

# Predictive Model for Cancer Detection Using Multifactorial Analysis: Examining the Impact of Lifestyle, Genetics, and Habits

Yajushrestha Shukla, Sanika Upasani and Swayam Pedgaonkar

Ms. Ruchi Jayaswal, Dr. Deepak Parashar

**Abstract**—The increasing worldwide worry about cancer gives us an immediate need for an ingenious prediction model for early diagnosis. This paper reveals a modern approach to cancer prediction based on multimodal analysis, using factors like diet (vegetarian/non-vegetarian), height, weight, blood groups, sex, marital status (married/ unmarried), smoking, and alcohol consumption.

Our model based on the above-mentioned factors predicts cancer risk, offering a complete and well-timed assessment. With the immense progress in data science, analytics, and machine learning, we aim to provide timely intervention and early detection capabilities.

This method visualizes a future where medical professionals and individuals themselves can far-sightedly assess cancer. It includes various elements that affect lifestyles, genetics, and behavioral distinctions to improve people's health and well-being.

By focusing on a user-centered approach and early mediation, our research stresses the importance of a comprehensive prediction model, holds promise for the development of innovative technologies, and provides for the development of existing healthcare systems.

This paper includes key research, methodological considerations, and expected outcomes for advancing public health through innovation, collaboration, and multidisciplinary predictive cancer diagnosis.

**Impact Statement** — In the realm of cancer prediction, the escalating global concern demands innovative solutions for early diagnosis. This paper addresses the pressing need by introducing a modern approach grounded in multimodal analysis,

encompassing lifestyle factors such as diet, height, weight, blood groups, marital status, smoking, and alcohol consumption. Unlike traditional models, our pioneering method goes beyond singular indicators, offering a holistic prediction framework. This novel approach not only significantly enhances early detection but also envisions a paradigm shift where individuals and healthcare professionals proactively assess and manage cancer risks. By leveraging data science and machine learning, our research propels the development of user-centered, comprehensive predictive models, promising transformative impact on public health and contributing to the evolution of healthcare systems.

**Index Terms**— Cancer Risk Prediction Model, Genetics and Cancer Prediction, Lifestyle and Cancer Risk, Multifactorial Analysis

## I. INTRODUCTION

Cancer, a word that has taken over 100 million human lives in the last decade [1] or so, is a disease that goes across borders and cultures, causing havoc on countless lives for generations. This persistent effort by this enemy has driven us to seek better and newer ways to fight it. With each passing day, we strive to understand it better, diagnose it earlier, and consequently conquer it.

In our fast moving, modern world, where time is of absolute importance and our health is cardinal, we found ourselves in need of an ingenious solution—a predictive model for early detection of cancer [2]. This research is the result of that need—the incorporation of our collective aspirations to face cancer on our terms.

The initiation of our exploration into this predictive model commences with the recognition that cancer manifests as a multifaceted entity, analogous to the complexities inherent in life itself, rather than existing in isolation. Its etiology is shaped by behavioral patterns, influenced by genetic predispositions, and intricately interwoven with lifestyle factors. Therefore, we must blend various factors to create a comprehensive prediction model for cancer.

Now, you may be curious about the nature of these diverse factors under discussion. They are quite known. We are talking about our eating habits (what we eat), whether we prefer a vegetarian or a non-vegetarian diet, our physical attributes like height and weight, our blood group, our gender, our marital status, and our other habits like whether we smoke and drink, and how often and how much.

With these factors in perspective, we start the road to prediction, offering an all-inclusive assessment of the risk of cancer [3]. The tools at our disposal have evolved, owing to the incredible developments in data science, analytics, and machine learning.

<sup>1</sup>This paragraph of the first footnote will contain the date on which you submitted your paper for review. It will also contain support information, including sponsor and financial support acknowledgment. For example, “This work was supported in part by the U.S. Department of Commerce under Grant BS123456.”

The next few paragraphs should contain the authors' current affiliations, including current address and e-mail. For example, F. A. Author is with the National Institute of Standards and Technology, Boulder, CO 80305 USA (e-mail: author@boulder.nist.gov).

S. B. Author, Jr., was with Rice University, Houston, TX 77005 USA. He is now with the Department of Physics, Colorado State University, Fort Collins, CO 80523 USA (e-mail: author@lamar.colostate.edu).

T. C. Author is with the Electrical Engineering Department, University of Colorado, Boulder, CO 80309 USA, on leave from the National Research Institute for Metals, Tsukuba, Japan (e-mail: author@nrim.go.jp).

This paragraph will include the Associate Editor who handled your paper.

# Predictive Model for Cancer Detection Using Multifactorial Analysis

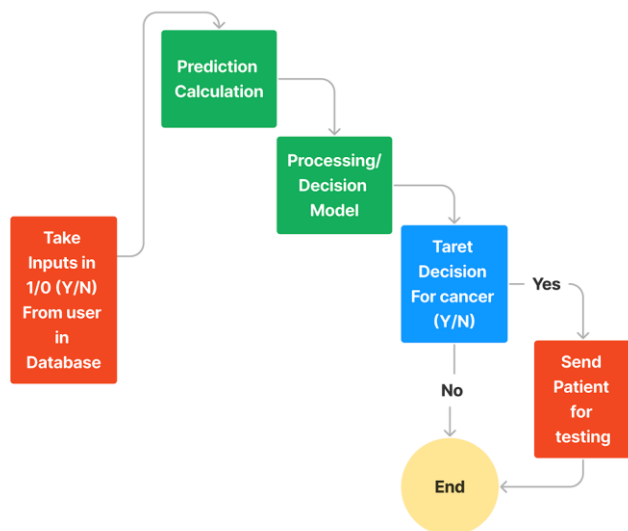
These technological advancements give us the opportunity to provide early detection capabilities and timely mediation.

Imagine a future where medical professionals and individuals themselves can have a perspective on their health. A future where the variations in our lives- the “what we eat, who we are, and how we live”—are examined carefully to provide us with a personalized alarm system for fatal diseases. We aim to move beyond the boundaries of a medicine based on the reaction to the arrival of conditions like cancer to a proactive approach where we anticipate its arrival and act accordingly.

Our approach to being user-centered believes that individuals and their well-being are to be placed as the topmost priority. We are not just attempting to predict cancer, we are also attempting to improve people's lives by contributing to their health and well-being. We want to traverse the gap between the intricacies of life, genetics, and behavioral separations, offering apprehensions that are both meaningful and practical.

Our research stresses the need for a comprehensive prediction model, one that summarizes the entirety of what makes us unique yet intertwined. Our control flow has been described in *Figure 1*. It is a call to visualize the possibilities, to see beyond the current line of advancements, and to consider a world where cancer is not a mountain we cannot climb but a challenge that we can tackle with a well-informed strategy.

In the pages that follow, we dive into this very approach to exploring the complexities of multifactorial predictive detection of cancer. We comb through existing research, throw light on existing methodological considerations, and reveal the expected outcomes of our journey. With an innovative mindset as our guide, collaboration as our compass, and serene public health as our final destination, we tread on this quest for a healthier and disease free future, a future where cancer is not an adversary but a challenge we can predict, prepare for, and tackle, a future where early detection is not a dream but a reality.



**Figure 1. Flow of the choice depending on the User Habits**

## II. LITERATURE REVIEW

[4] Xie et al. developed an advanced 2D CNN-based system for pulmonary nodule detection. By utilizing an upgraded 2D CNN

architecture to lower false positives and integrating a faster R-CNN for nodule classification, the system improved the CT reading process. Their approach achieved a sensitivity of 86.42% when evaluated in the LUNA16 dataset, presenting a substantial step forward in early lung cancer diagnosis. Jiang et al. [5] effectively delineated lung nodules through the utilization of a multi-patch methodology, amalgamating expertise from radiologists alongside the implementation of a four-channel neural network model. The method delivered a 94% sensitivity with 15.1 false positives per sample, signifying significant progress in lung nodule detection for early cancer diagnosis. Asuntha and Srinivasan [6] presented a deep learning strategy that involved extracting attributes using multiple methods and selecting the best features with the Fuzzy Particle Swarm Optimization algorithm. Their resulting dataset was employed to train a CNN with reduced computational complexity, achieving a remarkable sensitivity of 97.93%. This research signifies the potential of advanced feature extraction and selection techniques for improving lung nodule detection and advancing early lung cancer diagnosis.

This review [7] delves into the intricacies of oral squamous cell carcinoma (OSCC), exploring the changing landscape of its incidence, particularly among young patients. While traditionally linked to tobacco and alcohol consumption, the research highlights the challenge of ascertaining the underlying causes of OSCC in younger individuals. Genetic factors, predispositions to genetic instability, and the presence of risk factors except tobacco and alcohol are examined in-depth. The review emphasizes the importance of investigating potential environmental carcinogens, familial cancer history, viral infections, and stress as contributors to OSCC. Additionally, it discusses the application of non-invasive optical and photodynamic diagnostic methods, shedding light on potential advancements in OSCC detection and diagnosis [8].

Early detection is crucial to reducing lung cancer morbidity and mortality [9]. High-risk ever-smokers may benefit from screening using low-dose computed tomography (LDCT), according to the National Lung Screening Trial (NLST). However, the requirements for NLST selection are suboptimal, excluding many early-stage lung cancer patients based on smoking history and age [10]. Just 26.7% of US lung cancer cases meet NLST criteria, indicating room for improving precision in patient selection [11]. Alternative approaches, like stereotactic ablative body radiotherapy (SABR), show promise, particularly for stage I lung cancer patients ineligible for surgery [12]. Determining the optimal frequency of CT scans remains a challenge, considering concerns about radiation-induced cancers. Mitigating risks involves refining risk assessment tools and imaging protocols [13].

LDCT has limitations, including high false-positive rates and late-stage cancer emergence between screenings, necessitating complementary tests. Biomarkers like plasma microRNAs and circulating tumor cells (CTCs) have potential for reducing false positives and aiding in the diagnosis of aggressive cancers [14]. A deeper understanding of early molecular events in lung tumorigenesis may lead to novel biomarker development [15]. Although sputum analysis has regained interest, its role in the early identification of lung cancer is unknown. Challenges include study costs and developing reproducible testing

## Predictive Model for Cancer Detection Using Multifactorial Analysis

methods for emerging biomarkers. The future of lung cancer screening is expected to combine radiological and molecular methods, with surgical resection and alternative interventions like SABR offering treatment options. For patients unsuitable for these treatments, early diagnosis can aid in care planning and engagement with cancer services.

The study [16] delves into genetic marker analysis with the aim of improving early and precise lung cancer diagnosis. Exploring potential biomarkers, such as microsatellite changes, DNA hypermethylation, gene mutations (p53 and KRAS), and microRNA expression, the research investigates avenues for more effective lung cancer detection. It highlights the significant promise of microRNAs, particularly their expression profiles, as potential genetic markers capable of predicting lung cancer even 24 months earlier than traditional methods. Standardizing the quantification of circulating microRNAs [17] is emphasized as a crucial step for future clinical applications. This study aids in the early detection of lung cancer cases and the facilitation of prompt interventions.

This study [18] focused on identifying metabolic subtypes of gastric cancer and assessing their prognostic significance. Researchers carefully analyzed data from the Cancer Genome Atlas (TCGA) database [19], and clinical follow-up information. A cohort of 375 tumor samples and 32 normal samples was assembled by means of rigorous data processing. Utilizing genes related to glycolysis and cholesterol from the MsigDB database, a clustering approach revealed the existence of four distinct metabolic subtypes: cholesterol, quiescent, glycolysis, and mixed [20, 21]. The analysis of patient prognosis unveiled notable disparities, with the cholesterol subtype displaying a less favorable outlook compared to the glycolysis subtype. These differences were further supported by observed variations in the expression patterns of cholesterol and glycolysis genes within these subtypes. Significantly, tumor samples exhibited elevated expressions of both cholesterol and glycolysis genes when contrasted with normal samples, implying a correlation between the upregulation of these genes and the development of gastric cancer. This study has uncovered the presence of discrete metabolic subtypes in gastric cancer and underlined their potential as prognostic markers, shedding light on the intricate interplay between cholesterol and glycolysis genes and their association with the disease. These findings hold promise for advancing our understanding of gastric cancer biology and potentially guiding the development of targeted therapeutic strategies.

For the purpose of assessing the risk of colorectal cancer (CRC), the study [22] utilized data from the Darmkrebs: Chancen der Verhütung durch Screening (DACHS) project, a continuous population-based case-control investigation that has been conducted in southwest Germany since 2003. Eligibility for participation was determined by having a confirmed initial diagnosis of CRC, with participants aged at least 30 years and conversant in German. The control group was selected randomly from population registries, employing age, sex, and county of residence as matching criteria. Thorough interviews were conducted to capture lifestyle and medical information, encompassing factors such as body mass index, eating habits,

drinking, smoking, and physical activity. Genomic data were obtained from either blood samples or buccal cells, according to the participants, with DNA extraction performed via standard procedures. A significant factor in the CRC risk assessment was the recent colonoscopy history [23]. To ascertain the association between CRC risk and lifestyle factors, polygenic risk scores, and recent colonoscopy, the study employed multiple logistic regression, taking sex and age into account. The investigation also encompassed the estimation of the 30-year absolute risk of developing CRC in 50-year-old individuals, founded on diverse risk profiles [24]. Age-specific cancer hazard rates from the German Centre for Cancer Registry Data, Robert Koch Institute, along with relative risk data extracted from case-control studies, were utilized to account for an individual's age and risk factors in estimating the probability of colorectal cancer (CRC) development over a 30-year timeframe. Sensitivity analyses were carried out with varying relative risks for colonoscopy history. In summary, the study, involving 4220 CRC patients and 3338 control participants, provided invaluable insights into the interconnectedness of lifestyle components, genetic predisposition, and colonoscopy history in the context of assessing CRC risk.

In this study [25], the total contribution of alcohol-based cancer in 2020 was found. The method involved the selection of types of cancer with a significant result of a causal relationship with alcohol consumption, using data from the IARC and GIS on Alcohol and Health. Population Attributable Fraction (PAFs) were calculated based on lifetime sobriety and 2010 alcohol consumption data, along with estimates from the WCRF Continuous Update Project. A systematic literature review was used for both current and past drinking to obtain relative estimates. The effect of alcohol consumption on cancer was assessed using PAF, factoring in different levels of alcohol consumption, and graded by country, sex, and cancer site. A Monte-Carlo like approach was used to estimate uncertainty, and it determined 95% uncertainty intervals. The findings revealed an astonishing estimated 741,300 new cancer cases in 2020 (4.1% of all new cases) were because of alcohol consumption, with the highest PAFs observed for oesophageal, pharyngeal, and lip cavity cancer. The study stresses the worldwide spread of cancer owing to alcohol consumption, highlighting the need for very effective and immediate policies and interventions to reduce alcohol-related cancer risks and overall alcohol consumption.

In the study [26], 3309 participants enrolled in a currently going prospective cohort study, and their data was analyzed, which focused on the possible connection between lifestyle and psychosocial job qualities associated with cancer risk. The study gauged job demands, job strain, and iso-strain to investigate their influence on cancer risk, along with lifestyle activities like smoking, consuming alcohol, eating fruits and vegetables, and engaging in physical activity [27]. The results did not reveal significant inconsistencies in the presence of lifestyle risk factors for cancer, among other work characteristics that were investigated [28]. Furthermore, there was little evidence to indicate a link between occupational strain, iso-strain, and lifestyles linked to cancer when

## Predictive Model for Cancer Detection Using Multifactorial Analysis

multivariate analysis was used. Finally, the outcomes of the study do not substantiate the notion that job strain, or iso-strain is associated with adopting a cancer related lifestyle.

The study population involved individuals who were employed at the time of completing the Health and Life Experiences Questionnaire, which is the reason for the higher representation of younger participants with an elevated level of education. The analysis of sociodemographic variables and lifestyle characteristics, stratified by gender, revealed considerable differences, particularly in alcohol consumption, where men displayed a notably higher daily intake than women. Furthermore, the results did not identify a pronounced association between lifestyle factors, occupational strain, and iso-strain across the years 1993 and 1996 [27, 28]. Conclusively, this study has probed the intricate relationship between psychosocial job characteristics and cancer-related lifestyle risk factors, yielding inconclusive findings. The study serves as a stepping stone for deeper research into other psychosocial aspects and their potential mediation in the complex interplay between work environment and health behaviors. While the investigation adds valuable insights to this understudied area, it also highlights the imperative need for future comprehensive research endeavors.

With a 5-year survival rate of less than 1/5th, esophageal cancer is a devastating and complex disease primarily brought on by delayed diagnosis. The vast landscape of esophageal cancer presents numerous challenges, with varied commonalities and risk factors across different populations. While there have been significant risk factors like tobacco, alcohol, and reflux esophagitis for a long time, dietary practices and nutrition have an influence. This review provided us with a thorough summary of how diet and nutrition affect the risk of esophageal cancer, which is still less understood.

Further study resulted in the understanding that the causation of the disease was attributed to a wide range of dietary elements like fruits, vegetables, vitamins and minerals, fats, meats, salted foods, carcinogens, nitrogen compounds, mycotoxins, and also the temperature of the food consumed. Using this derived understanding is crucial for advancing our strategies to tackle esophageal cancer and devising effective prevention plans. By highlighting the importance of tobacco and alcohol cessation, a diet absolutely rich in vegetables and fruits, and the significant role of nutrition in at-risk populations, this paper [29] highlights the need for comprehensive efforts to reduce the impact of esophageal cancer on a broader population.

In this study [30], they aimed to scrutinize the relationship between extended exposure to air pollution and the degree of lung cancer among Koreans, since it is Korea's primary cause of cancer-related mortality. They observed an increasing trend in lung cancer cases, especially with shifts in biological types, such as an increase in the presence of adenocarcinomas, and a decrease in squamous cell carcinomas. Even though the key cause of lung cancer is tobacco use, high lung cancer rates amongst non and never smokers have raised questions about other factors like occupational exposures, environmental

tobacco dust, lower socioeconomic status, and air pollution. Notably, studies have connected the risk of lung cancer to air pollution, which includes nitrogen dioxide (NO<sub>2</sub>) and particulate matter with a diameter of less than 10 micrometers (PM<sub>10</sub>). In a California cohort study of persons who did not smoke, there was a positive correlation found between male incident lung cancer cases and PM<sub>10</sub> levels. and European cohorts reported a similar link, emphasizing the impact even below certain PM<sub>10</sub> levels. Additionally, a case-control study in Canada demonstrated that exposure to NO<sub>2</sub> was associated with increased odds of lung cancer. Importantly, this study delved into the risk factors and considered specific histological subtypes of lung cancer. This study offers critical insights into the role of air pollution in the emergence of lung cancer, emphasizing the importance of considering histological subtypes when assessing its association with air pollution.

**Table 1: Summarized Form of Literature Review**

## Predictive Model for Cancer Detection Using Multifactorial Analysis

Reference number	Method used	Dataset used	Finding/ Limitation/ Future Scope/Performance parameters
[4]	3 CNN in MCNN	LIDC-IDRI	1. Enhanced nodule detection. 2. High sensitivity, low FP. 3. Diagnostic accuracy: 86.84%.
[7]	Cytologic diagnosis during OSCC development	Web-based search initiated using Medline/PubMed	Genetics, alcohol, smoking, carcinogens, and various factors contribute to oral cancer.
[9]	Low Dose CT scan	NLST selections	Increasing population selection accuracy; SABR shows potential for non-surgical patients.
[16]	Genetic Marker Analysis	Biomarkers Studied	Potential microRNA biomarkers
[18]	Data processing, classification.	TCGA gastric cancer samples dataset.	Metabolic subtypes impact via gene expression variations
[22]	Data analysis and calculations	DACHS population-based study	Lifestyle, colonoscopy reduce CRC risk with risk reduction measures.
[26]	Cross-sectional survey analysis	Dutch cohort study	Logistic, linear regression found no job-cancer link; future explores psychosocial factors
[30]	Case-control study, regression models	908 lung cancer patients	Air pollution links, odds ratios, confidence intervals for lung cancer

### III. PROBLEM STATEMENT

The problem at hand is the need for a prediction model that should be capable of early detection of cancer by utilizing multifactorial analysis, comprising variables like dietary habits, weight, blood group, gender, marital status, smoking, and alcohol consumption. This model must ease all-inclusive assessments of cancer risk, bridging the gap between lifestyle, habits, and genetics. Leveraging the development of data science, analytics, and machine learning, this research is directed at providing a complete, user centered approach to predicting and preventing the onset of cancer. The final goal is to entitle individuals and medical professionals to a timely mediation plan of action, thus contributing to public welfare and well-being and completely changing the landscape of cancer diagnosis and prevention.

### IV. OBJECTIVE

This research aims to develop a multifactorial cancer prediction model, enabling early detection and user-centered

interventions, ultimately enhancing public health through personalized risk assessment and timely action.

- Development of a Multifaceted Cancer Prediction Model:** The primary objective is to create and introduce a novel predictive model that encompasses multiple predictive factors, including diet, weight, height, blood type, gender, marital status, smoking, and alcohol intake. This model will offer a multifaceted approach to assessing
- Cancer risk prediction with a focus on predicting the probability of cancer based on the mentioned factors.
- Early Detection and Timely Intervention:** This research contributes to the world of early detection of cancer based on the advancements in the fields of data science, analytics, and machine learning. The aim of this study is to empower medical professionals and individuals themselves to ardently assess their risk of cancer, therefore allowing them to intervene when necessary.
- User-Centered Approach:** The study emphasizes the significance of a user-centered approach, acknowledging that individual users are at the core of the healthcare process. By examining lifestyle, genetics, and behavior, this research aims to offer personalized cancer risk insights. This user-focused approach is anticipated to enhance public health by

# Predictive Model for Cancer Detection Using Multifactorial Analysis

making cancer prediction more relevant and accessible to individuals.

## V. EXPERIMENTAL SETUP

The research necessitates a specialized experimental configuration comprising specific hardware components: an Intel Core i5 or AMD Ryzen 5 5500 processor with six cores, a minimum of 16 GB RAM, and a basic GPU with at least 2 GB of VRAM, and a storage capacity of at least 1 GB. Python is designated as the essential programming language for seamless integration with prevalent machine learning and data analysis libraries. This tailored configuration ensures that computational operations are carried out as efficiently as possible, creating an atmosphere that is favorable to the successful investigation and examination of research goals, ultimately guaranteeing the acquisition of reliable results.

## VI. DATASET DESCRIPTION

The Health Assessment Dataset comprises information gathered from 8322 individuals, serving as a robust resource for evaluating various medical conditions and associated risks. A patient ID serves as a unique identifier for each dataset item, which includes a wide range of lifestyle, environmental, and health-related characteristics. These comprise numerical data like alcohol consumption levels and BMI, as well as categorical variables like age, gender, and smoking status. The information also contains binary indicators for genetic vulnerabilities to certain diseases and environmental exposures like air pollution. As categorical variables, symptoms pertaining to cardiovascular and respiratory health are also recorded. As categorical variables, symptoms pertaining to cardiovascular and respiratory health are also recorded. The total health assessment level of each individual is shown as ordinal data and is divided into three categories: low, medium, and high. With a manageable size of 1000 KB, this dataset serves as a valuable resource for conducting analyses, discerning patterns, and predicting potential health risks based on various factors.

Risk Factors for lung cancer are as follows:

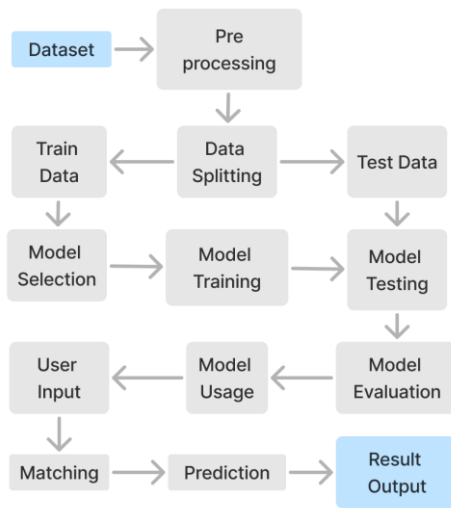
- **Gender:** This field represents the gender of individuals in the dataset, typically categorized as male or female. Gender can be a significant factor in assessing lung cancer risk.
- **Age:** Age is a number characteristic that represents each person's age within the dataset. Given that the incidence of lung cancer tends to rise with age, age has a significant role in determining one's risk for the disease.
- **Smoking:** This categorical attribute indicates whether an individual is a smoker, typically classified as a binary choice, such as "smoker" or "non-smoker." Smoking is a well-established risk factor for lung cancer.
- **Yellow fingers:** This attribute may describe whether an individual has yellow fingers, which can be associated with smoking due to nicotine staining.

- **Anxiety:** Anxiety is a categorical attribute that signifies whether an individual experiences anxiety, a psychological factor that may contribute to smoking and other behaviors related to lung cancer risk.
- **Peer pressure:** Peer pressure is a categorical attribute that can reveal whether or not a person is persuaded by their peers to take up risky behaviors like smoking or other habits.
- **Chronic disease:** This field may describe whether an individual has any chronic diseases or conditions, which can impact their overall health and, indirectly, their lung cancer risk.
- **Fatigue:** Fatigue is a categorical attribute that can signal whether an individual experiences excessive tiredness, which might be linked to various lifestyle factors and health conditions.
- **Allergy:** This attribute denotes whether an individual has allergies, which can impact respiratory health and potentially influence lung cancer risk.
- **Wheezing:** Wheezing is a categorical attribute that can signify whether an individual experiences wheezing or breathing difficulties, which could be related to lung health.
- **Alcohol consumption:** This field indicates whether an individual consumes alcohol. Alcohol consumption can interact with smoking and other factors in lung cancer risk assessment.
- **Coughing:** Coughing is a categorical attribute that represents whether an individual experiences persistent coughing, which can be related to respiratory health.
- **Shortness of breath:** Shortness of breath is a categorical attribute that can denote whether an individual experiences difficulty breathing, which is another important respiratory health factor.
- **Swallowing difficulty:** This attribute may describe whether an individual has difficulty swallowing, which could be linked to various health issues.
- **Chest pain:** Chest pain is a categorical attribute that signifies whether an individual experiences chest pain, which can have various causes, including those related to lung health.
- **Lung cancer:** This binary attribute is the target variable, indicating whether an individual has been diagnosed with lung cancer or not. It's the variable of primary interest in assessing the risk factors associated with lung cancer.

## VII. PROPOSED METHODOLOGY



# Predictive Model for Cancer Detection Using Multifactorial Analysis



**Figure 2. Proposed Methodology**

The initial phase of this research study is marked by outlining the scope of the study, objectives and primary sources of data. This initial phase is described in *Figure 2*. The subsequent phase involves collecting essential patient data through a structured questionnaire. This data encompasses critical elements such as age, dietary patterns, smoking history, fatigue levels, allergies, familial cancer history, respiratory conditions, and other pertinent variables. The comprehensiveness and quality of this data directly influence the model's precision and performance.

Once the data has been gathered, it is divided into training and testing subsets to prevent overfitting and enable the model to generalize to new, unseen data effectively. Data preprocessing is a pivotal step in cleaning and preparing the data for analysis. This process includes handling missing values, addressing outliers, and standardizing data formats. For instance, in medical data, addressing missing values in patient records is crucial for a complete dataset.

Categorical data exploration entails determining the distribution and linkages within the dataset. Visualizing the spread of different cancer kinds, for example, can reveal trends that might help with early identification. Analyzing age-related data reveals information about age groups' susceptibility to cancer and other important patterns, which informs feature selection and model design.

A comparison of age and categorical variables reveals correlations that influence the model's predictive ability. These visualizations show how risk variables change with age and allow for more educated feature selection. Creating a heatmap to visualize variable correlations can help with feature selection and identifying the features that have a major influence on cancer prediction. For example, research might indicate associations between lifestyle decisions and cancer risk.

Addressing class imbalance by oversampling guarantees that the model does not favor the dominant class, which is crucial for cancer prediction in areas where positive cases are rare. Data

splitting into training and testing sets provides model validation, whereas data scaling is critical for certain model types.

Choosing an appropriate machine learning model entails testing with several algorithms, followed by training and fine-tuning for optimum performance. Model evaluation is comparing the model's predictions to actual outcomes using measures such as precision, accuracy, recall, and F1-score to determine its usefulness in predicting cancer cases. In the last step, we explore future development opportunities, such as incorporating real-time patient data and adding additional features to improve prediction precision and relevance.

## VIII. RESULT ANALYSIS

Our findings are summarized in the following tables, graphs, and figures. Each provides a different aspect of the analysis and model performance.

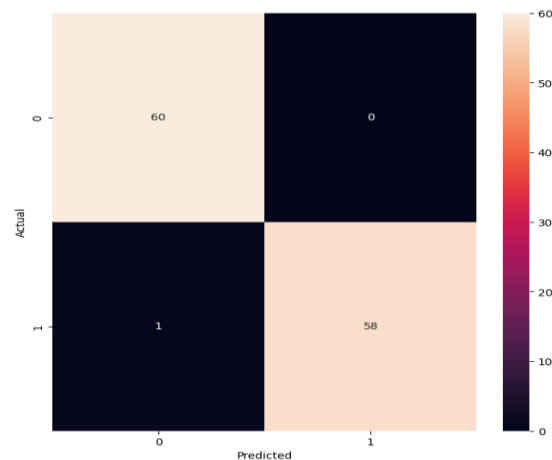
The analysis compares various classification models, each evaluated based on accuracy, F1 score, precision, and recall.

### Model Comparison/Selection

#### 1) Support Vector Machine (SVM):

Type: Supervised Learning

The Support Vector Machine (SVM) model demonstrates exceptional accuracy (**99%**), achieving near-perfect precision (**98%**) and recall (**99%**) for both class 0 and class 1 instances. This makes it highly suitable for binary and multiclass classification in high-dimensional spaces.



**Figure 3. SVM Benchmarks**

The above *Figure 3* refers to classes 0 and 1 as SVM benchmarks, which have been discussed

#### 2) K-Nearest Neighbors (KNN):

Type: Supervised Learning

K-Nearest Neighbors (KNN) displays excellent accuracy and a

strong balance between precision and recall, particularly excelling in accurately identifying class 0 instances. Logistic Regression offers respectable performance with a balanced trade-off between precision and recall.

3)Random Forest:

Type: Ensemble Learning

The Random Forest model excels in achieving excellent recall and precision for both class 0 and class 1 instances. Its ensemble approach, utilizing multiple decision trees, effectively reduces overfitting while maintaining high accuracy. This makes Random Forest particularly robust and suitable for complex classification tasks.

4)Gradient Boosting:

Type: Ensemble Learning

theGradient Boosting model demonstrates exceptional accuracy and precision, with a focus on minimizing prediction errors through the additive model it builds. It achieves a good balance between precision and recall for both class 0 and class 1 instances, showcasing its effectiveness in various predictive

Model	Precision	Recall	F-1 Score	ROC AUC
SVM	98%	99%	98%	98.31%
RANDOM	95%	94%	94%	98.39%
GRADIENT	95%	95%	95%	97.46%

modeling tasks.

Table 2: Model Comparison

Based on the model comparison presented in Table 2, the Support Vector Machine (SVM) stands out as the most suitable choice for the classification task. SVM achieves the highest precision, recall, and F1 score among all models, with precision and recall rates of **98%** and **99%** respectively, resulting in an overall F1 score of **98%**. Additionally, SVM demonstrates a robust performance with an ROC AUC of **98.31%**.

This exceptional performance in precision, recall, and F1 score indicates SVM's ability to effectively classify instances across different classes with high accuracy while maintaining a strong balance between precision and recall. The high ROC AUC further validates its capability to distinguish between positive and negative classes efficiently.

Therefore, based on its superior performance across multiple metrics, particularly in precision and recall, **SVM** is chosen as the right model to use for the classification task. It offers a reliable and accurate solution for effectively categorizing instances within the dataset.

IX. COMPARATIVE ANALYSIS

In our paper, we utilized the Surveyed\_Cancer dataset, employing Python and Orange as our tools of choice. We implemented **SVM** with plans to further enhance its performance through **ADA Boost**. The achieved accuracy of **98.13%** underscores the effectiveness of our approach. Our dataset offers several advantages, notably being daily life-based, which enhances its relevance and applicability. Additionally, it is easily accessible, contributing to its widespread usability. Moreover, the dataset's characteristics lean towards being more user-dependent rather than machine-dependent, facilitating easier interpretation and utilization of the findings.

In [31], the study focused on the Wisconsin Dataset using WEKA as the tool and implementing the Random Forest technique, resulting in an accuracy of 92.20%. This dataset is characterized by its substantial size, containing millions of patient data records. While this extensive dataset offers a rich source of information, its sheer volume can pose challenges in terms of processing and analysis, potentially requiring more computational resources.

Similarly, [31] examined the Wisconsin Diagnostic Breast Cancer dataset using Python. Various techniques were compared, including SVM, KNN, decision trees, and Naive Bayes, with an accuracy of 96.91%. This dataset is notable for its amalgamation of similar databases under the Wisconsin datasets umbrella, providing a comprehensive repository of relevant information for breast cancer diagnosis and research.

Comparatively, our dataset stands out due to its more user-friendly data fields and findings. Unlike the heavy dataset of the Wisconsin Dataset, our dataset offers a more accessible and manageable dataset. Its focus on daily life-based data makes it more relatable and applicable to real-world scenarios, potentially leading to more actionable insights and outcomes in cancer research and diagnosis.

Paper	Dataset Used	Technique Used	Accuracy
Our Paper	Surveyed_Cancer dataset	SVM	98.13%
		(ADABOOST In further Steps)	
[31]	Wisconsin Dataset	Random Forest	92.20%
[31]	Wisconsin Diagnostic Breast Cancer dataset	SVM vs KNN,	96.91%
		Decision Trees, Naive Bayes	

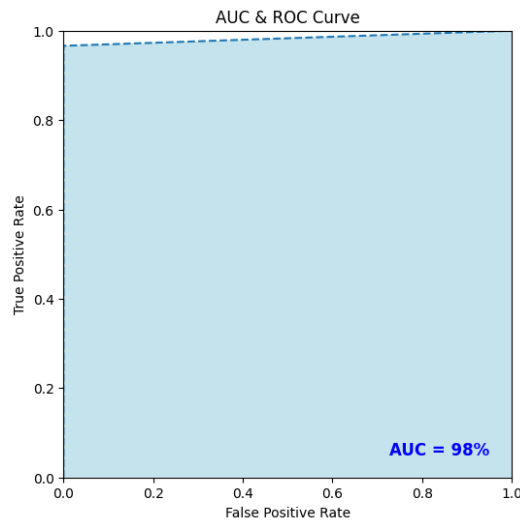


## X. RESULTS AND DISCUSSION

We employed various data preprocessing techniques, tested multiple machine learning models, and selected the optimal model based on precision. The tables, graphs, plots, heat maps, ROC curves, and box plots provide a comprehensive overview of our results.

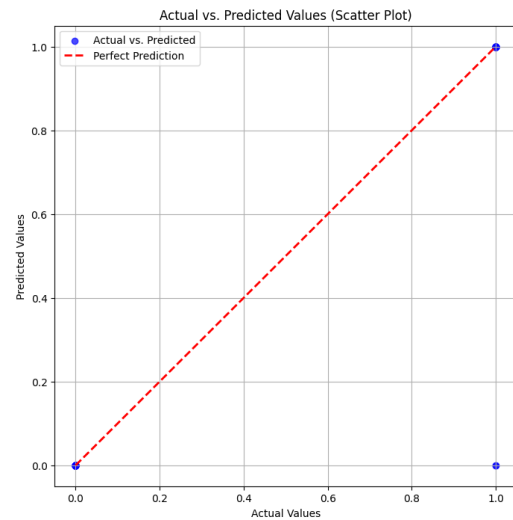
**ROC Curve (Receiver Operating Characteristic Curve):** The ROC curve is graphical depiction of a model's capacity for class differentiation. Within the Random Forest framework, the ROC curve facilitates the visualization of the trade-off between true positive rate (sensitivity) and false positive rate (1-specificity) by modifying the model's classification threshold.

**AUC (Area Under the ROC Curve):** AUC is a single metric that quantifies a Random Forest model's overall performance. The area under the ROC curve is measured, with a greater AUC indicating better model performance. AUC values vary from 0 to 1, with 0.5 representing chance and 1 representing a perfect classifier. A high AUC in the context of Random Forest indicates that the model is competent at differentiating between classes and producing accurate predictions



**Figure 4. AUC & ROC Curve For SVM**

The ROC curve and AUC analysis, as shown in *Figure 9*, illustrate the excellent performance of the SVM model. The curve's trajectory, which leans towards the top-left corner, indicates the model's strong ability to distinguish between classes. The AUC value, quantified at 98%, underscores the SVM model's exceptional precision in classifying data. This visual representation, along with the high AUC percentage, confirms the model's robustness in distinguishing between the two classes, making it a reliable choice for your research paper.



**Figure 5 Actual vs Predicted Values (Scatter Plot)**

In the scatter plot, as shown in *Figure 10*, actual target values were compared with the predictions made by the SVM model. The close alignment of the blue dots with the red dashed "Perfect Prediction" line suggests the SVM model's strong predictive accuracy. Additionally, the SVM model demonstrates an impressive ROC AUC (Receiver Operating Characteristic Area Under the Curve) value of **98.31%**. This high AUC value signifies the model's excellent ability to distinguish between different classes. These visualizations and metrics highlight the SVM model's reliability and precision for the classification task, making it a valuable choice for your research paper.

## XI. CONCLUSION

In summary, our research was dedicated to the development of a robust model for the early detection of cancer in patients. Our method was anchored by thorough data preprocessing, feature selection, and a systematic examination of several machine learning algorithms. We systematically tested and compared several models, including the LGBM, Random Forest, Logistic Regression, Support Vector Machine (SVM), and Gradient Boosting, with the objective of identifying the most effective model. After a thorough assessment, our findings revealed that the Support Vector Machine (SVM) model exhibited superior performance when compared to other models. It achieved a remarkable accuracy rate of 99% and a recall rate of 98%, which indicates its exceptional ability to accurately discern patients with cancer while maintaining a very low misclassification rate. Furthermore, our study included a comprehensive analysis of ROC curves and AUC scores. The AUC score of 98.31% exemplifies the SVM model's outstanding predictive capabilities and its capacity to effectively differentiate between cancer and non-cancer cases. Our research demonstrates that the meticulous application of machine learning techniques, combined with rigorous data analysis, holds significant promise in advancing early cancer detection. These results underline the potential for these models to provide invaluable support to healthcare professionals, contributing to timely and precise diagnoses, thereby enhancing patient outcomes. This research provides a robust foundation

# Predictive Model for Cancer Detection Using Multifactorial Analysis

for further exploration within the realm of medical diagnostics and underscores the transformative potential of data-driven methodologies in the field of healthcare.

## XII. FUTURE SCOPE

Building on current progress, we anticipate advancements in early cancer detection and personalized healthcare. Our team aims to create an innovative diagnostic gadget that analyzes blood, saliva, and sweat for a comprehensive health assessment. The objective is a non-invasive, user-friendly instrument for real-time analysis, revolutionizing early detection, reducing testing costs, and enhancing healthcare accessibility, especially in remote or impoverished areas.

To construct this device, we'll integrate cutting-edge technologies like artificial intelligence and sensors. Early diagnosis and timely intervention can significantly impact outcomes, potentially saving lives. We strive to usher in an era where individuals monitor their health proactively, and healthcare professionals rely on improved, non-invasive instruments for early disease detection. Our commitment is to lead this disruptive approach, ultimately improving patient outcomes and contributing to a more sustainable healthcare system.

## REFERENCES

- [1] <https://www.who.int/news-room/fact-sheets/detail/cancer>
- [2] David Crosby et al. .Early detection of cancer. Science375, eaay9040(2022) .DOI:10.1126/science.aay9040
- [3] Ritchie, J. B., Welch, B. M., Allen, C. G., Frey, L. J., Morrison, H., Schiffman, J. D., Alekseyenko, A. V., Dean, B., Halbert, C. H., & Bellcross, C. (2022, May 16). Comparison of a Cancer Family History Collection and Risk Assessment Tool – ItRunsInMyFamily – with Risk Assessment by Health-Care Professionals. *Public Health Genomics*, 25(3-4), 80–88. <https://doi.org/10.1159/000520001>
- [4] Saba, T. (2020). Recent advancement in cancer detection using machine learning: Systematic survey of decades, comparisons and challenges. *Journal of Infection and Public Health*, 13(9), 1274–1289. <https://doi.org/10.1016/j.jiph.2020.06.033>.
- [5] H. Jiang, H. Ma, W. Qian, M. Gao and Y. Li, "An Automatic Detection System of Lung Nodule Based on Multigroup Patch-Based Deep Learning Network," in *IEEE Journal of Biomedical and Health Informatics*, vol. 22, no. 4, pp. 1227–1237, July 2018, doi: 10.1109/JBHI.2017.2725903.
- [6] Asuntha, A., Srinivasan, A. Deep learning for lung Cancer detection and classification. *Multimed Tools Appl* 79, 7731–7762 (2020). <https://doi.org/10.1007/s11042-019-08394-3>
- [7] Zygogianni, A.G., Kyrgias, G., Karakitsos, P. et al. Oral squamous cell cancer: early detection and the role of alcohol and smoking. *Head Neck Oncol* 3, 2 (2011). <https://doi.org/10.1186/1758-3284-3-2>
- [8] Harris, D. M., & Werkhaven, J. (1987). Endogenous porphyrin fluorescence in tumors. *Lasers in surgery and medicine*, 7(6), 467–472. <https://doi.org/10.1002/lsm.190007060>
- [9] Knight, S., Crosbie, P., Balata, H., Chudziak, J., Hussell, T., & Dive, C. (2017). Progress and prospects of early detection in lung cancer. *Open Biology*, 7(9), 170070. <https://doi.org/10.1098/rsob.170070>
- [10] Taiwo, E. O., Yorio, J. T., Yan, J., & Gerber, D. E. (2012). How have we diagnosed early-stage lung cancer without radiographic screening? A contemporary single-center experience. *PloS one*, 7(12), e52313. <https://doi.org/10.1371/journal.pone.0052313>
- [11] Pinsky, P. F., & Berg, C. D. (2012). Applying the National Lung Screening Trial eligibility criteria to the US population: what percent of the population and of incident lung cancers would be covered?. *Journal of medical screening*, 19(3), 154–156. <https://doi.org/10.1258/jms.2012.012010>
- [12] Zheng, X., Schipper, M., Kidwell, K., Lin, J., Reddy, R., Ren, Y., Chang, A., Lv, F., Orringer, M., & Spring Kong, F. M. (2014). Survival outcome after stereotactic body radiation therapy and surgery for stage I non-small cell lung cancer: a meta-analysis. *International journal of radiation oncology, biology, physics*, 90(3), 603–611. <https://doi.org/10.1016/j.ijrobp.2014.05.055>
- [13] Rampinelli, C., De Marco, P., Origgi, D., Maisonneuve, P., Casiraghi, M., Veronesi, G., Spaggiari, L., & Bellomi, M. (2017). Exposure to low dose computed tomography for lung cancer screening and risk of cancer: secondary analysis of trial data and risk-benefit analysis. *BMJ (Clinical research ed.)*, 356, j347. <https://doi.org/10.1136/bmj.j347>
- [14] Sozzi, G., Boeri, M., Rossi, M., Verri, C., Suatoni, P., Bravi, F., Roz, L., Conte, D., Grassi, M., Sverzellati, N., Marchiano, A., Negri, E., La Vecchia, C., & Pastorino, U. (2014). Clinical utility of a plasma-based miRNA signature classifier within computed tomography lung cancer screening: a correlative MILD trial study. *Journal of clinical oncology : official journal of the American Society of Clinical Oncology*, 32(8), 768–773. <https://doi.org/10.1200/JCO.2013.50.4357>
- [15] Kwon, M. C., & Berns, A. (2013). Mouse models for lung cancer. *Molecular oncology*, 7(2), 165–177. <https://doi.org/10.1016/j.molonc.2013.02.010>
- [16] Wadowska K, Bil-Lula I, Trembecki Ł, Śliwińska-Mossoń M. Genetic Markers in Lung Cancer Diagnosis: A Review. *International Journal of Molecular Sciences*. 2020; 21(13):4569. <https://doi.org/10.3390/ijms21134569>
- [17] Shen, C., Wang, X., Tian, L., & Che, G. (2014). Microsatellite alteration in multiple primary lung cancer. *Journal of thoracic disease*, 6(10), 1499–1505. <https://doi.org/10.3978/j.issn.2072-1439.2014.09.14>
- [18] Zhu, Z., Qin, J., Dong, C., Yang, J., Yang, M., Tian, J., & Zhong, X. (2021). Identification of four gastric cancer subtypes based on genetic analysis of cholesterogenic and glycolytic pathways. *Bioengineered*, 12(1), 4780–4793. <https://doi.org/10.1080/21655979.2021.1956247>
- [19] Wang, Z., Jensen, M. A., & Zenklusen, J. C. (2016). A Practical Guide to The Cancer Genome Atlas (TCGA). *Methods in molecular biology (Clifton, N.J.)*, 1418, 111–141. [https://doi.org/10.1007/978-1-4939-3578-9\\_6](https://doi.org/10.1007/978-1-4939-3578-9_6)
- [20] Jiang, J., Zheng, Q., Zhu, W., Chen, X., Lu, H., Chen, D., Zhang, H., Shao, M., Zhou, L., & Zheng, S. (2020). Alterations in glycolytic/cholesterogenic gene expression in hepatocellular carcinoma. *Aging*, 12(11), 10300–10316. <https://doi.org/10.18632/aging.103254>
- [21] Karasinska, J. M., Topham, J. T., Kalloger, S. E., Jang, G. H., Denroche, R. E., Culibrk, L., Williamson, L. M., Wong, H. L., Lee, M. K. C., O'Kane, G. M., Moore, R. A., Mungall, A. J., Moore, M. J., Warren, C., Metcalfe, A., Notta, F., Knox, J. J., Gallinger, S., Laskin, J., Marra, M. A., ... Schaeffer, D. F. (2020). Altered Gene Expression along the Glycolysis-Cholesterol Synthesis Axis Is Associated with Outcome in Pancreatic Cancer. *Clinical cancer research : an official journal of the American Association for Cancer Research*, 26(1), 135–146. <https://doi.org/10.1158/1078-0432.CCR-19-1543>
- [22] Carr, P. R., Weigl, K., Edelmann, D., Jansen, L., Chang-Claude, J., Brenner, H., & Hoffmeister, M. (2020). Estimation of Absolute Risk of Colorectal Cancer Based on Healthy Lifestyle, Genetic Risk, and Colonoscopy Status in a Population-Based Study. *Gastroenterology*, 159(1), 129–138.e9. <https://doi.org/10.1053/j.gastro.2020.03.016>
- [23] Brenner, H., Stock, C., & Hoffmeister, M. (2014). Effect of screening sigmoidoscopy and screening colonoscopy on colorectal cancer incidence

# Predictive Model for Cancer Detection Using Multifactorial Analysis

and mortality: systematic review and meta-analysis of randomized controlled trials and observational studies. *BMJ*, 348, g2467. doi:10.1136/bmj.g2467

[24] Pfeiffer, R. M., & Petracci, E. (2011). Variance computations for functionals of absolute risk estimates. *Statistics & Probability Letters*, 81(7), 807-812. <https://doi.org/10.1016/j.spl.2011.02.002>

[25] Rumgay, H., Shield, K., Charvat, H., Ferrari, P., Sornpaisarn, B., Obot, I., et al. (2021). Global burden of cancer in 2020 attributable to alcohol consumption: a population-based study. *The Lancet Oncology*, 22(8), 1071-1080. [https://doi.org/10.1016/S1470-2045\(21\)00279-5](https://doi.org/10.1016/S1470-2045(21)00279-5)

[26] A Jeanne M van Loon, Marja Tijhuis, Paul G Surtees, Johan Ormel, Lifestyle risk factors for cancer: the relationship with psychosocial work environment, *International Journal of Epidemiology*, Volume 29, Issue 5, October 2000, Pages 785–792, <https://doi.org/10.1093/ije/29.5.785>

[27] Pols MA, Peeters PHM, Ocké MC, Slimani N, Bueno de Mesquita HB, Collette HJA. Estimation of reproducibility and relative validity of the questions included in the EPIC physical activity questionnaire. *Int J Epidemiol* 1997;26:S181–89.

[28] Baecke JAH, Burema J, Frijters JER. A short questionnaire for the measurement of habitual physical activity in epidemiological studies. *Am J Clin Nutr* 1982;36:936–

[29] A. G. Palladino-Davis, B. M. Mendez, P. M. Fisichella, C. S. Davis, Dietary habits and esophageal cancer, *Diseases of the Esophagus*, Volume 28, Issue 1, 1 January 2015, Pages 59–67, <https://doi.org/10.1111/dote.12097>

[30] Lamichhane, D. K., Kim, H. C., Choi, C. M., Shin, M. H., Shim, Y. M., Leem, J. H., Ryu, J. S., Nam, H. S., & Park, S. M. (2017). Lung Cancer Risk and Residential Exposure to Air Pollution: A Korean Population-Based Case-Control Study. *Yonsei medical journal*, 58(6), 1111–1118. <https://doi.org/10.3349/ymj.2017.58.6.1111>

[31] Keles, M. Kaya, "Breast Cancer Prediction and Detection Using Data Mining Classification Algorithms: A Comparative Study." *Tehnicki Vjesnik - Technical Gazette*, vol. 26, no. 1, 2019, p. 149+.

[32] B. Akbugday, "Classification of Breast Cancer Data Using Machine Learning Algorithms," 2019 Medical Technologies Congress (TIPTEKNO), Izmir, Turkey, 2019, pp. 1-4.

[33] M. R. Al-Hadidi, A. Alarabeyyat and M. Alhannah, "Breast Cancer Detection Using K-Nearest Neighbor Machine Learning Algorithm," 2016 9th International Conference on Developments in eSystems Engineering (DeSE), Liverpool, 2016, pp. 35-39.

[34] Kibeom Jang, Minsoon Kim, Candace A Gilbert, Fiona Simpkins, Tan A Ince, Joyce M Slingerland "WEGFA activates an epigenetic pathway regulating ovarian cancer initiating cells" *Embo Molecular Medicines* Volume 9 Issue 3 (2017)

[35] Joseph A. Cruz and David S. Wishart "Applications of Machine Learning in cancer prediction and prognosis *Cancer informatics*" 2(3):59-77 · February 2007

[36] M. R. Al-Hadidi, A. Alarabeyyat and M. Alhanahnah, "Breast Cancer Detection Using K-Nearest Neighbor Machine Learning Algorithm," 2016 9th International Conference on Developments in eSystems Engineering (DeSE), Liverpool, 2016, pp. 35-39.

[37] Abdelghani Bellaachia, Erhan Guven, "Predicting Breast Cancer Survivability Using Data Mining Techniques"

[38] Ch. Shravya, K. Pravalika, Shaik Subhani, "Prediction of Breast Cancer Using Supervised Machine Learning Techniques", *International Journal of Innovative Technology and Exploring Engineering (IJITEE)*, Volume-8 Issue-6, April 2019.