules to discovered interesting pattern A data mining method usually follows a particular model/framework for designing their function approach. The Apriori algorithm, simple and easy to understand it may be, is limited. The main disadvantage is that it costs too much and takes more time on keeping many candidate sets with frequent item sets, low minimum support or large amount of items.

## IV. RESULTS AND DISCUSSION

### 4.1 Results of Descriptive Statics of Study Variables

Figure 4.1 - GUI

Figure 5.2- Result

```

10
11 35
12 34
13 34
14 34
15 34
16 34
17 37
18 38
19 38
20 38
21 38
22 40
23 40
24 40
25 40
26 40
27 40
28 40
29 40
30 40
31 40
32 40
33 40
34 40
35 40
36 40
37 40
38 40
39 40
40 40
41 40
42 40
43 40
44 40
45 40
46 40
47 40
48 40
49 40
50 40
51 40
52 40
53 40
54 40
55 40
56 40
57 40
58 40
59 40
60 40
61 40
62 40
63 40
64 40
65 40
66 40
67 40
68 40
69 40
70 40
71 40
72 40
73 40
74 40
75 40
76 40
77 40
78 40
79 40
80 40
81 40
82 40
83 40
84 40
85 40
86 40
87 40
88 40
89 40
90 40
91 40
92 40
93 40
94 40
95 40
96 40
97 40
98 40
99 40
100 40
101 40
102 40
103 40
104 40
105 40
106 40
107 40
108 40
109 40
110 40
111 40
112 40
113 40
114 40
115 40
116 40
117 40
118 40
119 40
120 40
121 40
122 40
123 40
124 40
125 40
126 40
127 40
128 40
129 40
130 40
131 40
132 40
133 40
134 40
135 40
136 40
137 40
138 40
139 40
140 40
141 40
142 40
143 40
144 40
145 40
146 40
147 40
148 40
149 40
150 40
151 40
152 40
153 40
154 40
155 40
156 40
157 40
158 40
159 40
160 40
161 40
162 40
163 40
164 40
165 40
166 40
167 40
168 40
169 40
170 40
171 40
172 40
173 40
174 40
175 40
176 40
177 40
178 40
179 40
180 40
181 40
182 40
183 40
184 40
185 40
186 40
187 40
188 40
189 40
190 40
191 40
192 40
193 40
194 40
195 40
196 40
197 40
198 40
199 40
200 40
201 40
202 40
203 40
204 40
205 40
206 40
207 40
208 40
209 40
210 40
211 40
212 40
213 40
214 40
215 40
216 40
217 40
218 40
219 40
220 40
221 40
222 40
223 40
224 40
225 40
226 40
227 40
228 40
229 40
230 40
231 40
232 40
233 40
234 40
235 40
236 40
237 40
238 40
239 40
240 40
241 40
242 40
243 40
244 40
245 40
246 40
247 40
248 40
249 40
250 40
251 40
252 40
253 40
254 40
255 40
256 40
257 40
258 40
259 40
260 40
261 40
262 40
263 40
264 40
265 40
266 40
267 40
268 40
269 40
270 40
271 40
272 40
273 40
274 40
275 40
276 40
277 40
278 40
279 40
280 40
281 40
282 40
283 40
284 40
285 40
286 40
287 40
288 40
289 40
290 40
291 40
292 40
293 40
294 40
295 40
296 40
297 40
298 40
299 40
300 40
301 40
302 40
303 40
304 40
305 40
306 40
307 40
308 40
309 40
310 40
311 40
312 40
313 40
314 40
315 40
316 40
317 40
318 40
319 40
320 40
321 40
322 40
323 40
324 40
325 40
326 40
327 40
328 40
329 40
330 40
331 40
332 40
333 40
334 40
335 40
336 40
337 40
338 40
339 40
340 40
341 40
342 40
343 40
344 40
345 40
346 40
347 40
348 40
349 40
350 40
351 40
352 40
353 40
354 40
355 40
356 40
357 40
358 40
359 40
360 40
361 40
362 40
363 40
364 40
365 40
366 40
367 40
368 40
369 40
370 40
371 40
372 40
373 40
374 40
375 40
376 40
377 40
378 40
379 40
380 40
381 40
382 40
383 40
384 40
385 40
386 40
387 40
388 40
389 40
390 40
391 40
392 40
393 40
394 40
395 40
396 40
397 40
398 40
399 40
400 40
401 40
402 40
403 40
404 40
405 40
406 40
407 40
408 40
409 40
410 40
411 40
412 40
413 40
414 40
415 40
416 40
417 40
418 40
419 40
420 40
421 40
422 40
423 40
424 40
425 40
426 40
427 40
428 40
429 40
430 40
431 40
432 40
433 40
434 40
435 40
436 40
437 40
438 40
439 40
440 40
441 40
442 40
443 40
444 40
445 40
446 40
447 40
448 40
449 40
450 40
451 40
452 40
453 40
454 40
455 40
456 40
457 40
458 40
459 40
460 40
461 40
462 40
463 40
464 40
465 40
466 40
467 40
468 40
469 40
470 40
471 40
472 40
473 40
474 40
475 40
476 40
477 40
478 40
479 40
480 40
481 40
482 40
483 40
484 40
485 40
486 40
487 40
488 40
489 40
490 40
491 40
492 40
493 40
494 40
495 40
496 40
497 40
498 40
499 40
500 40
501 40
502 40
503 40
504 40
505 40
506 40
507 40
508 40
509 40
510 40
511 40
512 40
513 40
514 40
515 40
516 40
517 40
518 40
519 40
520 40
521 40
522 40
523 40
524 40
525 40
526 40
527 40
528 40
529 40
530 40
531 40
532 40
533 40
534 40
535 40
536 40
537 40
538 40
539 40
540 40
541 40
542 40
543 40
544 40
545 40
546 40
547 40
548 40
549 40
550 40
551 40
552 40
553 40
554 40
555 40
556 40
557 40
558 40
559 40
560 40
561 40
562 40
563 40
564 40
565 40
566 40
567 40
568 40
569 40
570 40
571 40
572 40
573 40
574 40
575 40
576 40
577 40
578 40
579 40
580 40
581 40
582 40
583 40
584 40
585 40
586 40
587 40
588 40
589 40
590 40
591 40
592 40
593 40
594 40
595 40
596 40
597 40
598 40
599 40
600 40
601 40
602 40
603 40
604 40
605 40
606 40
607 40
608 40
609 40
610 40
611 40
612 40
613 40
614 40
615 40
616 40
617 40
618 40
619 40
620 40
621 40
622 40
623 40
624 40
625 40
626 40
627 40
628 40
629 40
630 40
631 40
632 40
633 40
634 40
635 40
636 40
637 40
638 40
639 40
640 40
641 40
642 40
643 40
644 40
645 40
646 40
647 40
648 40
649 40
650 40
651 40
652 40
653 40
654 40
655 40
656 40
657 40
658 40
659 40
660 40
661 40
662 40
663 40
664 40
665 40
666 40
667 40
668 40
669 40
670 40
671 40
672 40
673 40
674 40
675 40
676 40
677 40
678 40
679 40
680 40
681 40
682 40
683 40
684 40
685 40
686 40
687 40
688 40
689 40
690 40
691 40
692 40
693 40
694 40
695 40
696 40
697 40
698 40
699 40
700 40
701 40
702 40
703 40
704 40
705 40
706 40
707 40
708 40
709 40
710 40
711 40
712 40
713 40
714 40
715 40
716 40
717 40
718 40
719 40
720 40
721 40
722 40
723 40
724 40
725 40
726 40
727 40
728 40
729 40
730 40
731 40
732 40
733 40
734 40
735 40
736 40
737 40
738 40
739 40
740 40
741 40
742 40
743 40
744 40
745 40
746 40
747 40
748 40
749 40
750 40
751 40
752 40
753 40
754 40
755 40
756 40
757 40
758 40
759 40
760 40
761 40
762 40
763 40
764 40
765 40
766 40
767 40
768 40
769 40
770 40
771 40
772 40
773 40
774 40
775 40
776 40
777 40
778 40
779 40
780 40
781 40
782 40
783 40
784 40
785 40
786 40
787 40
788 40
789 40
790 40
791 40
792 40
793 40
794 40
795 40
796 40
797 40
798 40
799 40
800 40
801 40
802 40
803 40
804 40
805 40
806 40
807 40
808 40
809 40
810 40
811 40
812 40
813 40
814 40
815 40
816 40
817 40
818 40
819 40
820 40
821 40
822 40
823 40
824 40
825 40
826 40
827 40
828 40
829 40
830 40
831 40
832 40
833 40
834 40
835 40
836 40
837 40
838 40
839 40
840 40
841 40
842 40
843 40
844 40
845 40
846 40
847 40
848 40
849 40
850 40
851 40
852 40
853 40
854 40
855 40
856 40
857 40
858 40
859 40
860 40
861 40
862 40
863 40
864 40
865 40
866 40
867 40
868 40
869 40
870 40
871 40
872 40
873 40
874 40
875 40
876 40
877 40
878 40
879 40
880 40
881 40
882 40
883 40
884 40
885 40
886 40
887 40
888 40
889 40
890 40
891 40
892 40
893 40
894 40
895 40
896 40
897 40
898 40
899 40
900 40
901 40
902 40
903 40
904 40
905 40
906 40
907 40
908 40
909 40
910 40
911 40
912 40
913 40
914 40
915 40
916 40
917 40
918 40
919 40
920 40
921 40
922 40
923 40
924 40
925 40
926 40
927 40
928 40
929 40
930 40
931 40
932 40
933 40
934 40
935 40
936 40
937 40
938 40
939 40
940 40
941 40
942 40
943 40
944 40
945 40
946 40
947 40
948 40
949 40
950 40
951 40
952 40
953 40
954 40
955 40
956 40
957 40
958 40
959 40
960 40
961 40
962 40
963 40
964 40
965 40
966 40
967 40
968 40
969 40
970 40
971 40
972 40
973 40
974 40
975 40
976 40
977 40
978 40
979 40
980 40
981 40
982 40
983 40
984 40
985 40
986 40
987 40
988 40
989 40
990 40
991 40
992 40
993 40
994 40
995 40
996 40
997 40
998 40
999 40
1000 40

```

The deployed morbidity forecasting system produced competitive results across several machine learning algorithms. This dataset was used to build 4 models — a Decision Tree model, Random Forest model, Naive Bayes Model and k-Nearest Neighbours. The scatter plots also facilitated a visual understanding of the distribution of symptoms for predicted diseases and help to interpret our results. This allows for the real-time predictions from various models on a user-friendly graphical interface(GUI) and great interaction with users. The reason to save user data as well as the predictions into an SQLite database was storing them for reviewing or future improvements of the system. The system acts as a great screening tool for initial disease prediction, and can provide even better results if trained on more data(i.e supplements) and minorly adjust the Machine learning models to increase automatically(HLC Grading up).

## V. CONCLUSION AND FUTURE SCOPE

In the medical field our application is for the first time and we are introducing this innovative application Prognostic Disease Analytics. It transforms predictive and preventive healthcare with the combination of healthcare expertise blended well with rigorous data analytics techniques. Its main use as a user-friendly system that streamlines the input of patient symptoms and additional learning algorithms — Decision Trees, Random Forests, Naive Bayes and k-Nearest Neighbors. The algorithmic diversity protects the robustness of predictions, being applicable to datasets and disease profiles. In addition, users should explore the data by using basic capabilities that will help more in understanding which symptom is most closely related to a particular disease. Furthermore the use of a database allows traceability and post-processing, making the system useful in time. This application is an important asset to the field of prognostic disease analytics, by focusing on usability, algorithm diversity and data management. It has a ton of room to grow. misfortunes of healthcare expulsion through accurate predictions for disease and help subtleties-informed decisions by medical experts In summary, this abstract highlights the encouraging trend of model building for data mining to intuitive interfaces supporting Health Care Needs.

### I. ACKNOWLEDGMENT

We would like to thank everyone who has supported me in the course of this study. I am grateful to my advisor Mr. Kannadasan R, our Head of Department Mrs. Umadevi Ks and our Dean of academics Mr. Ramesh Babu for their advice and knowledge in the disease outcome prediction field. I would also like to acknowledge colleagues and friends, who had taken the time to read my drafts and provide a lot of helpful critiques that allowed me to hone these thoughts.

We are also thankful to the Vellore Institute of Technology for providing its departmental resources and facilities. Finally, a huge thank you to my family and friends for never-ending support & understanding along the way.

### REFERENCES

- [1] - "Disease prediction system using support vector machine and multi-linear regression" - Vijayarani, S. and Dhayanand
- [2] - "A comparative study on disease prediction using supervised machine learning algorithms" - Rahman, A.S., Shamrat, F.J.M., Tasnim, Z., Roy, J. and Hossain
- [3] - "A hybrid intelligent system framework for heart disease prediction" - Shraddha Subhash Shirsath, Prof. Shubhangi Patil
- [4] - "Disease prediction by machine learning" - Thomas H. Cormen, Charles E. Leiserson, and Ronald L. Rivest
- [5] - "Improved algorithm for mining Apriori algorithm based on association rules" - Shabtay, L., Fournier-Viger, P., Yaari, R., & Dattne