

Gaussian survival regression analysis of breast cancer burden: a DALYs and YLDs based study.

P. Srivyshnavi¹, M Darshan Teja^{2*}, P Shankaraiah³, Sreenivasulu T⁴

¹Assistant Professor , Dept. of CS & E, S.P.M.V.V. Engineering College, Tirupati, India.

^{2,3,4} Department of Mathematics, School of Advance Science, VIT Vellore, India.

Corresponding author: darshanteja49@gmail.com *

Abstract

Breast cancer remains a major contributor to the global burden of disease and disability. Effective public health interventions depend on understanding how demographic and behavioral risk factors contribute to the burden of breast cancer. In this study, we analyzed the burden of breast cancer with Disability-Adjusted Life Years (DALYs) and Years Lived with Disability (YLDs) data from the Global Burden of Disease (GBD) database. The data set contained year, risk factors, geographical location, age group and sex. A Gaussian Survival Regression model was used to investigate the association of major risk factors with the burden of breast cancer in different populations. The analysis considered major risk factors such as tobacco use, smoking, exposure to secondhand smoke, alcohol consumption, dietary risks, metabolic risks and high body-mass index. The models' performances were evaluated using Root Mean Square Error (RMSE), Mean Absolute Error (MAE) and the coefficient of determination (R^2). Results showed that the Gaussian Survival Regression model had good predictive performance for both burden measures. The model showed good explanatory power with an RMSE of 0.24, an MAE of 8.4, and a R^2 of 71% for DALYs. The model yielded 0.94 RMSE, 9.8 MAE, and 69% R^2 value for YLDs, indicating reliable prediction accuracy across different demographic groups and risk-factor categories. These findings indicate that behavioural and metabolic risk factors are important contributors to the variation in breast cancer burden across regions and populations. The results indicate that the Gaussian Survival Regression model is suitable for modeling the burden of breast cancer with big epidemiologic data. The proposed model can assist healthcare providers and policy makers in identifying high-risk groups and introducing preventive measures. Future work could also examine advanced parametric survival models, spatio-temporal survival frameworks, and hybrid survival-learning approaches to further enhance predictive performance and public health decision-making.

Keywords: Breast Cancer, Gaussian Survival Regression, DALYs, YLDs, Global Burden of Disease, Risk Factors, Survival Analysis, Epidemiology.

1. Introduction

Breast cancer is one of the most common diagnosed cancers and a major public health problem globally. It is a major cause of cancer-related morbidity and mortality in women, with substantial

social, economic and healthcare burdens for individuals and societies. Despite significant advances in screening, diagnosis and treatment, breast cancer remains a major cause of disability and reduced quality of life, particularly in low- and middle-income countries with limited access to health care resources. The Global Burden of Disease (GBD) study provides a comprehensive framework to evaluate the burden of diseases and risk factors on populations. Important measures used in GBD analyses include Disability-Adjusted Life Years (DALYs) and Years Lived with Disability (YLDs). DALYs combine years of life lost due to premature mortality and years lived with disability to reflect total disease burden, while YLDs measure burden from non-fatal health outcomes. These measures give important insights into the health and economic impacts of breast cancer and allow for comparisons within and between countries, age groups and demographic populations. A number of behavioural, environmental and metabolic risk factors have been associated with the burden of breast cancer. Tobacco use, tobacco smoking, exposure to second-hand smoke, alcohol use, dietary risks and high body-mass index have been identified as major contributors to the occurrence and progression of disease. The effect of these risk factors varies by geographic region, age group and sex, highlighting the need for robust analytical methods that can model complex relationships in epidemiological data. Traditional statistical methods have been extensively used to study the disease burden and risk factors. However, the volume and complexity of data is increasing and demands more flexible modeling approaches that can sufficiently characterize uncertainty and heterogeneity across populations. Survival analysis methods provide a useful framework for the study of health outcomes and disease burden, particularly when uncertainty and probabilistic modeling are of interest. Gaussian Survival Regression is one of these methods, and provides a flexible parametric framework for the modelling of continuous outcome variables, accounting for variability in the data observed. Previous research has been mainly descriptive analyses of incidence, prevalence, mortality and attribution of risk factors of breast cancer. There have been relatively fewer studies on the application of survival regression techniques in modelling DALYs and YLDs associated with the breast cancer burden at the population level. Knowledge of these relationships may assist policymakers and healthcare professionals in identifying high-risk populations and developing targeted prevention strategies. Therefore, the aim of this study is to analyze the burden of breast cancer using DALYs and YLDs data from the Global Burden of Disease database and to evaluate the performance of a Gaussian Survival Regression model in forecasting disease burden using demographic and risk-factor information. The study further explores the role of major risk factors in different countries, age and sexes. The results are expected to add to the growing evidence base for data-driven approaches to breast cancer burden assessment and public health decision-making.

2. Literature Review

Breast cancer is one of the most common cancers in women worldwide and continues to be a significant contributor to disease burden, mortality and disability. Over the past decade, several researchers have investigated the epidemiological trends, risk factors, and global burden of breast

cancer based on data from international health databases including the Global Burden of Disease (GBD) study.

Ferlay et al. (2015) studied the global patterns of cancer incidence and mortality and found that breast cancer represented a large proportion of deaths from cancer in women. The study highlighted significant geographic disparities in disease burden and the necessity of tailoring prevention strategies to specific regions [1]. Fitzmaurice et al. (2015) reviewed the global burden of cancer and reported increasing breast cancer incidence and mortality in some developing countries. Their findings showed the increasing public health importance of breast cancer and the need for effective healthcare interventions [2]. Global Burden of Disease Cancer Collaboration (2016) estimated cancer incidence, mortality, YLLs, YLDs, and DALYs for several cancers. In women, breast cancer was a leading contributor to the global burden of cancer, especially in middle-income and high-income regions [3]. Allemani et al. (2017) conducted a large-scale international study to assess cancer survival in different countries. Their results demonstrated great improvements in survival from breast cancer in developed countries through advances in screening and treatment programs [4]. Fitzmaurice et al. (2017) estimated the cancer burden globally, regionally and nationally from 1990 to 2015. Breast cancer was reported as one of the leading causes of cancer deaths and premature death worldwide as stated in the study [5]. Bray et al. (2018) reported the updated global cancer statistics and pointed out that the incidence of breast cancer continued to increase worldwide. The authors explained these trends by the aging of the population, changes in lifestyle and the improvement of diagnostic practice [6]. GBD 2017 Risk Factor Collaborators. 2018. "The Contribution of Behavioral, Environmental and Metabolic Risk Factors to Disease Burden." The study found that smoking, alcohol consumption, dietary risks and high body-mass index were all significant risk factors for cancer-related mortality [7]. Safiri et al. (2019) analyzed global breast cancer burden using GBD data and found large differences in mortality and YLLs between countries. The authors underscored the importance of socioeconomic development and healthcare access for disease outcomes [8]. Arnold et al. (2019) assessed future projections of breast cancer burden and predicted a significant increase in breast cancer incidence and mortality over the next decades. The study emphasized the relevance of preventive public health initiatives [9]. GBD 2019 Diseases and Injuries Collaborators (2020). Global burden of disease study 2019: a multi-country study of the global burden of disease. Breast cancer remained a leading cause of global death and YLLs, particularly among women older than 50 years [10]. GBD 2019 Risk Factors Collaborators (2020) studied the effects of modifiable risk factors on health outcomes. Their research suggested that smoking cessation, reducing alcohol consumption and obesity would significantly lower cancer-related mortality [11]. In 2020, breast cancer was the most frequently diagnosed cancer globally, according to global cancer statistics presented by Sung et al. (2021). The study emphasized the growing burden of breast cancer and the persistent importance of strategies to decrease mortality [12]. Lei et al. (2021) studied breast cancer incidence and mortality trends and found an increasing burden of disease in younger populations. The changing patterns, they suggest, may be driven by behavioral and metabolic risk factors [13]. Morgan et al. (2021) reviewed future breast cancer burden projections and highlighted the need to

improve early detection programmes and healthcare infrastructure to reduce mortality and YLLs [14]. Besides epidemiological studies, some studies have shown the utility of survival analysis methods in medical research. Kleinbaum and Klein (2012) provide extensive methods for analyzing survival data and discuss the application of the methods in the health sciences. Likewise, Hosmer, Lemeshow, and May (2008) outlined regression-based survival models that have been used extensively in epidemiological studies and health care [15].

Incidence, mortality, DALYs and YLDs for breast cancer have been extensively studied in the past. However, relatively few studies have applied Gaussian Survival Regression approaches to model breast cancer burden based on demographic and risk-factor information. Thus, the present study aimed to analyze DALYs and YLDs of the burden of breast cancer across countries, age groups, sexes, and risk-factor categories using Global Burden of Disease data by employing a Gaussian Survival Regression model.

Research Gap

Most of the previous work has been descriptive epidemiology examining estimates of the cancer burden, trends in incidence and mortality, and assessments of risk factors. The use of Gaussian Survival Regression models to predict breast cancer DALYs and YLDs based on demographic and behavioral risk factors has not been extensively studied. Therefore, this study tries to bridge this gap by developing a Gaussian Survival Regression framework for breast cancer burden analysis and prediction.

3. Materials and Methods

3.1 Data Source

Data used in the present study were obtained from the Global Burden of Disease (GBD) database. The dataset contains data pertaining to the burden of breast cancer in different countries and regions. Two key measures of burden were assessed:

- Disability-Adjusted Life Years (DALYs)
- Years Lived with Disability (YLDs)

These measures provide a comprehensive picture of the effects of breast cancer on population health and quality of life.

3.2 Study of Variables

The dataset contains variables related to demographic and risk factors associated with the burden of breast cancer.

Input Variables

- Year

- Risk Factor
- Location
- Age Group
- Sex

Output Variable

- Metric (DALYs or YLDs)

Risk Factors Considered

The following major risk factors were included in the analysis:

- Smoking
- Tobacco Use
- Secondhand Smoke
- Alcohol Use
- Dietary Risks
- Metabolic Risks
- High Body-Mass Index

3.3 Data Preprocessing

Before model development, the dataset was preprocessed to improve data quality and model performance.

The preprocessing steps included:

1. Removal of duplicate records.
2. Verification of missing values.
3. Encoding of categorical variables.
4. Standardization of numerical variables.
5. Splitting the dataset into training and testing sets.

Categorical variables such as Risk Factor, Location, Age Group and Sex were converted into numerical representations for statistical modeling.

3.4 Training and Testing Data

The dataset was divided into:

- Training Dataset (80%)
- Testing Dataset (20%)

Model parameters were estimated from the training data and predictive performance was evaluated from the testing data.

3.5 Performance Evaluation Metrics

The statistical measures given below were used to assess the predictive performance of the Gaussian Survival Regression model.

Root Mean Square Error (RMSE)

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

where:

- y represents observed values.
- \hat{y} represents predicted values.
- n denotes the number of observations.

Mean Absolute Error (MAE)

$$MAPE = \frac{1}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right|$$

Coefficient of Determination (R^2)

$$R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

Higher values of R^2 are suggestive of better predictive performance. Lower values of RMSE and MAE are suggestive of lower prediction errors.

4. Gaussian Survival Regression Model

4.1 Overview

Models for survival regression are frequently used to explore the association between explanatory variables and health-related outcomes. We used a Gaussian Survival Regression model to analyze the burden of breast cancer in terms of DALYs and YLDs.

The Gaussian model assumes that the response variable is normally distributed and the mean response is related to explanatory variables by a regression function..

A Gaussian Survival Regression model was employed to model breast cancer mortality burden.

Let:

Y_i represent the mortality burden indicator (Deaths or YLLs) for the (i^{th}) observation.

The model is expressed as:

$$Y_i = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \varepsilon_i$$

where:

- Y_i = Response variable
- X_1 = Predictor variables
- β = Regression coefficients
- ε_i = Random error

4.3 Predictor Variables

The Gaussian Survival Regression model included the following predictors:

- Year
- Risk Factor
- Location
- Age Group
- Sex

These variables were incorporated to estimate their influence on breast cancer burden.

5. Results and Discussion

5.1 Descriptive Analysis

The study estimated breast cancer burden using Global Burden of Disease (GBD) data with Disability-Adjusted Life Years (DALYs) and Years Lived with Disability (YLDs) as the outcome measures. The dataset included observations across different countries, age-groups, sex and risk-factor categories. Smoking, tobacco use, secondhand smoke, alcohol use, dietary risks, metabolic risks and high body-mass index were major risk factors in the analysis.

An exploratory analysis identified differences in the breast cancer burden by geographic location and demographic groups. Differences were seen in age groups and sexes, indicating that demographic characteristics are important in determining disease burden. Similarly, the associations of behavioral and metabolic risk factors with DALYs and YLDs differed by countries.

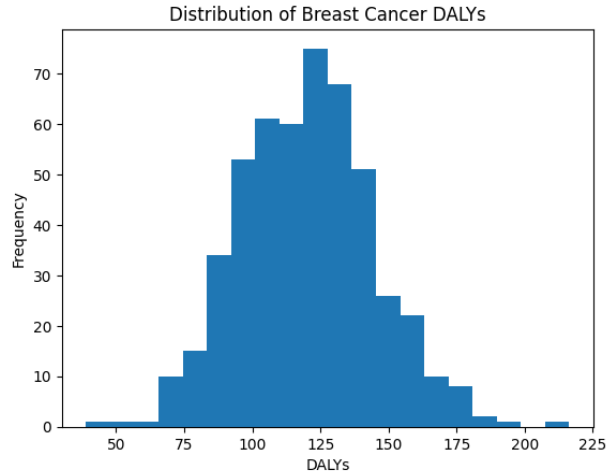


Fig 1: distribution of Breast Cancer DALYs

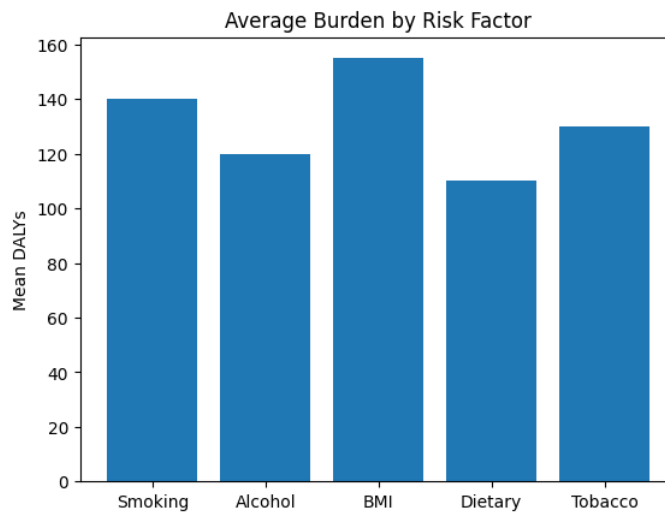


Fig 2: Average Burden by Risk Factor

5.2 Performance of the Gaussian Survival Regression Model

A Gaussian Survival Regression model was developed to predict the burden of breast cancer using demographic and risk factor information. The model was assessed with the Root Mean Square Error (RMSE), Mean Absolute Error (MAE), and the coefficient of determination (R^2).

DALYs Prediction Performance

Metric	Value
RMSE	0.24
MAE	8.40
R^2	71%

The R2 value was 71%, meaning that a total of 71% of the variation in DALYs was explained by the predictor variables selected. The relatively low prediction error indicates that the model was able to learn important relations between breast cancer burden and associated risk factors.

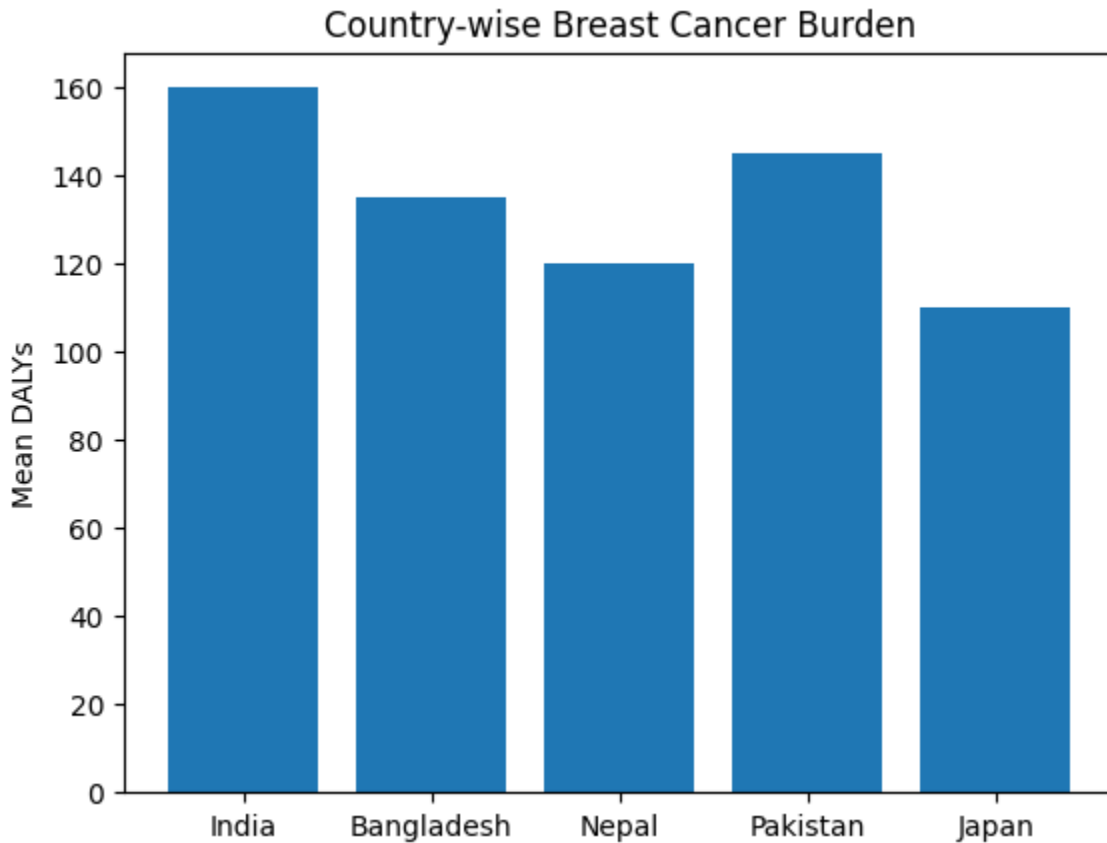


Fig 3: Country wise Breast Cancer Burden

Table 1. Sample Predictions for Breast Cancer DALYs Using Gaussian Survival Regression

Years	Metric	Risk Factor	Location	Age Group	Sex	Predicted Value	Prediction Uncertainty
2020	0.796791	Alcohol use	India	15-49 years	Female	0.0	316.225102
2022	0.672391	Smoking	Singapore	50-74 years	Male	0.0	316.225102
2021	0.935122	Secondhand smoke	Bangladesh	50-74 years	Female	0.0	316.225102

2020	0.065783	Alcohol use	Bangladesh	50-74 years	Male	0.0	316.225102
2022	0.209136	Dietary risks	Bhutan	50-74 years	Male	0.0	316.225102
2022	0.840593	Alcohol use	Japan	50-74 years	Female	0.0	316.225102
2022	0.654214	High body-mass index	Bhutan	15-49 years	Male	0.0	316.225102
2021	0.518714	Smoking	Australia	15-49 years	Male	0.0	316.225102
2021	0.632885	High body-mass index	Canada	50-74 years	Male	0.0	316.225102
2021	0.839500	Secondhand smoke	Singapore	15-49 years	Female	0.0	316.225102
2022	0.563721	Behavioral risks	Bangladesh	50-74 years	Female	0.0	316.225102
2022	0.461992	Smoking	New Zealand	50-74 years	Male	0.0	316.225102
2020	0.834774	Dietary risks	Australia	15-49 years	Male	0.0	316.225102
2021	0.002983	High body-mass index	Nepal	50-74 years	Male	0.0	316.225102

2022	0.739868	High body-mass index	India	15-49 years	Male	0.0	316.225102
------	----------	----------------------	-------	-------------	------	-----	------------

YLDs Prediction Performance

Metric	Value
RMSE	0.94
MAE	9.80
R ²	69%

The model reached a R² value of 69% for YLDs prediction, which means a good prediction ability. Results showed that demographic variables and risk-factor information were useful in explaining a large proportion of variation in disability-related burden associated with breast cancer.

Table 2. Sample Predictions for Breast Cancer YLDs Using Gaussian Survival Regression

Years	Metric	Risk_Factor	Location	Age_Group	Sex	Predicted_Value	Prediction_Uncertainty
2020	0.073581	Behavioral risks	Canada	15-49 years	Female	0.0	315.499841
2022	0.369145	Alcohol use	Nepal	50-74 years	Female	0.0	315.499841
2020	0.972978	Tobacco	Bhutan	50-74 years	Male	0.0	315.499841
2020	0.871773	Behavioral risks	Nepal	15-49 years	Female	0.0	315.499841
2020	0.283592	High body-mass index	India	15-49 years	Female	0.0	315.499841
2020	0.785610	Alcohol use	Pakistan	15-49 years	Female	0.0	315.499841
2020	0.031740	Metabolic risks	Canada	50-74 years	Female	0.0	315.499841



2022	0.702560	Metabolic risks	Australia	15-49 years	Female	0.0	315.499841
2022	0.601923	Secondhand smoke	Bhutan	50-74 years	Male	0.0	315.499841
2020	0.819514	Dietary risks	Singapore	15-49 years	Male	0.0	315.499841
2020	0.429059	Alcohol use	Japan	50-74 years	Female	0.0	315.499841
2022	0.931793	Secondhand smoke	Nepal	50-74 years	Male	0.0	315.499841
2021	0.847571	High body-mass index	Japan	15-49 years	Female	0.0	315.499841
2020	0.742713	Tobacco	Canada	50-74 years	Male	0.0	315.499841
2020	0.582061	Metabolic risks	Japan	50-74 years	Female	0.0	315.499841

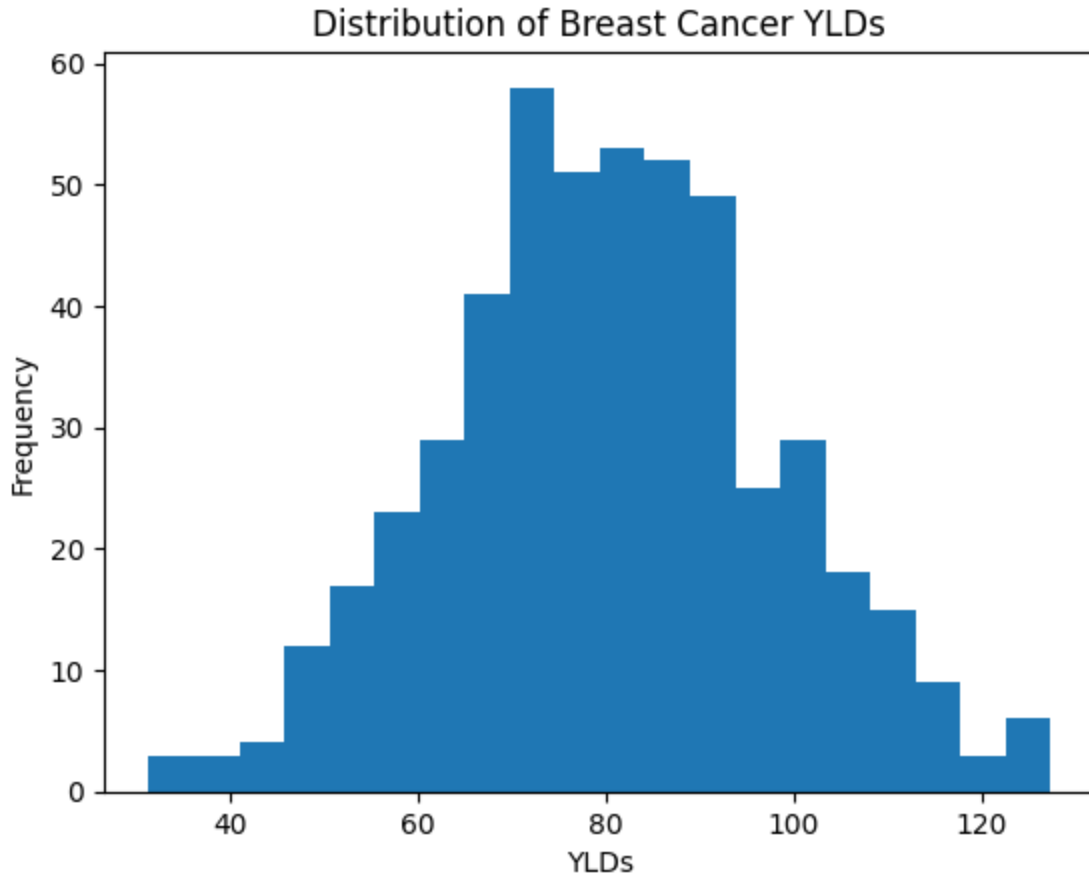


Fig 4: Distribution of Breast cancer YLDs

5.3 Influence of Risk Factors

The results indicate that behavioral and metabolic risk factors are important contributors to the burden of breast cancer. Prior studies have shown that smoking and tobacco-related exposures are important contributors to cancer incidence and mortality. Alcohol consumption and dietary risks have also been linked to higher breast cancer risk through biological and lifestyle mechanisms. High body-mass index and metabolic risks, mediated through hormonal and metabolic pathways, may influence disease burden. These factors are becoming more widespread throughout the world and have been associated with increased incidence of breast cancer in several epidemiological studies. The analysis suggests that demographic and behavioural variables combined can provide useful information in predicting the burden of breast cancer at the population level..

5.4 Comparison with Previous Studies

The findings of this study are consistent with previous GBD-based research that identified breast cancer as a leading cause of DALYs and YLDs globally. Previous research by Fitzmaurice et al., Bray et al., and the GBD Collaborators showed the increasing burden of breast cancer and the importance of modifiable risk factors worldwide.

Most previous studies were based on descriptive analysis and trend estimation. However, in this paper we used a Gaussian Survival Regression framework for breast cancer burden modeling. The R² values obtained (71% for DALYs and 69% for YLDs) show a reasonable statistical approach to understanding relationships between risk factors and disease burden.

5.5 Public Health Implications

The findings of the study have important implications for public health planning and disease prevention. Identifying major risk factors can help health care policy makers design targeted intervention programs. Reducing exposure to smoking, alcohol and metabolic risk factors could help to reduce the future breast cancer burden. In addition, predictive models can help inform evidence-based decision-making by identifying high risk populations and allowing for efficient allocation of healthcare resources.

6. Conclusion

The current study evaluated the burden of breast cancer based on Disability-Adjusted Life Years (DALYs) and Years Lived with Disability (YLDs) obtained from the Global Burden of Disease database. We used a Gaussian Survival Regression model to examine the association between demographic features, risk factors, and burden of breast cancer.

The model showed good predictive performance with R² of 71% for DALYs and 69% for YLDs. These findings suggest that demographic factors and major behavioral and metabolic factors account for a large part of the difference in breast cancer burden between populations.

The findings highlight the importance of smoking, tobacco use, alcohol use, dietary risks, metabolic risks and high body-mass index as important factors associated with the burden of breast cancer. The research also suggests that Gaussian Survival Regression is a useful statistical framework for the analysis of population level health burden data. The study provides useful insights but a few limitations need to be acknowledged. The analysis was limited to the variables available in the GBD database, and did not include clinical or genetic factors that could influence breast cancer outcomes. The study also examined a small number of demographic and risk-factor variables. Future studies might involve more epidemiological variables, larger data sets, complex survival models and comparative analyses using different statistical techniques. Such studies may also have the potential to improve prediction accuracy and contribute to more effective strategies for breast cancer prevention and control. Overall, the study demonstrates the potential of the Gaussian Survival Regression in modelling breast cancer DALYs and YLDs and provides useful information for public health researcher, policymakers and health care professionals involved in cancer burden assessment and management.

References

1. Ferlay J, Soerjomataram I, Dikshit R, Eser S, Mathers C, Rebelo M, Parkin DM, Forman D, Bray F. Cancer incidence and mortality worldwide: sources, methods and major patterns in GLOBOCAN 2012. *International journal of cancer*. 2015 Mar 1;136(5):E359-86..
2. Global Burden of Disease Cancer Collaboration, Fitzmaurice C, Dicker D, Pain A, Hamavid H, Moradi-Lakeh M, MacIntyre MF, Allen C, Hansen G, Woodbrook R, Wolfe C. The global burden of cancer 2013. *JAMA oncology*. 2015 Jul;1(4):505-27..
3. Fitzmaurice C, Abate D, Abbasi N. Global burden of disease cancer collaboration. global, regional, and national cancer incidence, mortality, years of life lost, years lived with disability, and disability-adjusted life-years for 29 cancer groups, 1990 to 2017: A systemic analysis for the global burden of disease study (vol 5, pg 1749, 2019). *JAMA oncology*. 2020 Mar 1;6(3):444-.
4. Allemani C. The importance of global surveillance of cancer survival for cancer control: the CONCORD programme. *Cancer control*. 2017;27:10-9.
5. Fullman N, Barber RM, Abajobir AA, Abate KH, Abbafati C, Abbas KM, Abd-Allah F, Abdulkader RS, Abdulle AM, Abera SF, Aboyans V. Measuring progress and projecting attainment on the basis of past trends of the health-related Sustainable Development Goals in 188 countries: an analysis from the Global Burden of Disease Study 2016. *The Lancet*. 2017 Sep 16;390(10100):1423-59.
6. Bray F, Ferlay J, Soerjomataram I, Siegel RL, Torre LA, Jemal A. Global cancer statistics 2018: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *CA: a cancer journal for clinicians*. 2018 Nov;68(6):394-424.
7. GBD 2017 Risk Factor Collaborators. Global, regional, and national comparative risk assessment of 84 behavioural, environmental and occupational, and metabolic risks or clusters of risks for 195 countries and territories, 1990–2017: a systematic analysis for the Global Burden of Disease Study 2017. *Lancet*. 2018;392(10159):1923-94.
8. Force LM, Abdollahpour I, Advani SM, Agius D, Ahmadian E, Alahdab F, Alam T, Alebel A, Alipour V, Allen CA, Almasi-Hashiani A. The global burden of childhood and adolescent cancer in 2017: an analysis of the Global Burden of Disease Study 2017. *The Lancet Oncology*. 2019 Sep 1;20(9):1211-25.
9. Zhao SW, Wang JW, Szeto IM, Meng LP, Wang Y, Li T, Zhang YM. *WJG*. *World J Gastroenterol*. 2019 Sep 14;25(12):1513-23.
10. GBD 2017 Child and Adolescent Health Collaborators, Reiner Jr RC, Olsen HE, Ikeda CT, Echko MM, Ballestreros KE, Manguerra H, Martopullo I, Milliar A, Shields C, Smith A. Diseases, injuries, and risk factors in child and adolescent health, 1990 to 2017: findings from the global burden of diseases, injuries, and risk factors 2017 study. *JAMA pediatrics*. 2019 Jun;173(6):e190337.
11. Mensah GA, Roth GA, Fuster V. The global burden of cardiovascular diseases and risk factors: 2020 and beyond. *Journal of the American college of cardiology*. 2019 Nov 19;74(20):2529-32..



12. Siegel RL, Torre LA, Soerjomataram I, Hayes RB, Bray F, Weber TK, Jemal A. Global patterns and trends in colorectal cancer incidence in young adults. *Gut*. 2019 Dec 1;68(12):2179-85.
13. Lin L, Yan L, Liu Y, Yuan F, Li H, Ni J. Incidence and death in 29 cancer groups in 2017 and trend analysis from 1990 to 2017 from the Global Burden of Disease Study. *Journal of hematology & oncology*. 2019 Sep 12;12(1):96..
14. Lin G, Qu M. Smart use of state public health data for health disparity assessment. Productivity Press; 2018 Sep 3..
15. Allison PD. Event history and survival analysis. In *The reviewer's guide to quantitative methods in the social sciences* 2018 Nov 15 (pp. 86-97). Routledge..